# Lecture 6:
# "Backpropagation Algorithm"

**Dr. Ehab Essa**

**Computer Science Department, Mansoura University**

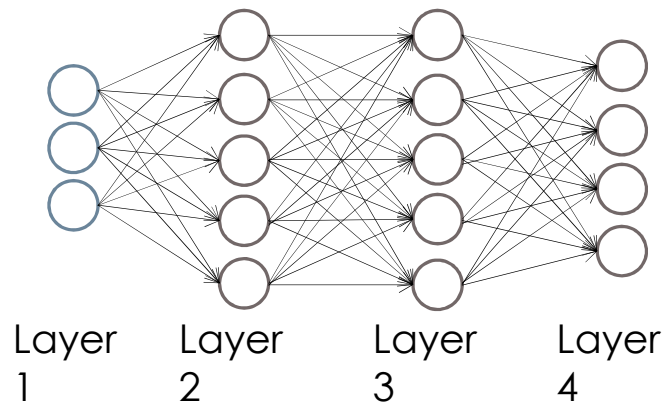# Today

▶ Backpropagation Algorithm

▶ Lecture Notes:

  ▶ **Chapter 6**, Fundamentals of Neural Networks: Architectures, Algorithms, and Applications by Laurene Fausett

Layer 1   Layer 2   Layer 3   Layer 4

$$\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \ldots, (x^{(m)}, y^{(m)})\}$$

$L = $ total no. of layers in network

$s_l = $ no. of units (not counting bias unit) in layer $l$

## Binary classification

$y = 0$ or $1$

1 output unit

## Multi-class classification (K classes)

$y \in \mathbb{R}^K$ E.g. $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$, $\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$, $\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$ $\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$

pedestrian  car  motorcycle  truck

K output units

- Assume that a set of examples $\Im=\{\mathbf{x}(n),\mathbf{t}(n)\}$, n=1,...,N is given. $\mathbf{x}(n)$ is the *input vector* of dimension $m_0$ and t(n) is the *desired response* vector of dimension M

- Thus an *error signal*, $e_{j(n)} = y_{j(n)} - t_j(n)$ can be defined for the output neuron j.

- We can derive a learning algorithm for an MLP by assuming an optimization approach which is based on the **steepest descent direction**, I.e.

$$\Delta\mathbf{w}(n) = -\alpha\mathbf{g}(n)$$

Where $\mathbf{g}$(n) is the gradient vector of the cost function and $\alpha$ is the *learning rate*.

- The algorithm that it is derived from the steepest descent direction is called **back-propagation**
- Assume that we define a SSE instantaneous cost function (I.e. per example) as follows:

$$E(n) = \frac{1}{2} \sum_{j \in C} e_j^{2}(n)$$

▶ Where C is the set of all *output neurons*.
- If we assume that there are N examples in the set ℑ then the *average squared error* is:

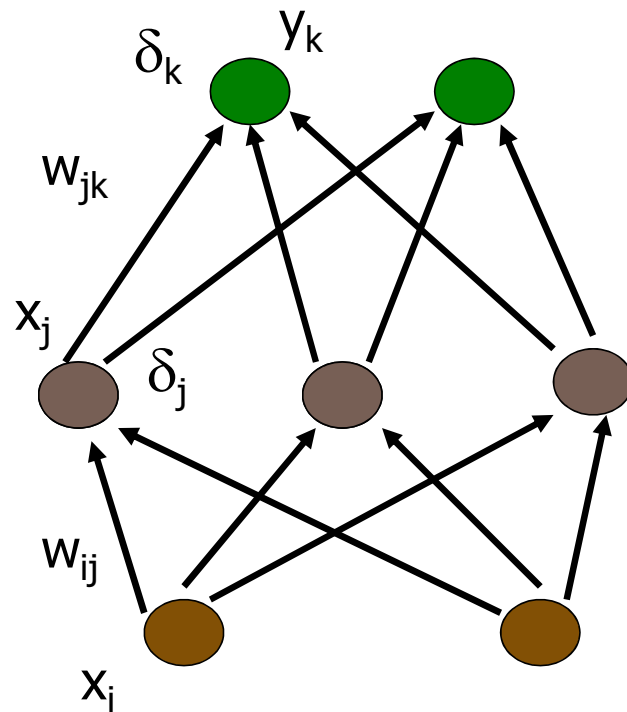$$E_{av} = \frac{1}{N} \sum_{n=1}^{N} E(n)$$

- We have two ways to calculate the gradient with respect to $E_{av}$ or $E(n)$.
  - **Batch** mode : In the first case we calculate the gradient per *epoch* (i.e. in all patterns N). $E_{av}$ is used. (average gradient over the examples in the batch)
  - **Online** or **Stochastic** mode: updates model parameters according to the gradient calculated from one *pattern*. $E(n)$ is used

- Assume that we use the online mode for the rest of the calculation. The gradient is defined as:

$$\nabla g(n) = \frac{\partial E(n)}{\partial w_{ji}(n)}$$

$\delta_k$ $\quad$ $y_k$

$w_{jk}$

$x_j$ $\quad$ $\delta_j$

$w_{ij}$

$x_i$

**Backward step**: propagate errors from output to hidden layer

**Forward step**: Propagate activation from input to output layer

$$\frac{d}{dx}g(x) = \frac{d}{dx}\left(\frac{1}{1+e^{-x}}\right)$$

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

$$= -(1+e^{-x})^{-2}(-e^{-x})$$

$$= \frac{e^{-x}}{(1+e^{-x})^2}$$

$$= \frac{(1+e^{-x})-1}{(1+e^{-x})^2}$$

$$= \frac{1+e^{-x}}{(1+e^{-x})^2} - \left(\frac{1}{1+e^{-x}}\right)^2$$

$$= g(x) - g(x)^2$$

$$g(x)' = g(x)(1-g(x))$$

# Chain Rule

▶ In calculus, the chain rule is a formula for computing the derivative of the composition of two or more functions

▶ Consider z to be a function of the variable y, which is itself a function of x (y and z are therefore dependent variables), and so, z becomes a function of x as well, The chain rule may be written:

$$\frac{dz}{dx} = \frac{dz}{dy} \cdot \frac{dy}{dx}.$$
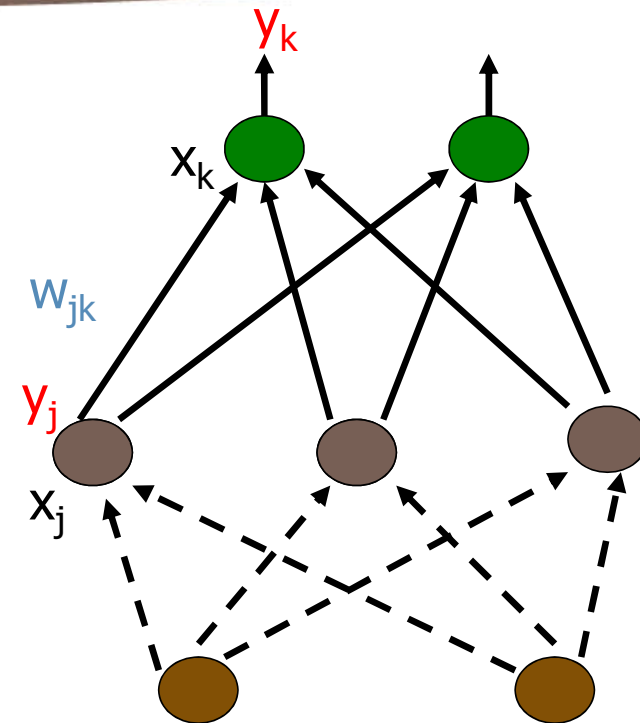
▶ If $z(y) = e^y$ and y = 2x, so that $z(y(x)) = e^{2x}$.

$$\frac{dz(y)}{dx} = \frac{dz(y)}{dy} \times \frac{dy(x)}{dx} = e^y \cdot 2$$

## Output layer node

$$\frac{\partial E}{\partial W_{jk}} = \frac{\partial E}{\partial y_k} \cdot \frac{\partial y_k}{\partial x_k} \cdot \frac{\partial x_k}{\partial W_{jk}}$$

$$\frac{\partial E}{\partial W_{jk}} = \frac{\partial}{\partial W_{jk}} \frac{1}{2} \sum_{k \in K} (y_k - t_k)^2$$

$$\frac{\partial E}{\partial W_{jk}} = (y_k - t_k) \frac{\partial}{\partial W_{jk}} y_k$$

$$\frac{\partial E}{\partial W_{jk}} = (y_k - t_k) \frac{\partial}{\partial W_{jk}} g(x_k)$$

$$\frac{\partial E}{\partial W_{jk}} = (y_k - t_k) g(x_k)(1 - g(x_k)) \frac{\partial}{\partial W_{jk}} x_k$$

$$\frac{\partial E}{\partial W_{jk}} = (y_k - t_k) y_k (1 - y_k) y_j$$

$y_k$

$x_k$

$W_{jk}$

$y_j$

$x_j$

**x** pre-activation
**y** after-activation

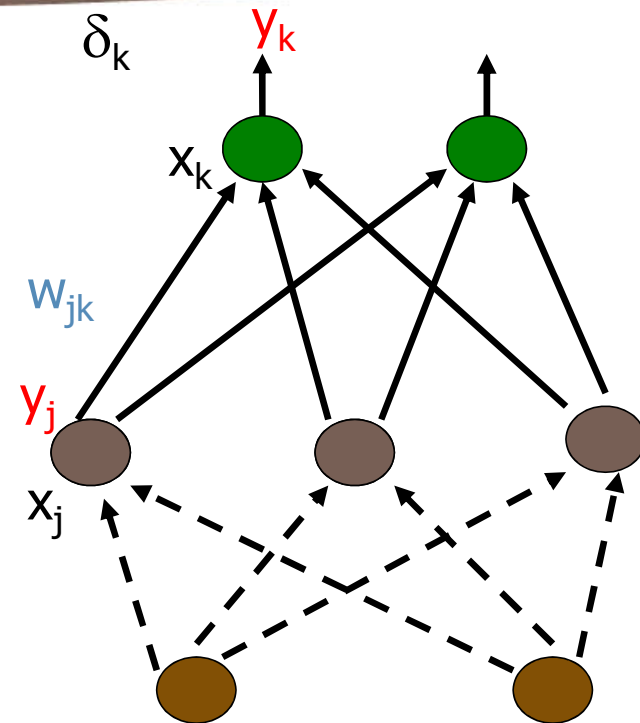▶ For notation purposes I will define $\delta_k$ to be the expression $(y_k - t_k) y_k (1 - y_k)$, so we can rewrite the equation above as

▶ $\frac{\partial E}{\partial W_{jk}} = (y_k - t_k)y_k(1 - y_k)y_j$
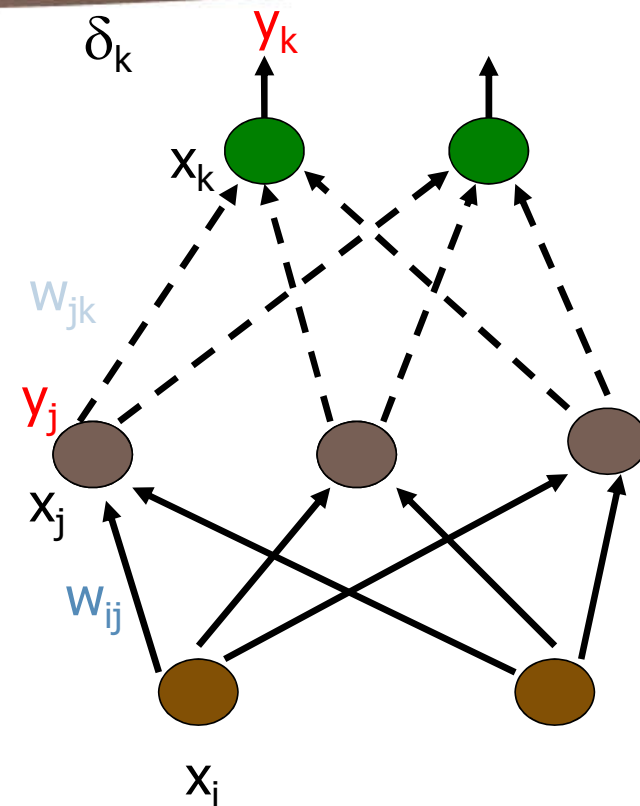
▶ $\frac{\partial E}{\partial W_{jk}} = \delta_k \; y_j$

Where $\delta_k = (y_k - t_k)y_k(1 - y_k)$,

$$\frac{\partial E}{\partial W_{ij}} = \frac{\partial E}{\partial y_k} \cdot \frac{\partial y_k}{\partial x_k} \cdot \frac{\partial x_k}{\partial y_j} \cdot \frac{\partial y_j}{\partial x_j} \cdot \frac{\partial x_j}{\partial W_{ij}}$$

$$\frac{\partial E}{\partial W_{ij}} = \frac{\partial}{\partial W_{ij}} \frac{1}{2} \sum_{k \in K} (y_k - t_k)^2$$

$$\frac{\partial E}{\partial W_{ij}} = \sum_{k \in K} (y_k - t_k) \frac{\partial}{\partial W_{ij}} y_k$$

$$\frac{\partial E}{\partial W_{ij}} = \sum_{k \in K} (y_k - t_k) \frac{\partial}{\partial W_{ij}} g(x_k)$$

$$\frac{\partial E}{\partial W_{ij}} = \sum_{k \in K} (y_k - t_k) \, g(x_k)(1 - g(x_k)) \frac{\partial}{\partial W_{ij}} x_k$$

$$\frac{\partial E}{\partial W_{ij}} = \sum_{k \in K} (y_k - t_k) \, y_k(1 - y_k) \frac{\partial x_k}{\partial y_j} \frac{\partial y_j}{\partial W_{ij}}$$

$$\frac{\partial E}{\partial W_{ij}} = \sum_{k \in K} (y_k - t_k) \, y_k(1 - y_k) W_{jk} \frac{\partial y_j}{\partial W_{ij}}$$
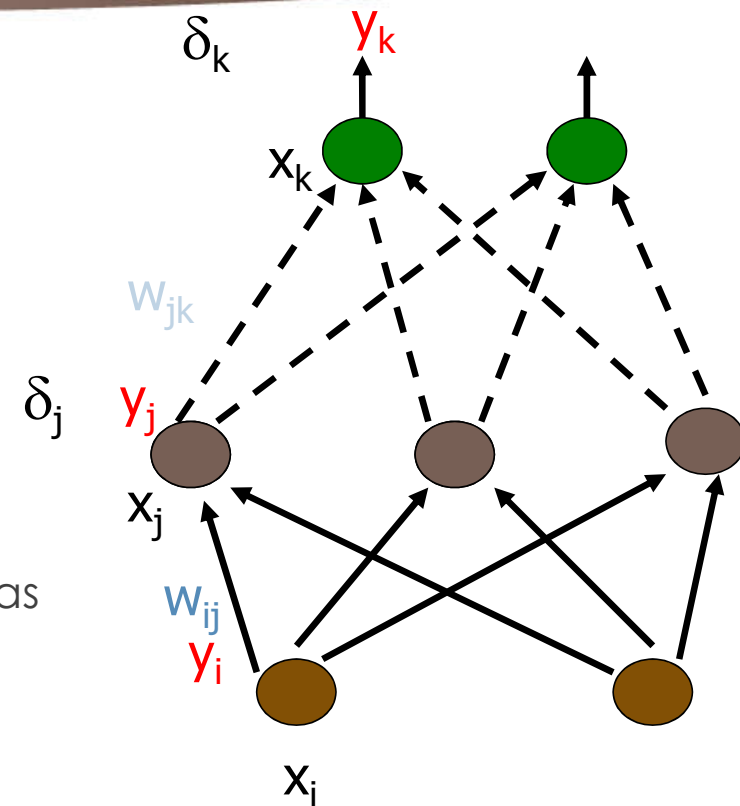
$\delta_k$

$y_k$

$x_k$

$W_{jk}$

$y_j$

$x_j$

$W_{ij}$

$x_i$

- $\frac{\partial E}{\partial W_{ij}} = \sum_{k \in K}(y_k - t_k)\, y_k(1 - y_k)W_{jk}\frac{\partial y_j}{\partial W_{ij}}$

- $\frac{\partial E}{\partial W_{ij}} = \frac{\partial y_j}{\partial W_{ij}}\sum_{k \in K}(y_k - t_k)\, y_k(1 - y_k)W_{jk}$

- $\frac{\partial E}{\partial W_{ij}} = y_j(1 - y_j)\frac{\partial x_j}{\partial W_{ij}}\sum_{k \in K}(y_k - t_k)\, y_k(1 - y_k)W_{jk}$

- $\frac{\partial E}{\partial W_{ij}} = y_j(1 - y_j)y_i\sum_{k \in K}(y_k - t_k)\, y_k(1 - y_k)W_{jk}$

- But, recalling our definition of $\delta_k$ we can write this as

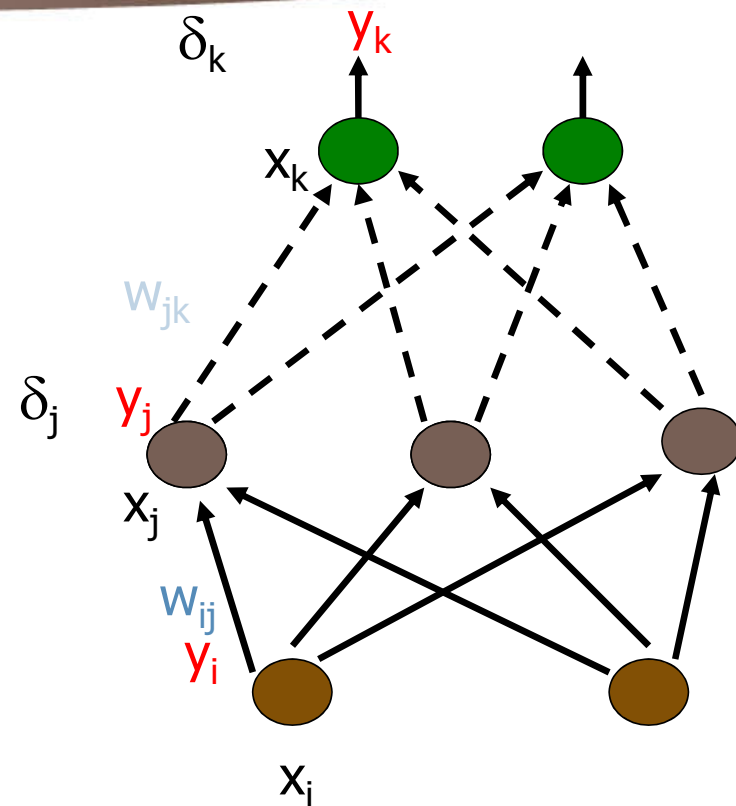$$\frac{\partial E}{\partial W_{ij}} = y_i y_j(1 - y_j)\sum_{k \in K}\delta_k\, W_{jk}$$

▶ Similar to before we will now define all terms besides the $y_i$ to be $\delta_j$, so we have

$$\frac{\partial E}{\partial W_{ij}} = y_i \delta_j$$

Where $\delta_j = y_j(1 - y_j) \sum_{k \in K} \delta_k W_{jk}$

$\delta_k$     $y_k$

$x_k$

$w_{jk}$

$\delta_j$     $y_j$

$x_j$

$w_{ij}$

$y_i$

$x_i$

▶ Initialize each $w_i$ to some small random value

▶ Until the termination condition is met, Do

    ▶ For each training example $<(x_1,\ldots x_n),t>$ Do

**Forward phase**

        ▶ Input the instance $(x_1,\ldots,x_n)$ to the network and compute the network outputs $y_k$

        ▶ For each output unit k

$$\delta_k = y_k(1\text{-}y_k)(y_k - t_k)$$

**Backward phase**

        ▶ For each hidden unit j

In case, weight are from input layer to hidden layer: $y_i = x_i$

$$\delta_j = y_j(1\text{-}y_j)\sum_k w_{j,k}\,\delta_k$$

        ▶ For each network weight $w_{i,j}$ Do

$$w_{i,j} = w_{i,j} + \Delta w_{i,j} \quad \text{where } \Delta wij = -\alpha\delta_j y_i$$

$$(\text{Bias}) \quad w_{0,j} = w_{0,j} + \Delta w_{0,j} \quad \text{where } \Delta w_{0j} = -\alpha\delta_j$$
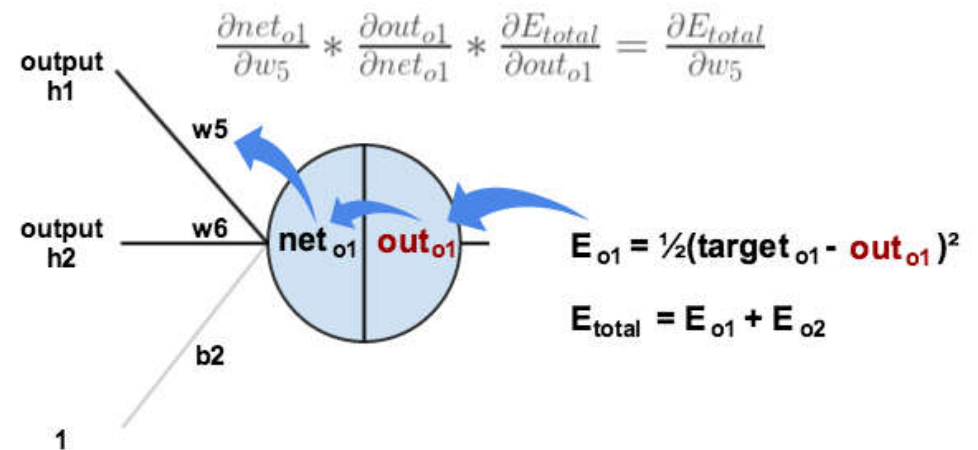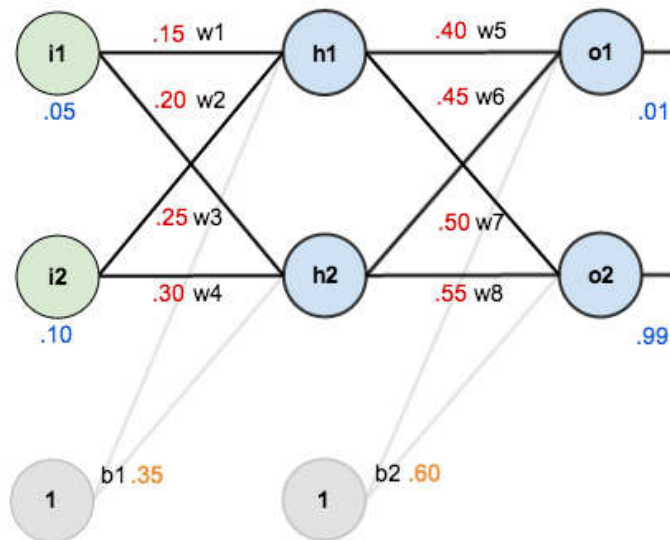
**Learning rules**

# Numerical Example

▶ Check the example

   ▶ https://mattmazur.com/2015/03/17/a-step-by-step-backpropagation-example/
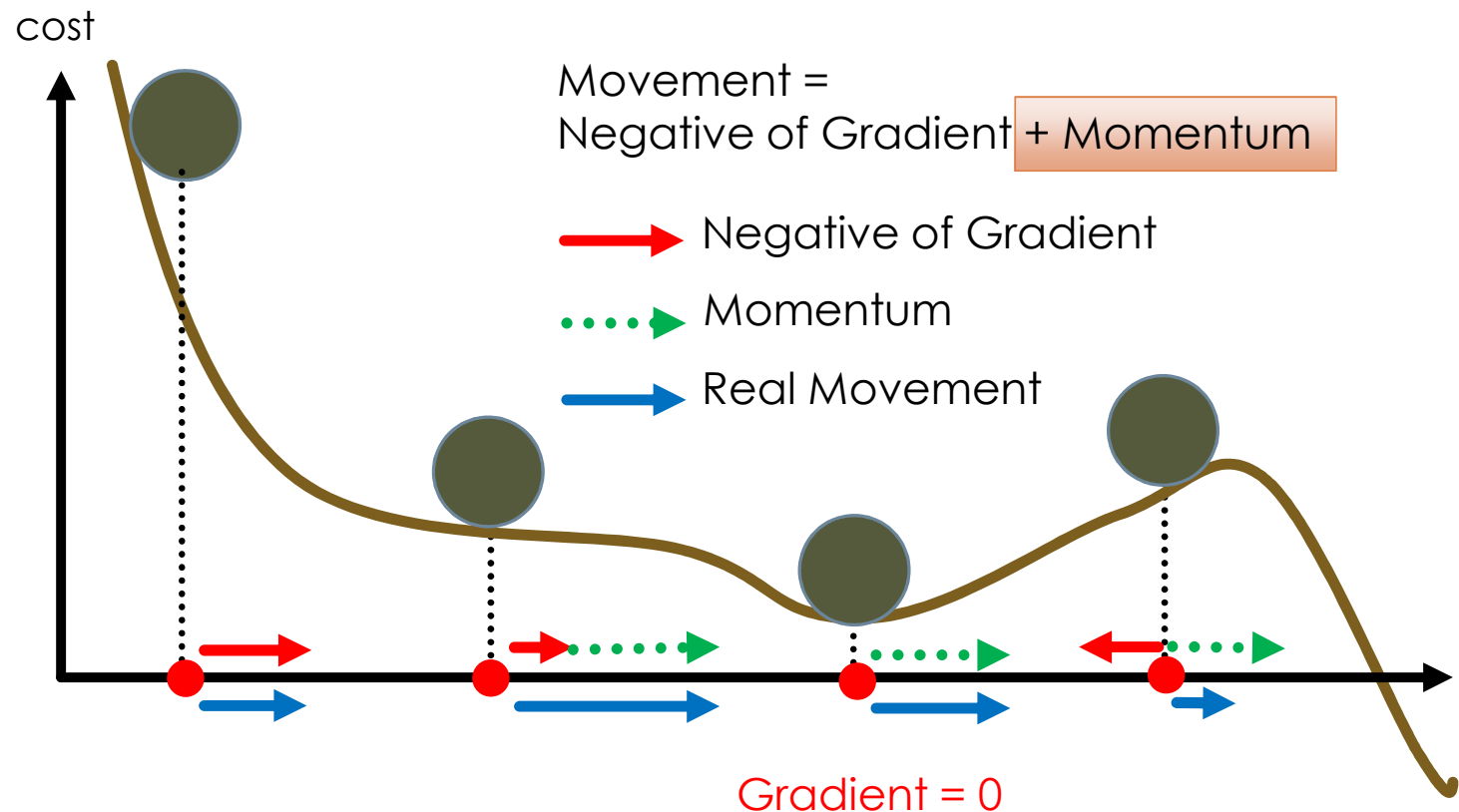
Adding **momentum term**

- Momentum is a method that helps accelerate SGD in the relevant direction and dampens oscillations.

- It does this by adding a fraction of the update vector of the past time step to the current update vector.

cost

Movement =
Negative of Gradient + Momentum

→ Negative of Gradient

····▸ Momentum

→ Real Movement

Gradient = 0

Adding **momentum term**

▶ Weights update at time t+1 contains the momentum of the previous updates, e.g.,
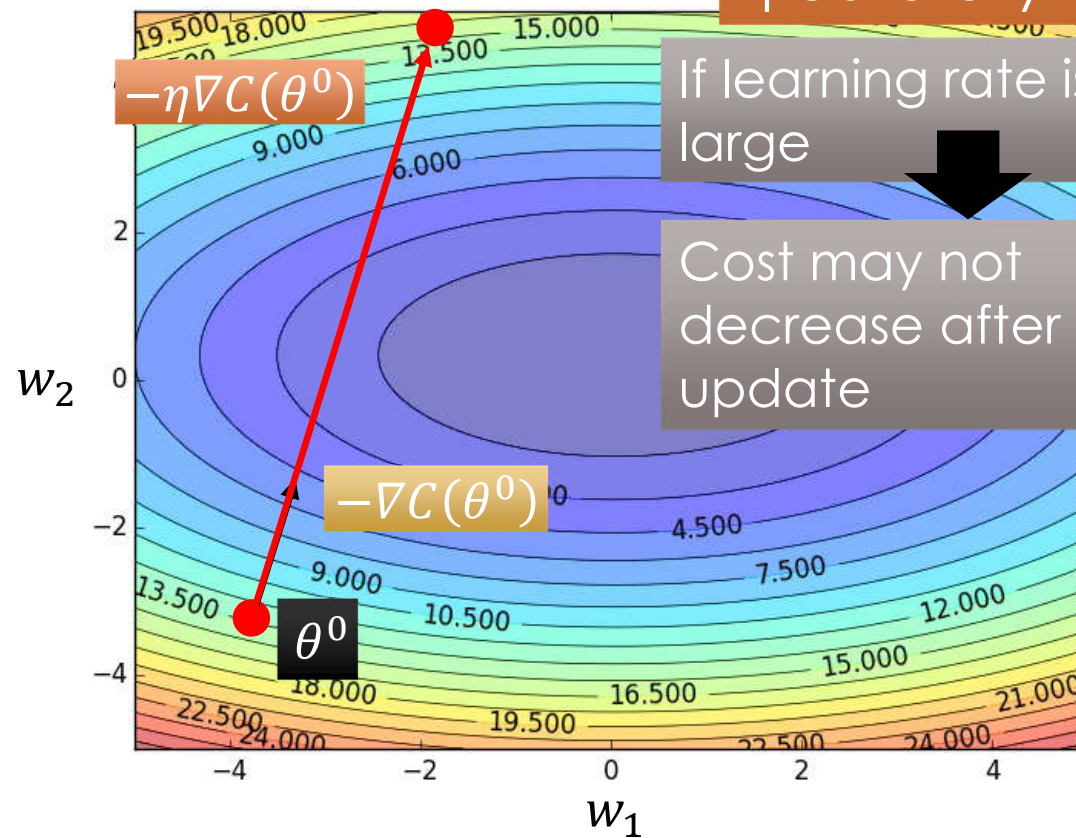
▶
$$\Delta w_{ij}(t+1) = \alpha \cdot \delta_j \cdot y_i + \mu \cdot \Delta w_{ij}(t), \text{ where } 0 < \mu < \alpha << 1$$

$$\text{then } \Delta w_{ij}(t+1) = \sum_{s=1}^{t} \mu^{t-s} \alpha \cdot \delta_j(s) \cdot y_i(s)$$

▶ an exponentially weighted sum of all previous updates

▶ Avoid sudden change of directions of weight update (smoothing the learning process)

▶ Error is no longer monotonically decreasing

# Learning Rate



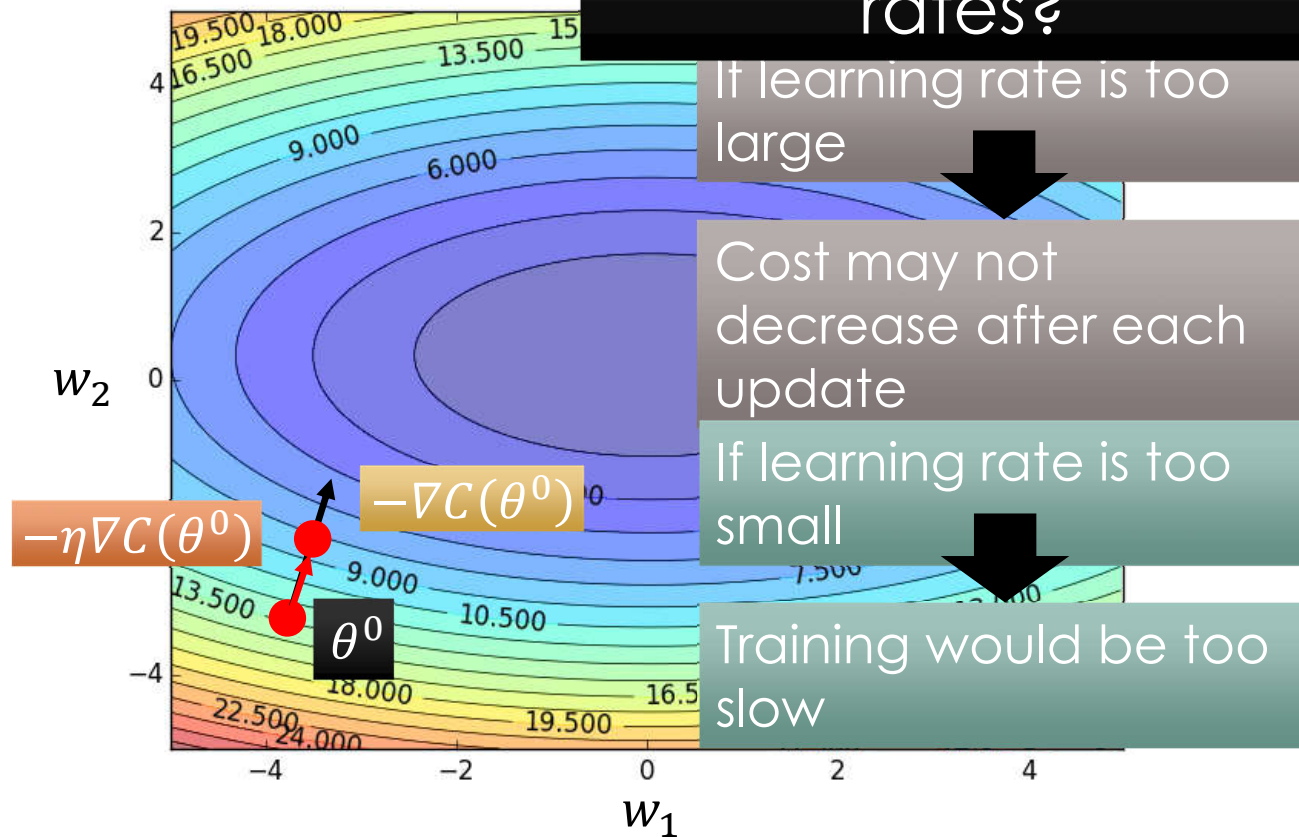Set the learning rate η carefully

If learning rate is too large

Cost may not decrease after each update

$-\eta \nabla C(\theta^0)$

$-\nabla C(\theta^0)$

$\theta^0$

$w_2$

$w_1$

# Learning Rate

Can we give different parameters different learning rates?

If learning rate is too large

Cost may not decrease after each update

If learning rate is too small

Training would be too slow

$-\eta \nabla C(\theta^0)$

$-\nabla C(\theta^0)$

$\theta^0$

$w_2$

$w_1$

# Adagrad

Original Gradient Descent

$$w^t \leftarrow w^{t-1} - \eta \nabla E(w^{t-1})$$

Each parameter w are considered separately

$$w^{t+1} \leftarrow w^t - \eta_w g^t \qquad g^t = \frac{\partial E(w^t)}{\partial w}$$

Parameter dependent learning rate

constant

$$\eta_w = \frac{\eta}{\sqrt{\sum_{i=0}^{t}(g^i)^2}}$$

Summation of the square of the previous derivatives

# Adagrad

$$\eta_w = \frac{\eta}{\sqrt{\sum_{i=0}^{t}(g^i)^2}}$$

$w_1$

| $g^0$ |
|-------|
| 0.1   |

$w_2$

| $g^0$ |
|-------|
| 20.0  |

Learning rate:

$$\frac{\eta}{\sqrt{0.1^2}} = \frac{\eta}{0.1}$$

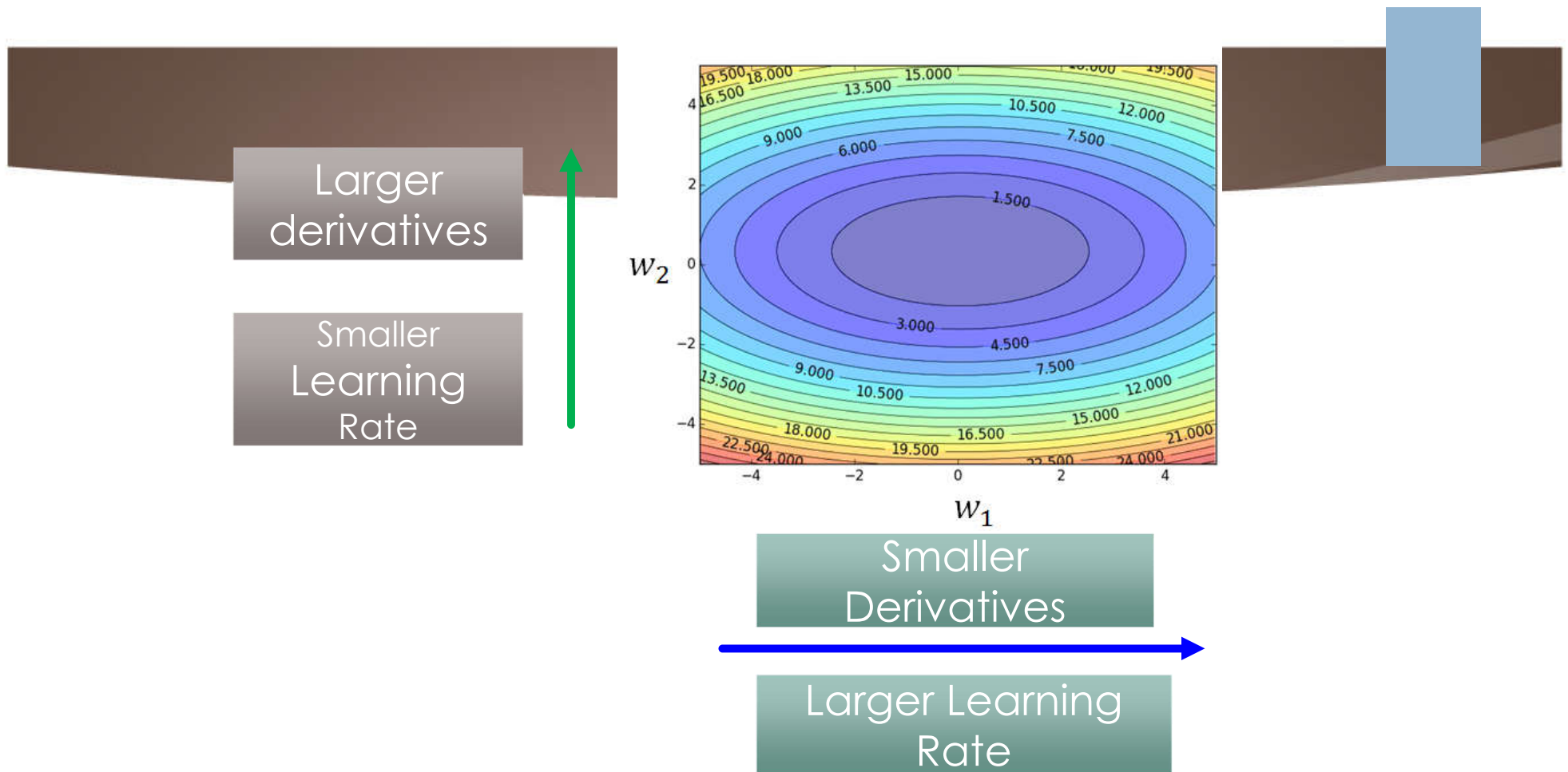$$\frac{\eta}{\sqrt{0.1^2 + 0.2^2}} = \frac{\eta}{0.22}$$

Learning rate:

$$\frac{\eta}{\sqrt{20^2}} = \frac{\eta}{20}$$

$$\frac{\eta}{\sqrt{20^2 + 10^2}} = \frac{\eta}{22}$$

**_Observation:_** 1. Learning rate is smaller and smaller for all parameters

2. Smaller derivatives, larger learning rate, and vice versa

Why?

Larger derivatives

Smaller Learning Rate

$w_2$

$w_1$

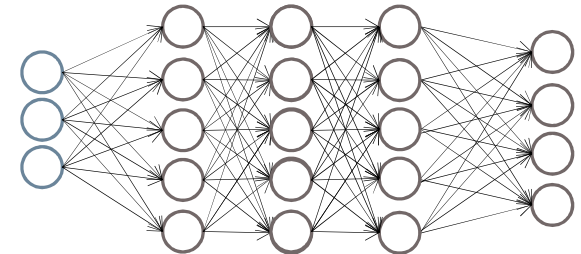Smaller Derivatives

Larger Learning Rate

2. Smaller derivatives, larger learning rate, and vice versa

Why?

# Not the whole story ......

- Adagrad [John Duchi, JMLR'11]
- RMSprop
    - https://www.youtube.com/watch?v=O3sxAc4hxZU
- Adadelta [Matthew D. Zeiler, arXiv'12]
- Adam [Diederik P. Kingma, ICLR'15]
- AdaSecant [Caglar Gulcehre, arXiv'14]
- "No more pesky learning rates" [Tom Schaul, arXiv'12]

Pick a network architecture (connectivity pattern between neurons)



No. of input units: Dimension of features $x^{(i)}$

No. output units: Number of classes

Reasonable default: 1 hidden layer, or if >1 hidden layer, have same no. of hidden units in every layer (usually the more the better)

Logistic regression:

$$J(\theta) = -\frac{1}{m}\left[\sum_{i=1}^{m} y^{(i)}\log h_\theta(x^{(i)}) + (1-y^{(i)})\log(1-h_\theta(x^{(i)}))\right] + \frac{\lambda}{2m}\sum_{j=1}^{n}\theta_j^2$$

**Regularization term**

Neural network:

$$h_\Theta(x) \in \mathbb{R}^K \quad (h_\Theta(x))_i = i^{th} \text{ output}$$

$$J(\Theta) = -\frac{1}{m}\left[\sum_{i=1}^{m}\sum_{k=1}^{K} y_k^{(i)}\log(h_\Theta(x^{(i)}))_k + (1-y_k^{(i)})\log(1-(h_\Theta(x^{(i)}))_k)\right]$$

$$+\frac{\lambda}{2m}\sum_{l=1}^{L-1}\sum_{i=1}^{s_l}\sum_{j=1}^{s_{l+1}}(\Theta_{ji}^{(l)})^2$$