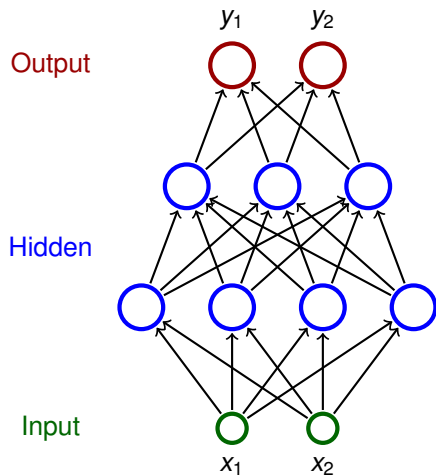MLP training – theory

# Architecture – Multilayer Perceptron (MLP)



- ▶ Neurons partitioned into **layers**; one input layer, one output layer, possibly several hidden layers
- ▶ layers numbered from 0; the input layer has number 0
  - ▶ E.g. three-layer network has two hidden layers and one output layer
- ▶ Neurons in the $i$-th layer are connected with all neurons in the $i + 1$-st layer
- ▶ Architecture of a MLP is typically described by numbers of neurons in individual layers (e.g. 2-4-3-2)

# MLP – architecture

**Notation:**

- ▶ Denote
  - ▶ $X$ a set of *input* neurons
  - ▶ $Y$ a set of *output* neurons
  - ▶ $Z$ a set of *all* neurons ($X, Y \subseteq Z$)

# MLP – architecture

**Notation:**

- ▶ Denote
    - ▶ $X$ a set of *input* neurons
    - ▶ $Y$ a set of *output* neurons
    - ▶ $Z$ a set of *all* neurons ($X, Y \subseteq Z$)
- ▶ individual neurons denoted by indices $i, j$ etc.
    - ▶ $\xi_j$ is the inner potential of the neuron $j$ *after the computation stops*

# MLP – architecture

**Notation:**

- Denote
    - $X$ a set of *input* neurons
    - $Y$ a set of *output* neurons
    - $Z$ a set of *all* neurons ($X, Y \subseteq Z$)
- individual neurons denoted by indices $i, j$ etc.
    - $\xi_j$ is the inner potential of the neuron $j$ *after the computation stops*
    - $y_j$ is the output of the neuron $j$ *after the computation stops*

    (define $y_0 = 1$ is the value of the formal unit input)

# MLP – architecture

**Notation:**

- ▶ Denote
    - ▶ $X$ a set of *input* neurons
    - ▶ $Y$ a set of *output* neurons
    - ▶ $Z$ a set of *all* neurons ($X, Y \subseteq Z$)
- ▶ individual neurons denoted by indices $i, j$ etc.
    - ▶ $\xi_j$ is the inner potential of the neuron $j$ *after the computation stops*
    - ▶ $y_j$ is the output of the neuron $j$ *after the computation stops*
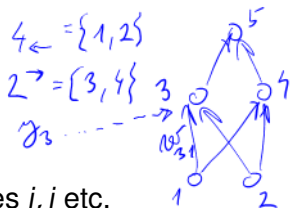
    (define $y_0 = 1$ is the value of the formal unit input)

- ▶ $w_{ji}$ is the weight of the connection **from $i$ to $j$**

    (in particular, $w_{j0}$ is the weight of the connection from the formal unit input, i.e. $w_{j0} = -b_j$ where $b_j$ is the bias of the neuron $j$)

# MLP – architecture

**Notation:**

- Denote
    - $X$ a set of *input* neurons
    - $Y$ a set of *output* neurons
    - $Z$ a set of *all* neurons ($X, Y \subseteq Z$)
- individual neurons denoted by indices $i, j$ etc.
    - $\xi_j$ is the inner potential of the neuron $j$ *after the computation stops*
    - $y_j$ is the output of the neuron $j$ *after the computation stops*

    (define $y_0 = 1$ is the value of the formal unit input)

- $w_{ji}$ is the weight of the connection **from** $i$ **to** $j$

    (in particular, $w_{j0}$ is the weight of the connection from the formal unit input, i.e. $w_{j0} = -b_j$ where $b_j$ is the bias of the neuron $j$)

- $j_{\leftarrow}$ is a set of all $i$ such that $j$ is adjacent from $i$
  (i.e. there is an arc **to** $j$ from $i$)

# MLP – architecture

**Notation:**

- Denote
  - $X$ a set of *input* neurons
  - $Y$ a set of *output* neurons
  - $Z$ a set of *all* neurons ($X, Y \subseteq Z$)
- individual neurons denoted by indices $i, j$ etc.
  - $\xi_j$ is the inner potential of the neuron $j$ *after the computation stops*
  - $y_j$ is the output of the neuron $j$ *after the computation stops*

  (define $y_0 = 1$ is the value of the formal unit input)

- $w_{ji}$ is the weight of the connection **from** $i$ **to** $j$

  (in particular, $w_{j0}$ is the weight of the connection from the formal unit input, i.e. $w_{j0} = -b_j$ where $b_j$ is the bias of the neuron $j$)

- $j_{\leftarrow}$ is a set of all $i$ such that $j$ is adjacent from $i$
  (i.e. there is an arc **to** $j$ from $i$)
- $j^{\rightarrow}$ is a set of all $i$ such that $j$ is adjacent to $i$
  (i.e. there is an arc **from** $j$ to $i$)

# MLP – activity

**Activity:**

▶ inner potential of neuron $j$:

$$\xi_j = \sum_{i \in j_\leftarrow} w_{ji} y_i$$

# MLP – activity

**Activity:**

▶ inner potential of neuron $j$:

$$\xi_j = \sum_{i \in j_\leftarrow} w_{ji} y_i$$

▶ activation function $\sigma_j$ for neuron $j$ (arbitrary differentiable)
[ e.g. logistic sigmoid $\sigma_j(\xi) = \frac{1}{1+e^{-\lambda_j \xi}}$ ]

# MLP – activity

**Activity:**

▶ inner potential of neuron $j$:

$$\xi_j = \sum_{i \in j_\leftarrow} w_{ji} y_i$$

▶ activation function $\sigma_j$ for neuron $j$ (arbitrary differentiable)
[ e.g. logistic sigmoid $\sigma_j(\xi) = \frac{1}{1+e^{-\lambda_j \xi}}$ ]

▶ State of non-input neuron $j \in Z \setminus X$ after the computation stops:
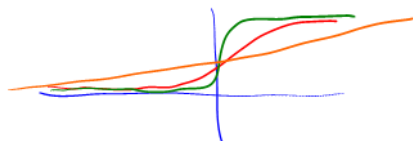
$$y_j = \sigma_j(\xi_j)$$

($y_j$ depends on the configuration $\vec{w}$ and the input $\vec{x}$, so we sometimes write $y_j(\vec{w}, \vec{x})$ )

# MLP – activity

**Activity:**

▶ inner potential of neuron $j$:

$$\xi_j = \sum_{i \in j_{\leftarrow}} w_{ji} y_i$$



▶ activation function $\sigma_j$ for neuron $j$ (arbitrary differentiable)
[ e.g. logistic sigmoid $\sigma_j(\xi) = \frac{1}{1 + e^{-\lambda_j \xi}}$ ]

▶ State of non-input neuron $j \in Z \setminus X$ after the computation stops:

$$y_j = \sigma_j(\xi_j)$$

($y_j$ depends on the configuration $\vec{w}$ and the input $\vec{x}$, so we sometimes write $y_j(\vec{w}, \vec{x})$ )

▶ The network computes a function $\mathbb{R}^{|X|}$ do $\mathbb{R}^{|Y|}$. Layer-wise computation: First, all input neurons are assigned values of the input. In the $\ell$-th step, all neurons of the $\ell$-th layer are evaluated.

## MLP – learning

**Learning:**

▶ Given a **training set** $\mathcal{T}$ of the form

$$\left\{ \left( \vec{x}_k, \vec{d}_k \right) \quad \middle| \quad k = 1, \ldots, p \right\}$$

Here, every $\vec{x}_k \in \mathbb{R}^{|X|}$ is an *input vector* end every $\vec{d}_k \in \mathbb{R}^{|Y|}$ is the desired network output. For every $j \in Y$, denote by $d_{kj}$ the desired output of the neuron $j$ for a given network input $\vec{x}_k$ (the vector $\vec{d}_k$ can be written as $\left( d_{kj} \right)_{j \in Y}$).

**Learning:**

▶ Given a **training set** $\mathcal{T}$ of the form

$$\left\{ \left( \vec{x}_k, \vec{d}_k \right) \quad \mid \quad k = 1, \ldots, p \right\}$$

Here, every $\vec{x}_k \in \mathbb{R}^{|X|}$ is an *input vector* end every $\vec{d}_k \in \mathbb{R}^{|Y|}$ is the desired network output. For every $j \in Y$, denote by $d_{kj}$ the desired output of the neuron $j$ for a given network input $\vec{x}_k$ (the vector $\vec{d}_k$ can be written as $\left( d_{kj} \right)_{j \in Y}$).

▶ **Error function:**

$$E(\vec{w}) = \sum_{k=1}^{p} E_k(\vec{w})$$

where

$$E_k(\vec{w}) = \frac{1}{2} \sum_{j \in Y} \left( y_j(\vec{w}, \vec{x}_k) - d_{kj} \right)^2$$

$$\mathcal{T} = \left\{ \left( \vec{x}_1, 6 \right), \left( \vec{x}_2, 7 \right) \right\}$$

$$E = E_1 + E_2$$

$$E_1 = \frac{1}{2} \left( y(\vec{w}, \vec{x}_1) - 6 \right)^2$$

$$E_2 = \frac{1}{2} \left( y(\vec{w}, \vec{x}_2) - 7 \right)^2$$

# MLP – learning algorithm

**Batch algorithm (gradient descent):**

The algorithm computes a sequence of weight vectors $\vec{w}^{(0)}, \vec{w}^{(1)}, \vec{w}^{(2)}, \ldots$.

- weights in $\vec{w}^{(0)}$ are randomly initialized to values close to 0
- in the step $t + 1$ (here $t = 0, 1, 2 \ldots$), weights $\vec{w}^{(t+1)}$ are computed as follows:

$$w_{ji}^{(t+1)} = w_{ji}^{(t)} + \Delta w_{ji}^{(t)}$$

# MLP – learning algorithm

**Batch algorithm (gradient descent):**

The algorithm computes a sequence of weight vectors
$\vec{w}^{(0)}, \vec{w}^{(1)}, \vec{w}^{(2)}, \ldots$.

- weights in $\vec{w}^{(0)}$ are randomly initialized to values close to 0
- in the step $t + 1$ (here $t = 0, 1, 2 \ldots$), weights $\vec{w}^{(t+1)}$ are computed as follows:

$$w_{ji}^{(t+1)} = w_{ji}^{(t)} + \Delta w_{ji}^{(t)}$$

where

$$\Delta w_{ji}^{(t)} = -\varepsilon(t) \cdot \frac{\partial E}{\partial w_{ji}}(\vec{w}^{(t)})$$

is a *weight update* of $w_{ji}$ in step $t + 1$ and $0 < \varepsilon(t) \leq 1$ is a *learning rate* in step $t + 1$.

# MLP – learning algorithm

**Batch algorithm (gradient descent):**

The algorithm computes a sequence of weight vectors
$\vec{w}^{(0)}, \vec{w}^{(1)}, \vec{w}^{(2)}, \ldots$.

- ▶ weights in $\vec{w}^{(0)}$ are randomly initialized to values close to 0
- ▶ in the step $t + 1$ (here $t = 0, 1, 2 \ldots$), weights $\vec{w}^{(t+1)}$ are computed as follows:

$$w_{ji}^{(t+1)} = w_{ji}^{(t)} + \Delta w_{ji}^{(t)}$$

where

$$\Delta w_{ji}^{(t)} = -\varepsilon(t) \cdot \frac{\partial E}{\partial w_{ji}}(\vec{w}^{(t)})$$

is a *weight update* of $w_{ji}$ in step $t + 1$ and $0 < \varepsilon(t) \leq 1$ is a *learning rate* in step $t + 1$.

Note that $\frac{\partial E}{\partial w_{ji}}(\vec{w}^{(t)})$ is a component of the gradient $\nabla E$, i.e. the weight update can be written as $\vec{w}^{(t+1)} = \vec{w}^{(t)} - \varepsilon(t) \cdot \nabla E(\vec{w}^{(t)})$.

# MLP – error function gradient

For every $w_{ji}$ we have

$$\frac{\partial E}{\partial w_{ji}} = \sum_{k=1}^{p} \frac{\partial E_k}{\partial w_{ji}}$$

# MLP – error function gradient

For every $w_{ji}$ we have

$$\frac{\partial E}{\partial w_{ji}} = \sum_{k=1}^{p} \frac{\partial E_k}{\partial w_{ji}}$$

where for every $k = 1, \ldots, p$ holds

$$\frac{\partial E_k}{\partial w_{ji}} = \frac{\partial E_k}{\partial y_j} \cdot \sigma'_j(\xi_j) \cdot y_i$$

For $j \in Y$: $\dfrac{\partial E_k}{\partial y_j} = \dfrac{\partial\left(\frac{1}{2}\left(\sum_{n \in Y}(y_n - d_{kn})^2\right)\right.}{\partial y_j} = y_j - d_{kj}$

For $j \notin Y$: $\dfrac{\partial E_k}{\partial y_j} = \sum_{n \in j^{\rightarrow}} \dfrac{\partial E_k}{\partial y_n} \cdot \dfrac{\partial y_n}{\partial \xi_n} \cdot \dfrac{\partial \xi_n}{\partial y_j}$

$\qquad = \sum_{n \in j^{\rightarrow}} \dfrac{\partial E_k}{\partial y_n} \cdot \sigma_n'(\xi_n) \cdot w_{nj}$

since $\dfrac{\partial y_n}{\partial \xi_n} = \dfrac{\partial(\sigma_n(\xi_n))}{\partial \xi_n} = \sigma_n'(\xi_n)$

$\dfrac{\partial \xi_n}{\partial y_j} = \dfrac{\partial\left(\sum_{s \in n^{\leftarrow}} w_{ns} \cdot y_s\right)}{\partial y_j} = w_{nj}$

For every $w_{ji}$ we have

$$\frac{\partial E}{\partial w_{ji}} = \sum_{k=1}^{p} \frac{\partial E_k}{\partial w_{ji}}$$

where for every $k = 1, \ldots, p$ holds

$$\frac{\partial E_k}{\partial w_{ji}} = \frac{\partial E_k}{\partial y_j} \cdot \sigma_j'(\xi_j) \cdot y_i$$

and for every $j \in Z \smallsetminus X$ we get

$$\frac{\partial E_k}{\partial y_j} = y_j - d_{kj} \qquad\qquad \text{for } j \in Y$$

$$\frac{\partial E_k}{\partial w_{ji}} = \frac{\partial E_k}{\partial y_{\partial}} \cdot \frac{\partial y_{\partial}}{\partial \xi_{\partial}} \cdot \frac{\partial \xi_{\partial}}{\partial w_{ji}} = \frac{\partial E_k}{\partial y_j} \cdot \sigma'_{\partial}(\xi_{\partial}) \cdot y_i$$

since $\dfrac{\partial y_{\partial}}{\partial \xi_{\partial}} = \dfrac{\partial (\sigma_{\partial}(\xi_{\partial}))}{\partial \xi_{\partial}} = \sigma'_{i}(\xi_{\partial})$

$$\frac{\partial \xi_{\partial}}{\partial w_{ji}} = \frac{\partial \left( \sum_{r \in \partial^{\leftarrow}} w_{\partial r} \cdot y_r \right)}{\partial w_{ji}} = y_i$$

For every $w_{ji}$ we have

$$\frac{\partial E}{\partial w_{ji}} = \sum_{k=1}^{p} \frac{\partial E_k}{\partial w_{ji}}$$

where for every $k = 1, \ldots, p$ holds

$$\frac{\partial E_k}{\partial w_{ji}} = \frac{\partial E_k}{\partial y_j} \cdot \sigma'_j(\xi_j) \cdot y_i$$

and for every $j \in Z \smallsetminus X$ we get

$$\frac{\partial E_k}{\partial y_j} = y_j - d_{kj} \qquad \text{for } j \in Y$$

$$\frac{\partial E_k}{\partial y_j} = \sum_{r \in j^{\rightarrow}} \frac{\partial E_k}{\partial y_r} \cdot \sigma'_r(\xi_r) \cdot w_{rj} \qquad \text{for } j \in Z \smallsetminus (Y \cup X)$$

(Here all $y_j$ are in fact $y_j(\vec{w}, \vec{x}_k)$).

## MLP – error function gradient

- If $\sigma_j(\xi) = \frac{1}{1+e^{-\lambda_j \xi}}$ for all $j \in Z$, then

$$\sigma_j'(\xi_j) = \lambda_j y_j (1 - y_j)$$

# MLP – error function gradient

▶ If $\sigma_j(\xi) = \frac{1}{1+e^{-\lambda_j \xi}}$ for all $j \in Z$, then

$$\sigma_j'(\xi_j) = \lambda_j y_j (1 - y_j)$$

and thus for all $j \in Z \smallsetminus X$:

$$\frac{\partial E_k}{\partial y_j} = y_j - d_{kj} \qquad\qquad \text{for } j \in Y$$

$$\frac{\partial E_k}{\partial y_j} = \sum_{r \in j^\rightarrow} \frac{\partial E_k}{\partial y_r} \cdot \lambda_r y_r (1 - y_r) \cdot w_{rj} \quad \text{for } j \in Z \smallsetminus (Y \cup X)$$

# MLP – error function gradient

▶ If $\sigma_j(\xi) = \frac{1}{1+e^{-\lambda_j \xi}}$ for all $j \in Z$, then
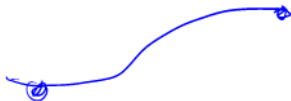
$$\sigma_j'(\xi_j) = \lambda_j y_j (1 - y_j)$$

and thus for all $j \in Z \smallsetminus X$:

$$\frac{\partial E_k}{\partial y_j} = y_j - d_{kj} \qquad\qquad \text{for } j \in Y$$

$$\frac{\partial E_k}{\partial y_j} = \sum_{r \in j^\rightarrow} \frac{\partial E_k}{\partial y_r} \cdot \lambda_r y_r (1 - y_r) \cdot w_{rj} \quad \text{for } j \in Z \smallsetminus (Y \cup X)$$

▶ If $\sigma_j(\xi) = a \cdot \tanh(b \cdot \xi_j)$ for all $j \in Z$, then

$$\sigma_j'(\xi_j) = \frac{b}{a}(a - y_j)(a + y_j)$$

Compute $\frac{\partial E}{\partial w_{ji}} = \sum_{k=1}^{p} \frac{\partial E_k}{\partial w_{ji}}$ as follows:

## MLP – computing the gradient

Compute $\frac{\partial E}{\partial w_{ji}} = \sum_{k=1}^{p} \frac{\partial E_k}{\partial w_{ji}}$ as follows:

Initialize $\mathcal{E}_{ji} := 0$
(By the end of the computation: $\mathcal{E}_{ji} = \frac{\partial E}{\partial w_{ji}}$)

## MLP – computing the gradient

Compute $\frac{\partial E}{\partial w_{ji}} = \sum_{k=1}^{p} \frac{\partial E_k}{\partial w_{ji}}$ as follows:

Initialize $\mathcal{E}_{ji} := 0$
(By the end of the computation: $\mathcal{E}_{ji} = \frac{\partial E}{\partial w_{ji}}$)

For every $k = 1, \ldots, p$ do:

## MLP – computing the gradient

Compute $\frac{\partial E}{\partial w_{ji}} = \sum_{k=1}^{p} \frac{\partial E_k}{\partial w_{ji}}$ as follows:

Initialize $\mathcal{E}_{ji} := 0$
(By the end of the computation: $\mathcal{E}_{ji} = \frac{\partial E}{\partial w_{ji}}$)

For every $k = 1, \ldots, p$ do:

    **1. forward pass:** compute $y_j = y_j(\vec{w}, \vec{x}_k)$ for all $j \in Z$

## MLP – computing the gradient

Compute $\frac{\partial E}{\partial w_{ji}} = \sum_{k=1}^{p} \frac{\partial E_k}{\partial w_{ji}}$ as follows:

Initialize $\mathcal{E}_{ji} := 0$
(By the end of the computation: $\mathcal{E}_{ji} = \frac{\partial E}{\partial w_{ji}}$)

For every $k = 1, \ldots, p$ do:

1. **forward pass:** compute $y_j = y_j(\vec{w}, \vec{x}_k)$ for all $j \in Z$

2. **backward pass:** compute $\frac{\partial E_k}{\partial y_j}$ for all $j \in Z$ using *backpropagation* (see the next slide!)

## MLP – computing the gradient

Compute $\frac{\partial E}{\partial w_{ji}} = \sum_{k=1}^{p} \frac{\partial E_k}{\partial w_{ji}}$ as follows:

Initialize $\mathcal{E}_{ji} := 0$
(By the end of the computation: $\mathcal{E}_{ji} = \frac{\partial E}{\partial w_{ji}}$)

For every $k = 1, \ldots, p$ do:

1. **forward pass:** compute $y_j = y_j(\vec{w}, \vec{x}_k)$ for all $j \in Z$

2. **backward pass:** compute $\frac{\partial E_k}{\partial y_j}$ for all $j \in Z$ using *backpropagation* (see the next slide!)

3. compute $\frac{\partial E_k}{\partial w_{ji}}$ for all $w_{ji}$ using

$$\frac{\partial E_k}{\partial w_{ji}} := \frac{\partial E_k}{\partial y_j} \cdot \sigma_j'(\xi_j) \cdot y_i$$

## MLP – computing the gradient

Compute $\frac{\partial E}{\partial w_{ji}} = \sum_{k=1}^{p} \frac{\partial E_k}{\partial w_{ji}}$ as follows:

Initialize $\mathcal{E}_{ji} := 0$
(By the end of the computation: $\mathcal{E}_{ji} = \frac{\partial E}{\partial w_{ji}}$)

For every $k = 1, \ldots, p$ do:

1. **forward pass:** compute $y_j = y_j(\vec{w}, \vec{x}_k)$ for all $j \in Z$

2. **backward pass:** compute $\frac{\partial E_k}{\partial y_j}$ for all $j \in Z$ using *backpropagation* (see the next slide!)

3. compute $\frac{\partial E_k}{\partial w_{ji}}$ for all $w_{ji}$ using

$$\frac{\partial E_k}{\partial w_{ji}} := \frac{\partial E_k}{\partial y_j} \cdot \sigma'_j(\xi_j) \cdot y_i$$

4. $\mathcal{E}_{ji} := \mathcal{E}_{ji} + \frac{\partial E_k}{\partial w_{ji}}$

The resulting $\mathcal{E}_{ji}$ equals $\frac{\partial E}{\partial w_{ji}}$.

## MLP – backpropagation

Compute $\frac{\partial E_k}{\partial y_j}$ for all $j \in Z$ as follows:

# MLP – backpropagation

Compute $\frac{\partial E_k}{\partial y_j}$ for all $j \in Z$ as follows:

- if $j \in Y$, then $\frac{\partial E_k}{\partial y_j} = y_j - d_{kj}$

## MLP – backpropagation

Compute $\frac{\partial E_k}{\partial y_j}$ for all $j \in Z$ as follows:

▶ if $j \in Y$, then $\frac{\partial E_k}{\partial y_j} = y_j - d_{kj}$

▶ if $j \in Z \smallsetminus Y \cup X$, then assuming that $j$ is in the $\ell$-th layer and assuming that $\frac{\partial E_k}{\partial y_r}$ has already been computed for all neurons in the $\ell + 1$-st layer, compute

$$\frac{\partial E_k}{\partial y_j} = \sum_{r \in j^{\rightarrow}} \frac{\partial E_k}{\partial y_r} \cdot \sigma_r'(\xi_r) \cdot w_{rj}$$

(This works because all neurons of $r \in j^{\rightarrow}$ belong to the $\ell + 1$-st layer.)

## Complexity of the batch algorithm

Computation of $\frac{\partial E}{\partial w_{ji}}(\vec{w}^{(t-1)})$ stops in time linear in the size of the network plus the size of the training set.

(assuming unit cost of operations including computation of $\sigma'_r(\xi_r)$ for given $\xi_r$)

## Complexity of the batch algorithm

Computation of $\frac{\partial E}{\partial w_{ji}}(\vec{w}^{(t-1)})$ stops in time linear in the size of the network plus the size of the training set.

(assuming unit cost of operations including computation of $\sigma'_r(\xi_r)$ for given $\xi_r$)

**Proof sketch:** The algorithm does the following $p$ times:

## Complexity of the batch algorithm

Computation of $\frac{\partial E}{\partial w_{ji}}(\vec{w}^{(t-1)})$ stops in time linear in the size of the network plus the size of the training set.

(assuming unit cost of operations including computation of $\sigma'_r(\xi_r)$ for given $\xi_r$)

**Proof sketch:** The algorithm does the following *p* times:
1. forward pass, i.e. computes $y_j(\vec{w}, \vec{x}_k)$

## Complexity of the batch algorithm

Computation of $\frac{\partial E}{\partial w_{ji}}(\vec{w}^{(t-1)})$ stops in time linear in the size of the network plus the size of the training set.
(assuming unit cost of operations including computation of $\sigma'_r(\xi_r)$ for given $\xi_r$)

**Proof sketch:** The algorithm does the following *p* times:

  **1.** forward pass, i.e. computes $y_j(\vec{w}, \vec{x}_k)$
  **2.** backpropagation, i.e. computes $\frac{\partial E_k}{\partial y_j}$

## Complexity of the batch algorithm

Computation of $\frac{\partial E}{\partial w_{ji}}(\vec{w}^{(t-1)})$ stops in time linear in the size of the network plus the size of the training set.

(assuming unit cost of operations including computation of $\sigma'_r(\xi_r)$ for given $\xi_r$)

**Proof sketch:** The algorithm does the following *p* times:

1. forward pass, i.e. computes $y_j(\vec{w}, \vec{x}_k)$

2. backpropagation, i.e. computes $\frac{\partial E_k}{\partial y_j}$

3. computes $\frac{\partial E_k}{\partial w_{ji}}$ and adds it to $\mathcal{E}_{ji}$ (a constant time operation in the unit cost framework)

## Complexity of the batch algorithm

Computation of $\frac{\partial E}{\partial w_{ji}}(\vec{w}^{(t-1)})$ stops in time linear in the size of the network plus the size of the training set.

(assuming unit cost of operations including computation of $\sigma'_r(\xi_r)$ for given $\xi_r$)

**Proof sketch:** The algorithm does the following *p* times:

1. forward pass, i.e. computes $y_j(\vec{w}, \vec{x}_k)$

2. backpropagation, i.e. computes $\frac{\partial E_k}{\partial y_j}$

3. computes $\frac{\partial E_k}{\partial w_{ji}}$ and adds it to $\mathcal{E}_{ji}$ (a constant time operation in the unit cost framework)

The steps 1. - 3. take linear time.

## Complexity of the batch algorithm

Computation of $\frac{\partial E}{\partial w_{ji}}(\vec{w}^{(t-1)})$ stops in time linear in the size of the network plus the size of the training set.

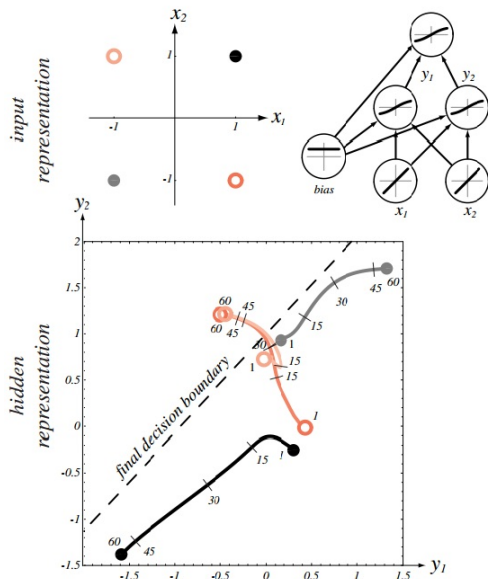(assuming unit cost of operations including computation of $\sigma'_r(\xi_r)$ for given $\xi_r$)

**Proof sketch:** The algorithm does the following *p* times:

1. forward pass, i.e. computes $y_j(\vec{w}, \vec{x}_k)$

2. backpropagation, i.e. computes $\frac{\partial E_k}{\partial y_j}$

3. computes $\frac{\partial E_k}{\partial w_{ji}}$ and adds it to $\mathcal{E}_{ji}$ (a constant time operation in the unit cost framework)

The steps 1. - 3. take linear time.

Note that the speed of convergence of the gradient descent cannot be estimated ...

# Illustration of the gradient descent – XOR

# MLP – learning algorithm

**Online algorithm:**

The algorithm computes a sequence of weight vectors
$\vec{w}^{(0)}, \vec{w}^{(1)}, \vec{w}^{(2)}, \ldots$.

- ▶ weights in $\vec{w}^{(0)}$ are randomly initialized to values close to 0
- ▶ in the step $t + 1$ (here $t = 0, 1, 2 \ldots$), weights $\vec{w}^{(t+1)}$ are computed as follows:

$$w_{ji}^{(t+1)} = w_{ji}^{(t)} + \Delta w_{ji}^{(t)}$$

where

$$\Delta w_{ji}^{(t)} = -\varepsilon(t) \cdot \frac{\partial E_k}{\partial w_{ji}}(w_{ji}^{(t)})$$

is the *weight update* of $w_{ji}$ in the step $t + 1$ and $0 < \varepsilon(t) \leq 1$ is the *learning rate* in the step $t + 1$.

There are other variants determined by selection of the training examples used for the error computation (more on this later).

# SGD

- ▶ weights in $\vec{w}^{(0)}$ are randomly initialized to values close to 0
- ▶ in the step $t + 1$ (here $t = 0, 1, 2 \ldots$), weights $\vec{w}^{(t+1)}$ are computed as follows:

    - ▶ Choose (randomly) a set of training examples $T \subseteq \{1, \ldots, p\}$
    - ▶ Compute

      $$\vec{w}^{(t+1)} = \vec{w}^{(t)} + \Delta \vec{w}^{(t)}$$

      where

      $$\Delta \vec{w}^{(t)} = -\varepsilon(t) \cdot \sum_{k \in T} \nabla E_k(\vec{w}^{(t)})$$

- ▶ $0 < \varepsilon(t) \leq 1$ is a *learning rate* in step $t + 1$
- ▶ $\nabla E_k(\vec{w}^{(t)})$ is the gradient of the error of the example $k$

Note that the random choice of the minibatch is typically implemented by randomly shuffling all data and then choosing minibatches sequentially.