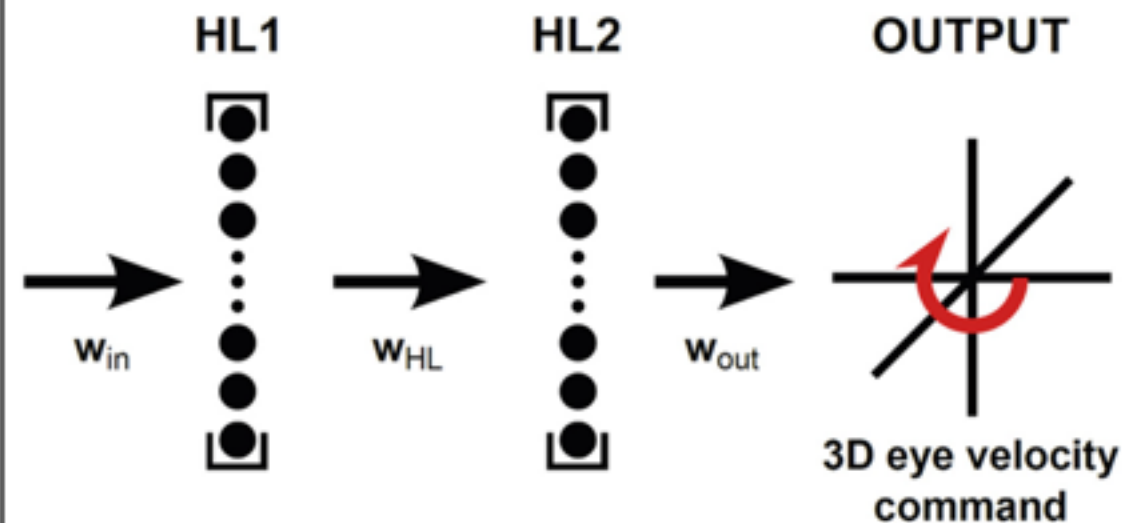
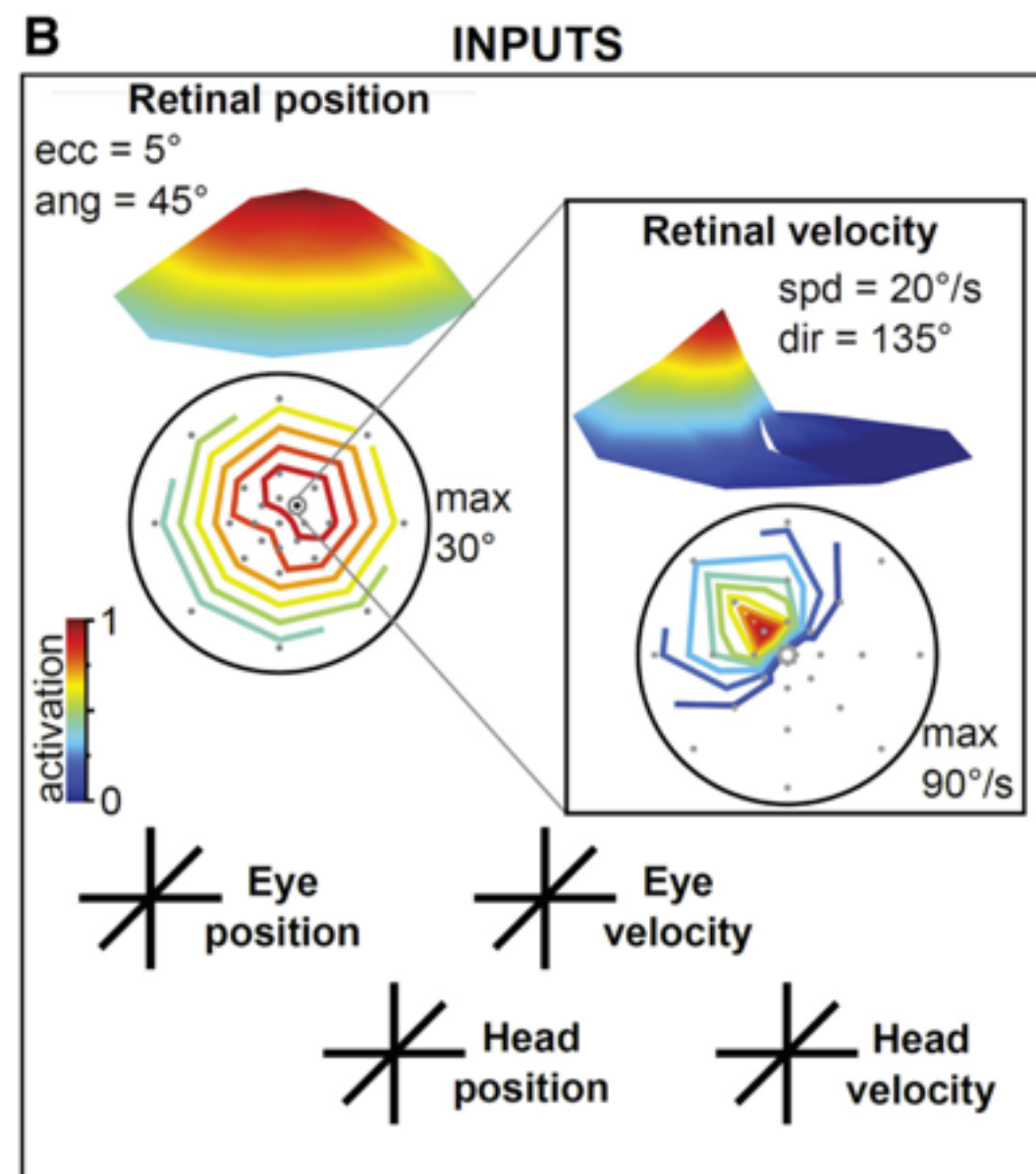


# Rate-based artificial neural networks and error backpropagation learning

Scott Murdison

Machine learning journal club

May 16, 2016



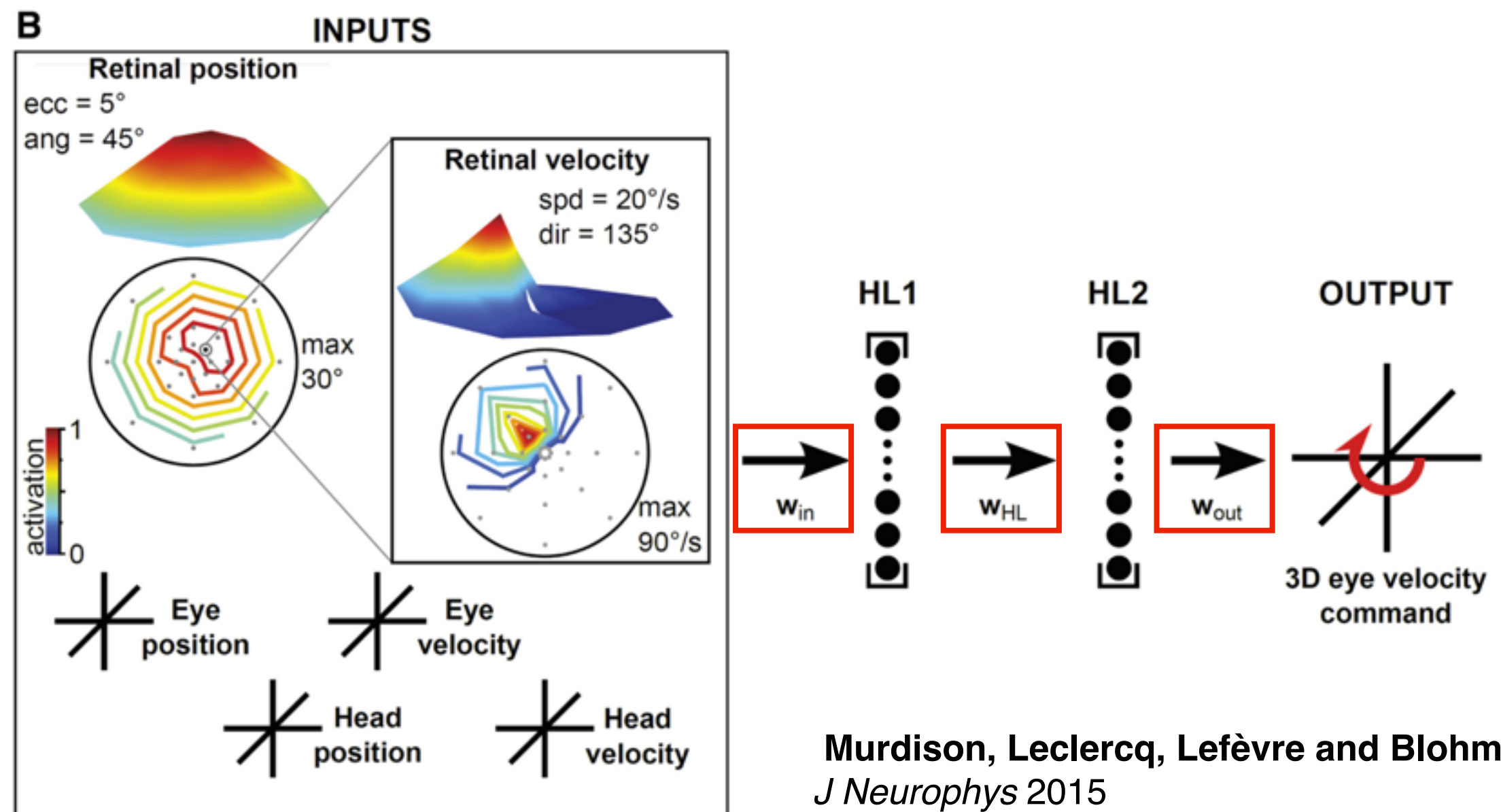
Murdison, Leclercq, Lefèvre and Blohm  
*J Neurophys* 2015

# Rate-based artificial neural networks and error backpropagation learning

Scott Murdison

Machine learning journal club

May 16, 2016



# Neural networks??? Why would we use those?!

# Neural networks??? Why would we use those?!

Some problems are just too complicated (i.e. nonlinear) to solve with simple programming...

# Neural networks??? Why would we use those?!

Some problems are just too complicated (i.e. nonlinear) to solve with simple programming...

# Neural networks??? Why would we use those?!

Some problems are just too complicated (i.e. nonlinear) to solve with simple programming...

## Computer vision/audition applications:

- digit recognition
- facial recognition
- reading
- speech recognition
- object recognition

# Neural networks??? Why would we use those?!

Some problems are just too complicated (i.e. nonlinear) to solve with simple programming...

## Computer vision/audition applications:

- digit recognition
- facial recognition
- reading
- speech recognition
- object recognition

## Others:

- detection of medical disorders
- industrial process control
- stock market predictions
- the list goes on...

# Neural networks??? Why would we use those?!

Some problems are just too complicated (i.e. nonlinear) to solve with simple programming...

## Computer vision/audition applications:

- digit recognition
- facial recognition
- reading
- speech recognition
- object recognition

## Others:

- detection of medical disorders
- industrial process control
- stock market predictions
- the list goes on...



# Neural networks??? Why would we use those?!

Some problems are just too complicated (i.e. nonlinear) to solve with simple programming...

## Computer vision/audition applications:

- digit recognition
- facial recognition
- reading
- speech recognition
- object recognition

## Others:

- detection of medical disorders
- industrial process control
- stock market predictions
- the list goes on...

Could theoretically write a program to solve each of these problems, but the rules governing such a program would be incredibly complicated!

# Neural networks??? Why would we use those?!

Some problems are just too complicated (i.e. nonlinear) to solve with simple programming...

## Computer vision/audition applications:

- digit recognition
- facial recognition
- reading
- speech recognition
- object recognition

## Others:

- detection of medical disorders
- industrial process control
- stock market predictions
- the list goes on...

Could theoretically write a program to solve each of these problems, but the rules governing such a program would be incredibly complicated!

## Geoff Hinton sweet intro vid

# Neural networks??? Why would we use those?!

Some problems are just too complicated to solve with simple programming...

## Computer vision/audition applications:

- facial recognition

# Neural networks??? Why would we use those?!

Some problems are just too complicated to solve with simple programming...

## Computer vision/audition applications:

- facial recognition

In the Fall of 1992, for a class project in Artificial Intelligence, I designed a neural network to locate facial features in images. The one hundred images I used came from the underclassmen section of the 1987 [University High School](#) yearbook. They were scanned in at 96 by 128 resolution. I set four of the images aside to comprise the testing set, and for the remaining ninety-six I manually specified the coordinates of the left eye, right eye, nose, and mouth.

Paul Debevec. *A Neural Network for Facial Feature Location*. UC Berkeley CS283 Project Report, December 1992. <http://www.debevec.org/FaceRecognition/>

# Neural networks??? Why would we use those?!

Some problems are just too complicated to solve with simple programming...

## Computer vision/audition applications:

- facial recognition



*Training Set Image*

manually located left eye,  
nose and mouth

# Neural networks??? Why would we use those?!

Some problems are just too complicated to solve with simple programming...

## Computer vision/audition applications:

- facial recognition



*Training Set Image*

manually located left eye,  
nose and mouth



*Log-polar maps of left eye, nose, and mouth*

polar-transformed images  
around each feature

# Neural networks??? Why would we use those?!

Some problems are just too complicated to solve with simple programming...

## Computer vision/audition applications:

- facial recognition



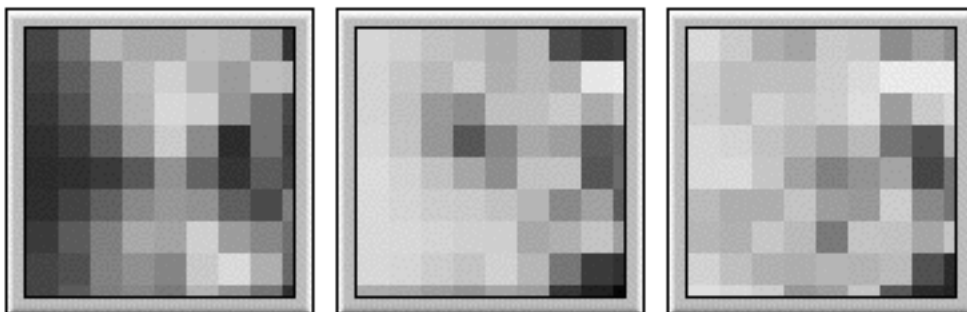
*Training Set Image*

manually located left eye,  
nose and mouth



*Log-polar maps of left eye, nose, and mouth*

polar-transformed images  
around each feature



*8 by 8 subsamples of the above maps*

subsampled images to  
feed to neural network ...Why?

# Neural networks??? Why would we use those?!

Some problems are just too complicated to solve with simple programming...

## Computer vision/audition applications:

- facial recognition



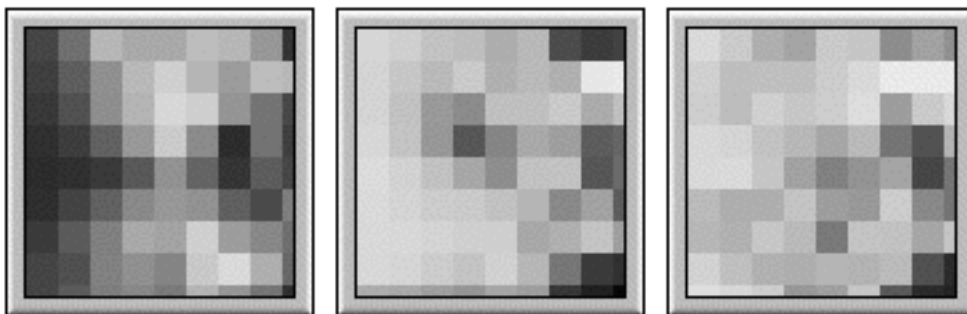
manually located left eye,  
nose and mouth

*Training Set Image*



*Log-polar maps of left eye, nose, and mouth*

polar-transformed images  
around each feature



*8 by 8 subsamples of the above maps*

subsampled images to  
feed to neural network ...Why?

A. Avoids local minima



# Neural networks??? Why would we use those?!

Some problems are just too complicated to solve with simple programming...

## Computer vision/audition applications:

- facial recognition



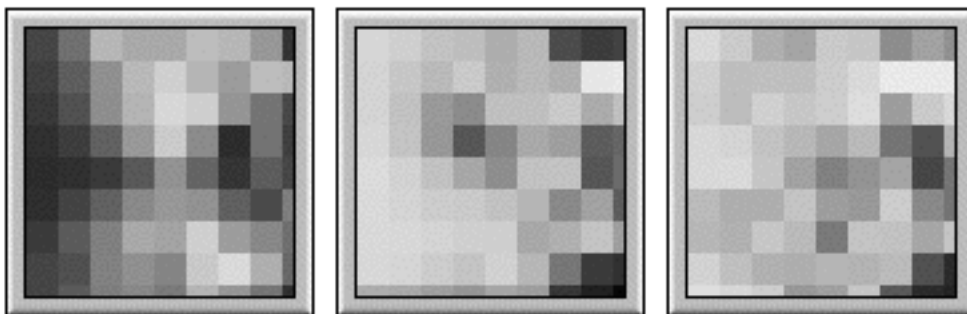
*Training Set Image*

manually located left eye,  
nose and mouth



*Log-polar maps of left eye, nose, and mouth*

polar-transformed images  
around each feature



*8 by 8 subsamples of the above maps*

subsampled images to  
feed to neural network ...Why?

A. Avoids local minima

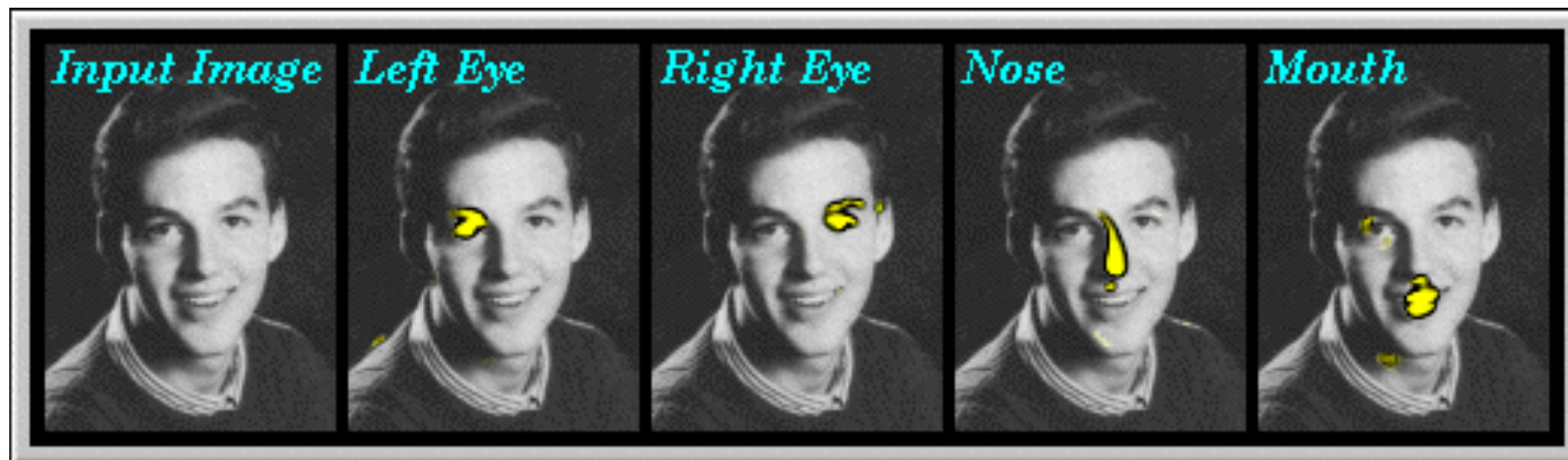
B. 4MB RAM in 1992 = \$150 USD

# Neural networks??? Why would we use those?!

Some problems are just too complicated to solve with simple programming...

## Computer vision/audition applications:

- facial recognition



*Neural network outputs for a previously unseen face*

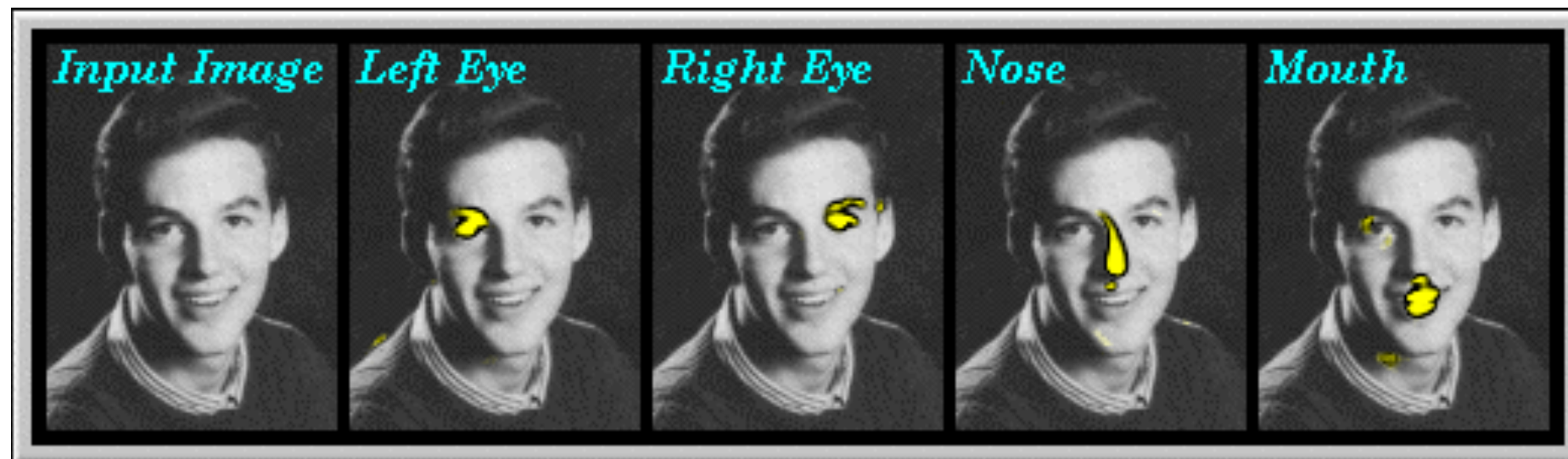
Paul Debevec. *A Neural Network for Facial Feature Location*. UC Berkeley CS283 Project Report, December 1992. <http://www.debevec.org/FaceRecognition/>

# Neural networks??? Why would we use those?!

Some problems are just too complicated to solve with simple programming...

## Computer vision/audition applications:

- facial recognition



*Neural network outputs for a previously unseen face*

After training a simple, feedforward network with backprop using yearbook photos it could successfully detect each eye, nose and mouth in previously unseen photo!

Paul Debevec. *A Neural Network for Facial Feature Location*. UC Berkeley CS283 Project Report, December 1992. <http://www.debevec.org/FaceRecognition/>

# The choice for rate-based over spiking for machine learning

# The choice for rate-based over spiking for machine learning



# The choice for rate-based over spiking for machine learning



## **Pros**

neuron-level resolution

time-resolved

(spiking dynamics, variability)

several levels of abstraction available

(H-H, leaky integrate-and-fire, Izhikevich)

# The choice for rate-based over spiking for machine learning



## **Pros**

neuron-level resolution

time-resolved

(spiking dynamics, variability)

several levels of abstraction available

(H-H, leaky integrate-and-fire, Izhikevich)

## **Cons**

computationally expensive for large

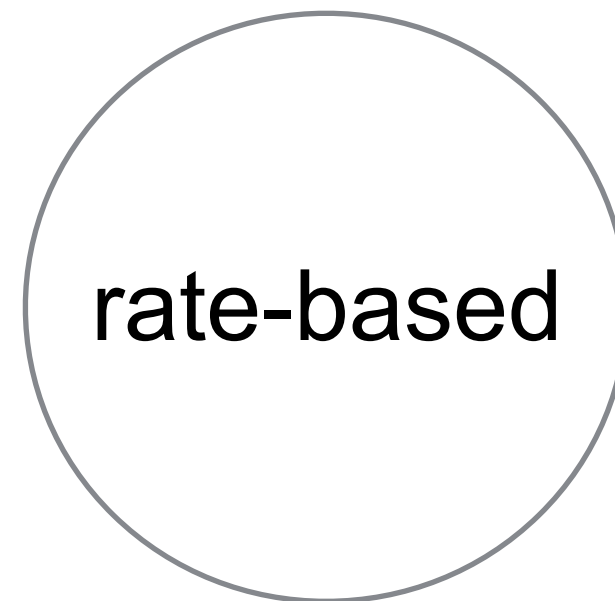
populations of neurons

intractable for simulations of whole brain

areas

**not obviously useful for machine learning**

# The choice for rate-based over spiking for machine learning



## **Pros**

neuron-level resolution

time-resolved

(spiking dynamics, variability)

several levels of abstraction available

(H-H, leaky integrate-and-fire, Izhikevich)

## **Cons**

computationally expensive for large

populations of neurons

intractable for simulations of whole brain

areas

**not obviously useful for machine learning**



# The choice for rate-based over spiking for machine learning



spiking

## **Pros**

neuron-level resolution

time-resolved

(spiking dynamics, variability)

several levels of abstraction available

(H-H, leaky integrate-and-fire, Izhikevich)

## **Cons**

computationally expensive for large

populations of neurons

intractable for simulations of whole brain

areas

**not obviously useful for machine learning**



rate-based

## **Pros**

describes average firing of functional

populations of neurons

computationally cheap

can address network structure/function

mathematical simplicity

**Universal function approximator**

# The choice for rate-based over spiking for machine learning



spiking

## Pros

neuron-level resolution  
time-resolved  
(spiking dynamics, variability)  
several levels of abstraction available  
(H-H, leaky integrate-and-fire, Izhikevich)

## Cons

computationally expensive for large  
populations of neurons  
intractable for simulations of whole brain  
areas  
**not obviously useful for machine learning**



rate-based

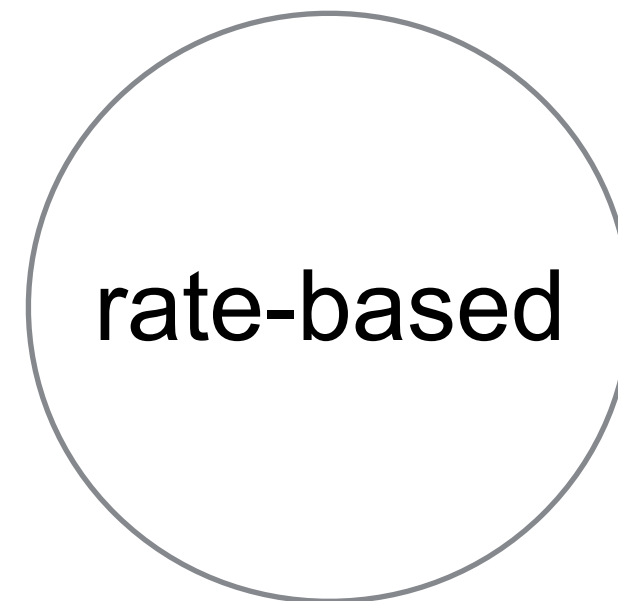
## Pros

describes average firing of functional  
populations of neurons  
computationally cheap  
can address network structure/function  
mathematical simplicity  
**Universal function approximator**

## Cons

lose benefits of spiking spatiotemporal  
resolution  
averages over timing/dynamics/variability/  
etc.  
biological analogs not always obvious

# The choice for rate-based over spiking for machine learning



## **Pros**

describes average firing of functional populations of neurons  
computationally cheap  
can address network structure/function  
mathematical simplicity

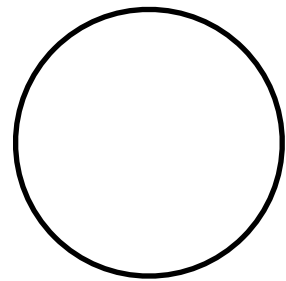
**Universal function approximator**

## **Cons**

lose benefits of spiking spatiotemporal resolution  
    averages over timing/dynamics/variability/  
    etc.  
biological analogs not always obvious

# General architecture

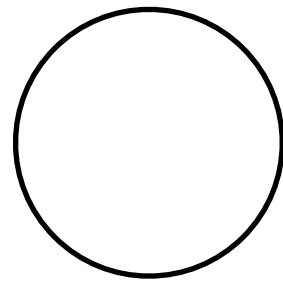
**Simplest network is the perceptron**



# General architecture

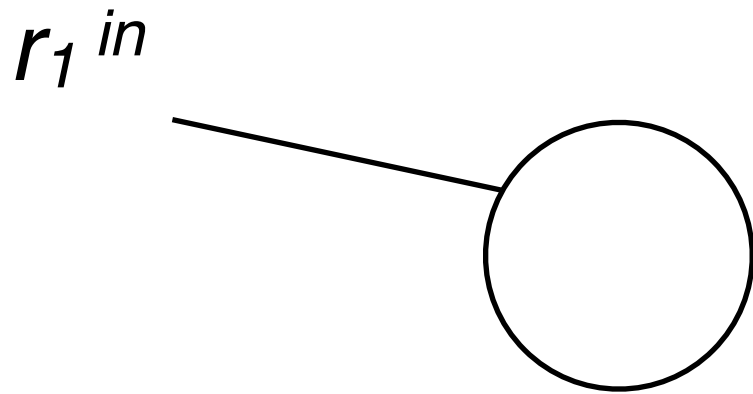
**Simplest network is the perceptron**

$r_1^{in}$



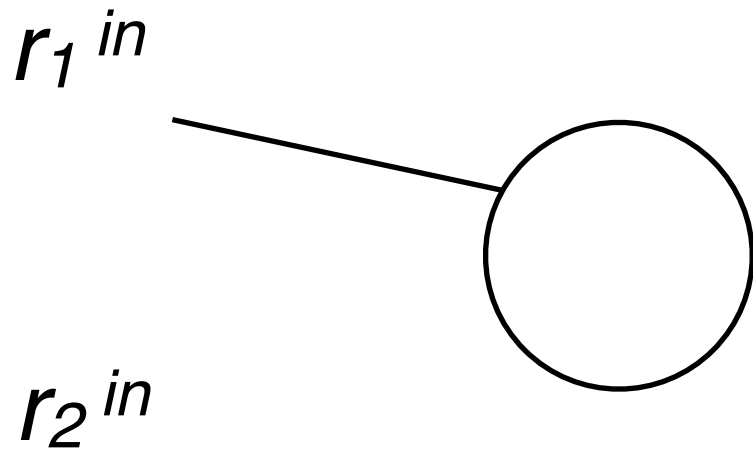
# General architecture

**Simplest network is the perceptron**



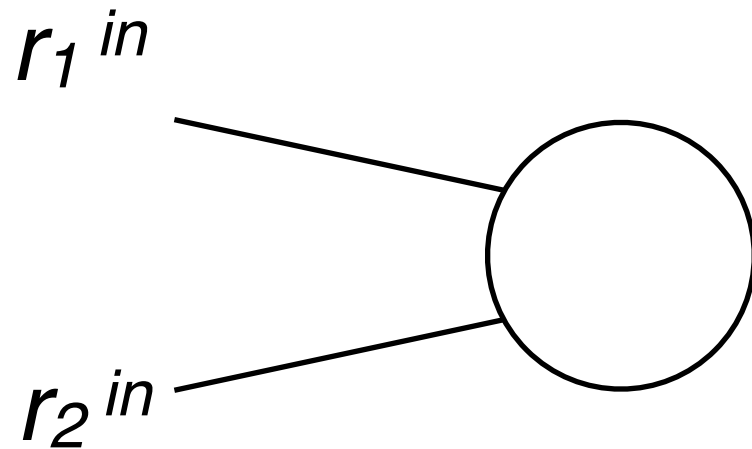
# General architecture

Simplest network is the perceptron



# General architecture

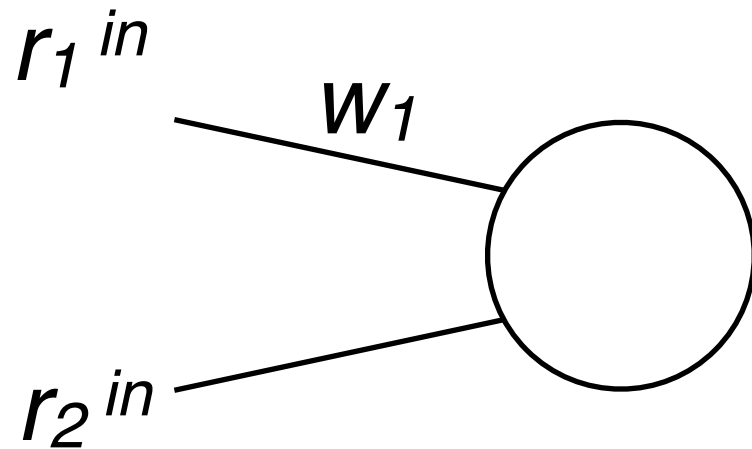
Simplest network is the perceptron





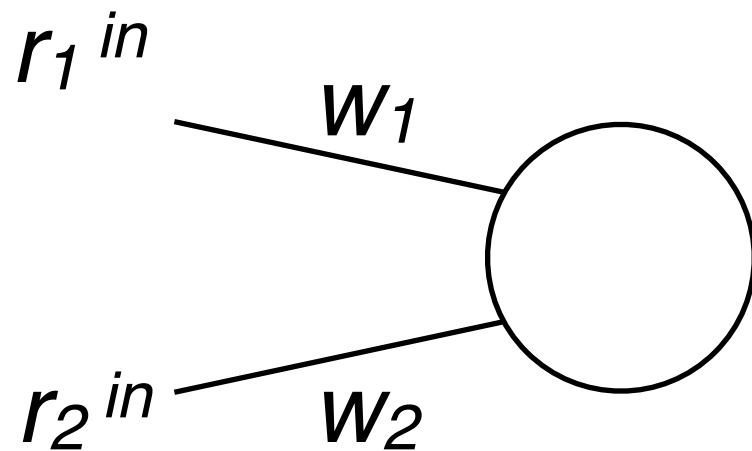
# General architecture

Simplest network is the perceptron



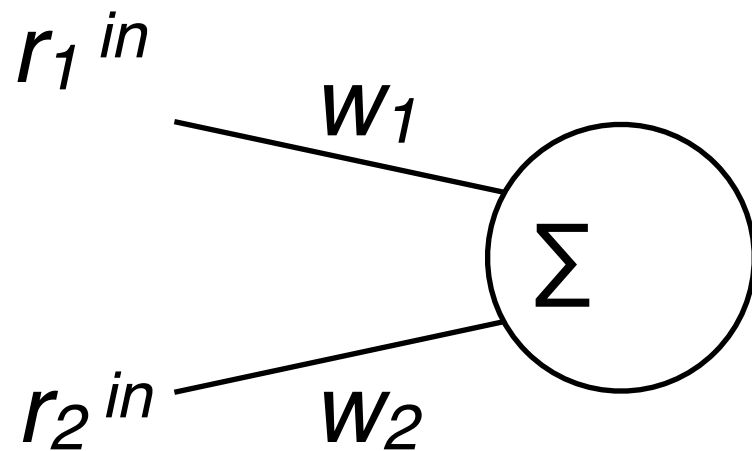
# General architecture

Simplest network is the perceptron



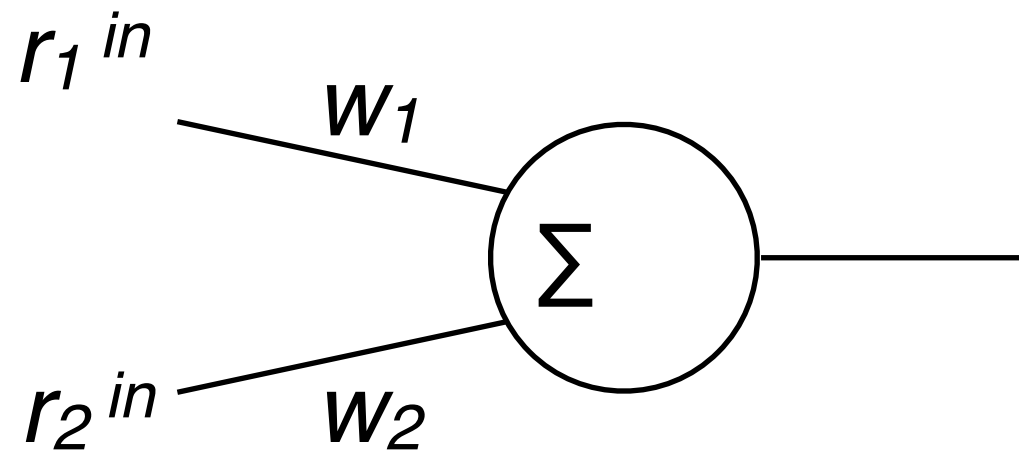
# General architecture

Simplest network is the perceptron



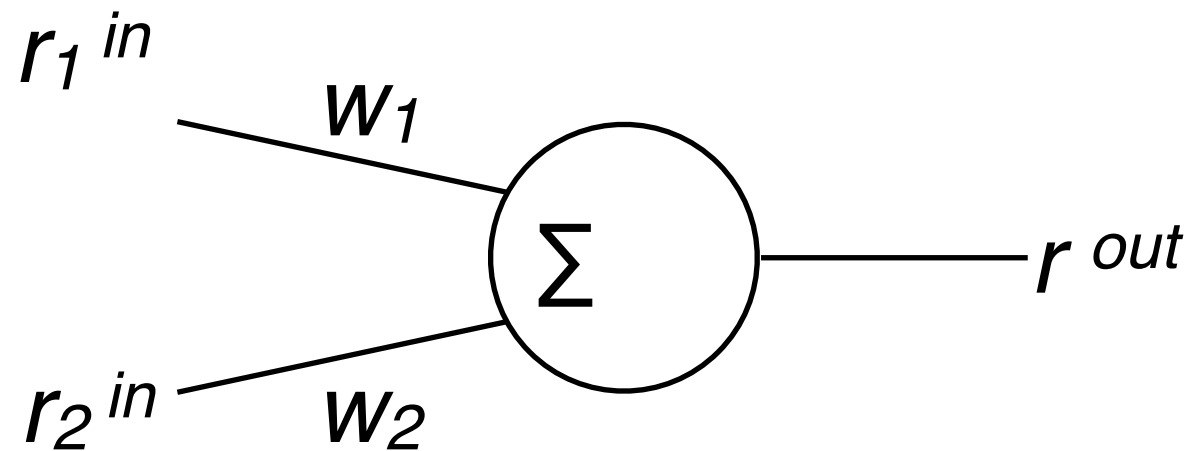
# General architecture

Simplest network is the perceptron



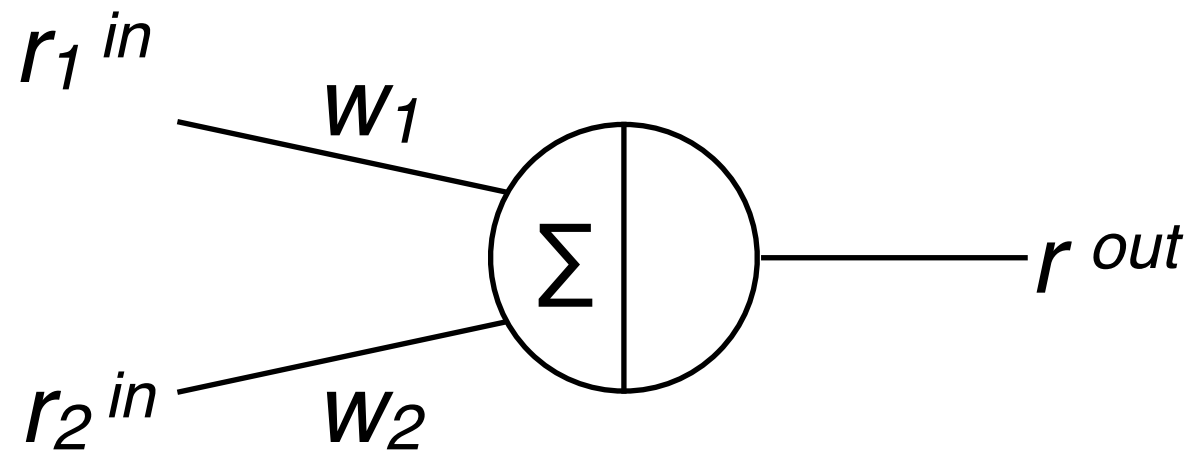
# General architecture

Simplest network is the perceptron



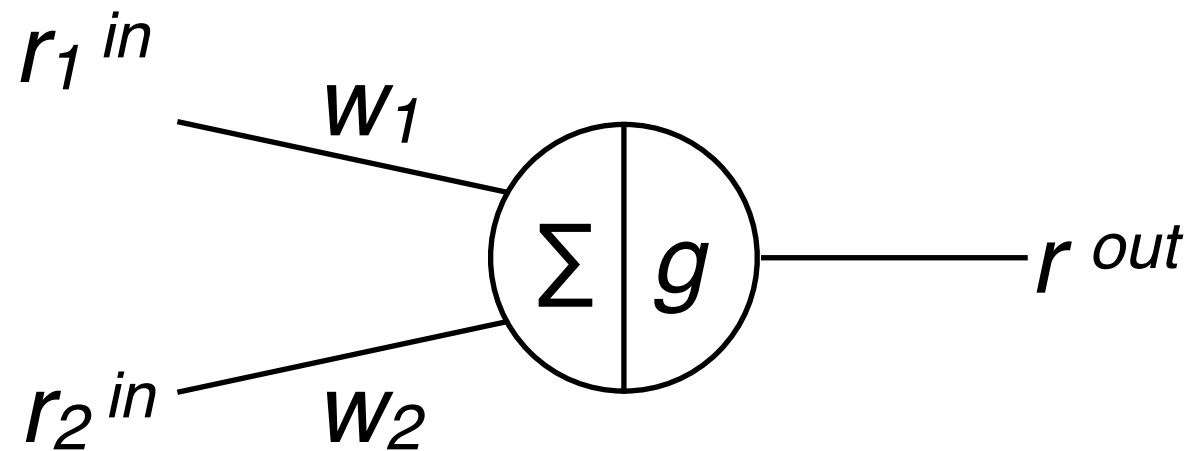
# General architecture

Simplest network is the perceptron



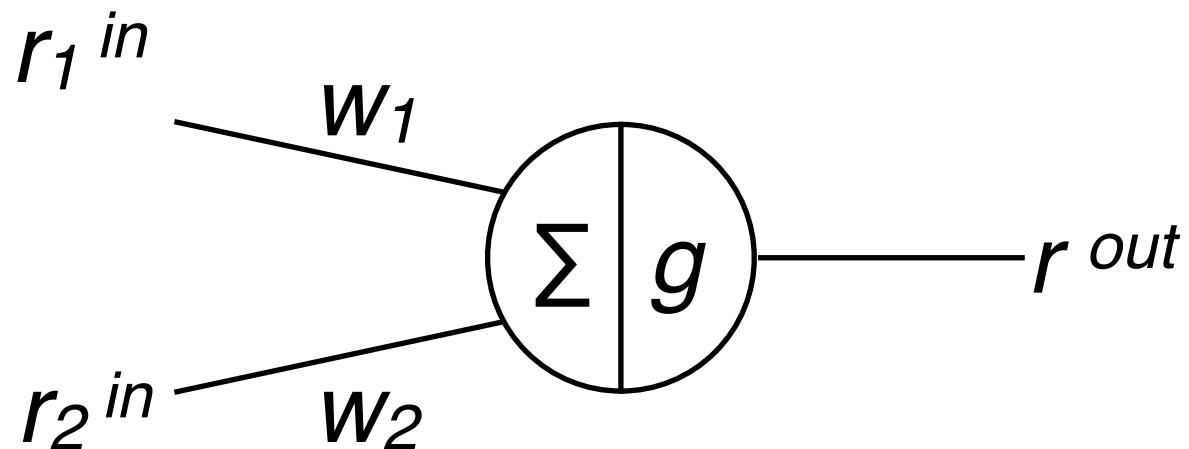
# General architecture

Simplest network is the perceptron



# General architecture

Simplest network is the perceptron



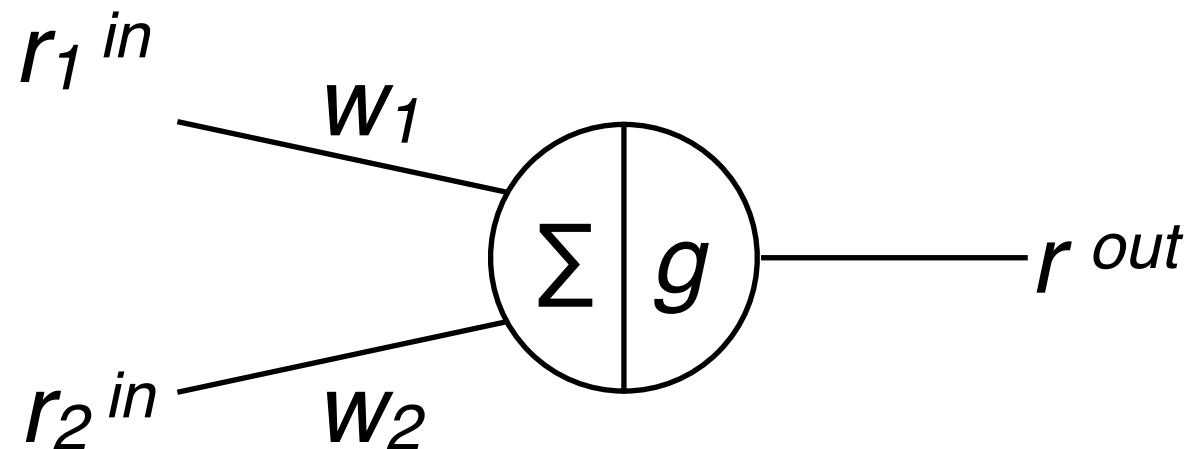
**Transfer function  $g(\Sigma)$**

general engineering term used to describe output of some processing unit as a function of input



# General architecture

Simplest network is the perceptron



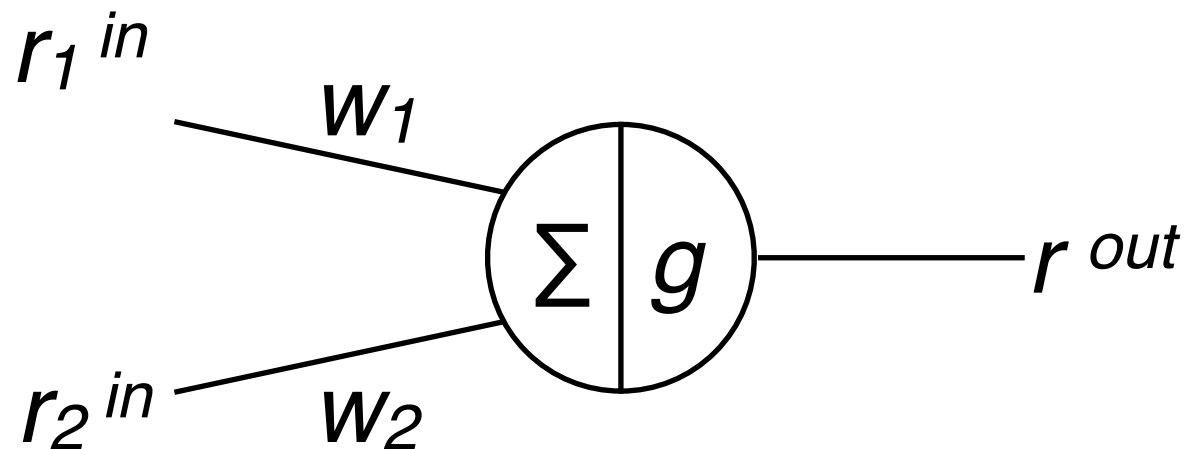
**Transfer function  $g(\Sigma)$**

general engineering term used to describe output of some processing unit as a function of input

if  $r_i^{in} = x_i$  and  $r^{out} = y$

# General architecture

Simplest network is the perceptron



**Transfer function  $g(\Sigma)$**

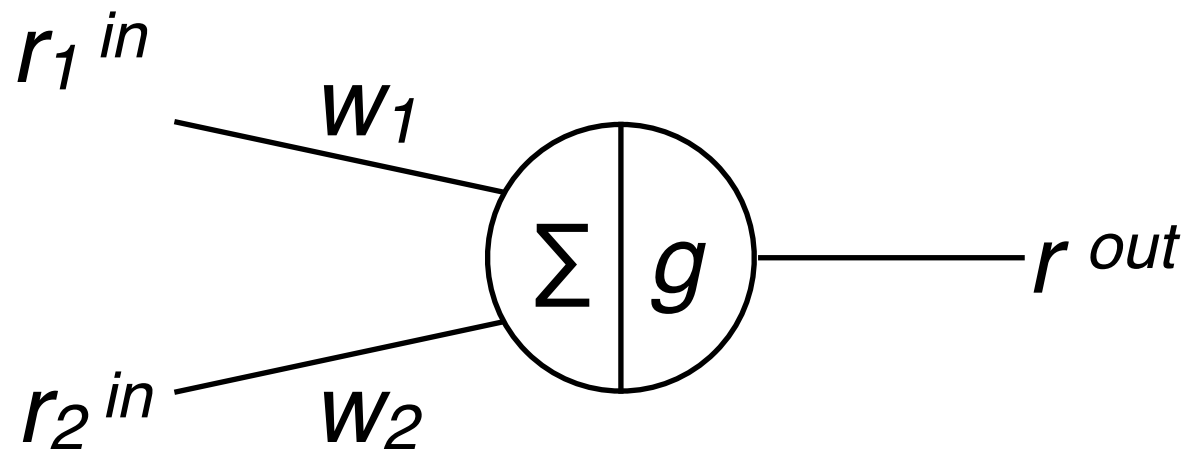
general engineering term used to describe output of some processing unit as a function of input

if  $r_i^{in} = x_i$  and  $r^{out} = y$

then  $y = g(w_1 \cdot x_1 + w_2 \cdot x_2)$

# General architecture

Simplest network is the perceptron



**Transfer function  $g(\Sigma)$**

general engineering term used to describe output of some processing unit as a function of input

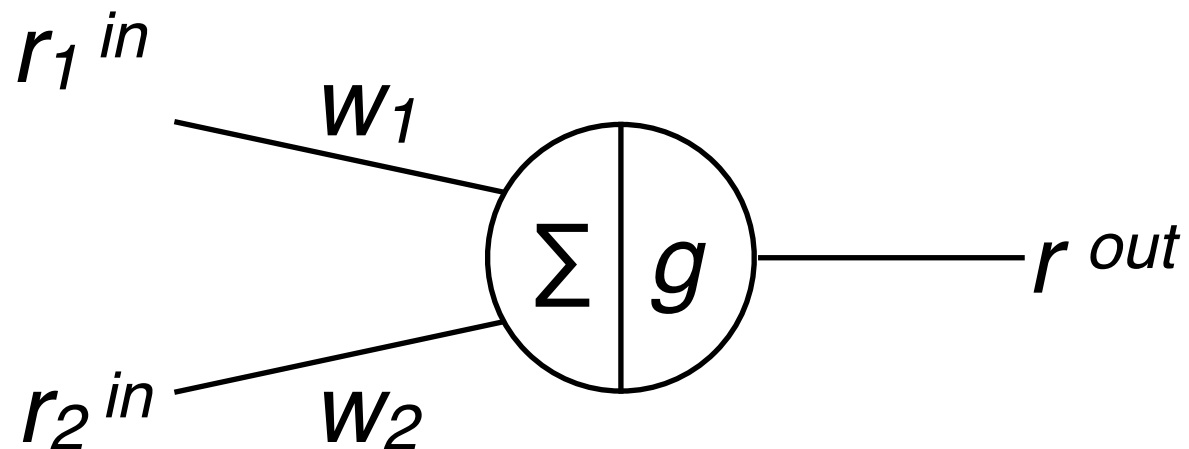
if  $r_i^{in} = x_i$  and  $r^{out} = y$

then  $y = g(w_1 \cdot x_1 + w_2 \cdot x_2)$

and if  $g(\Sigma) = \Sigma$  (i.e.  $g$  is purely linear)

# General architecture

Simplest network is the perceptron



**Transfer function  $g(\Sigma)$**

general engineering term used to describe output of some processing unit as a function of input

if  $r_i^{in} = x_i$  and  $r^{out} = y$

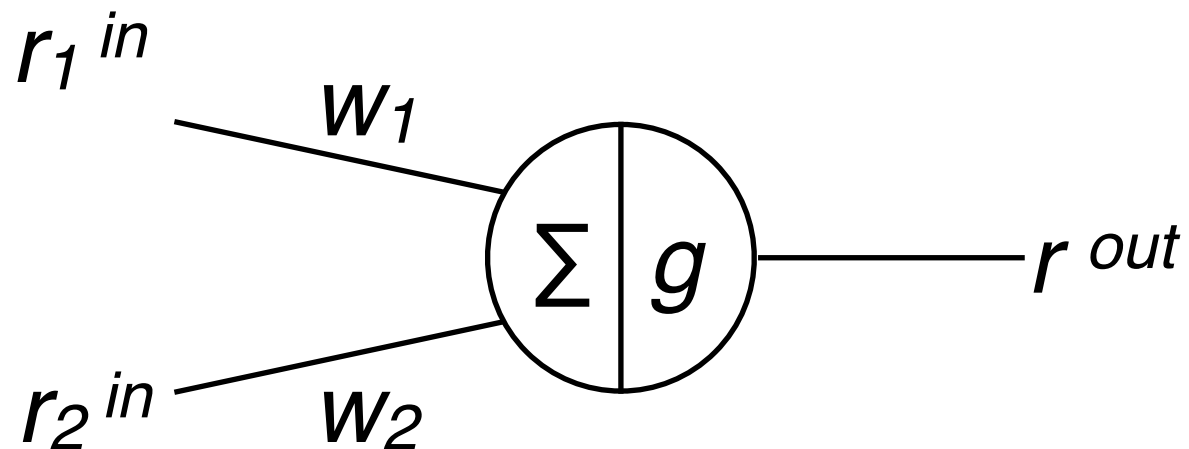
then  $y = g(w_1 \cdot x_1 + w_2 \cdot x_2)$

and if  $g(\Sigma) = \Sigma$  (i.e.  $g$  is purely linear)

then  $y = w_1 \cdot x_1 + w_2 \cdot x_2$

# General architecture

Simplest network is the perceptron



## Transfer function $g(\Sigma)$

general engineering term used to describe output of some processing unit as a function of input

if  $r_i^{in} = x_i$  and  $r^{out} = y$

then  $y = g(w_1 \cdot x_1 + w_2 \cdot x_2)$

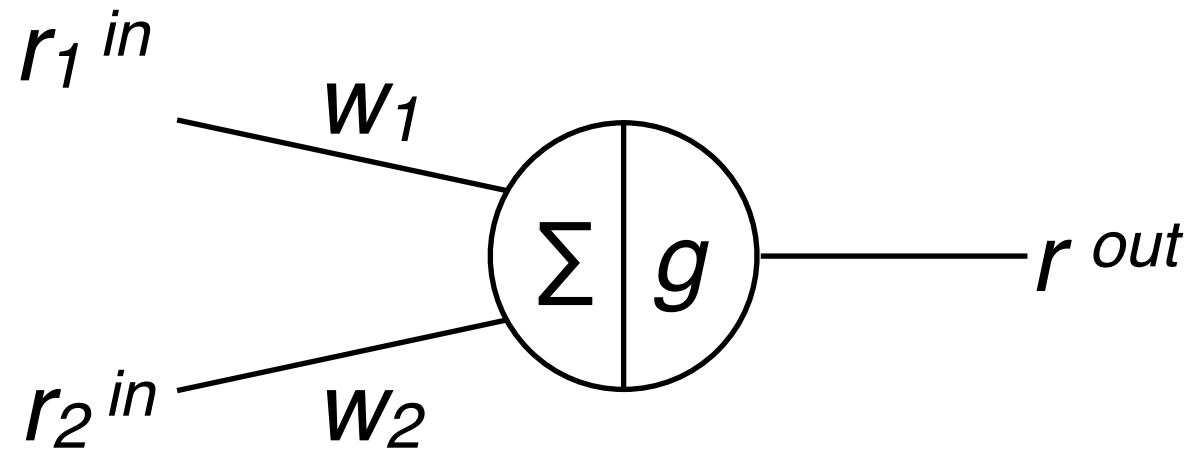
and if  $g(\Sigma) = \Sigma$  (i.e.  $g$  is purely linear)

then  $y = w_1 \cdot x_1 + w_2 \cdot x_2$

## General single-layer mapping

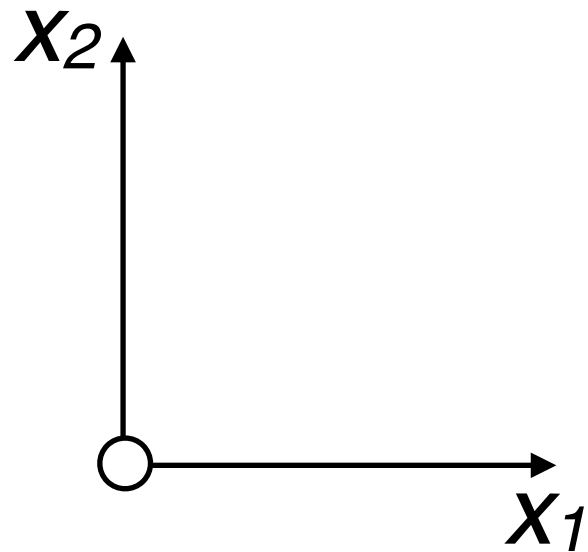
$$r_i^{out} = g\left(\sum_j w_{ij} r_j^{in}\right) \Leftrightarrow \mathbf{r}^{out} = g(\mathbf{W}\mathbf{r}^{in})$$

# The XOR problem

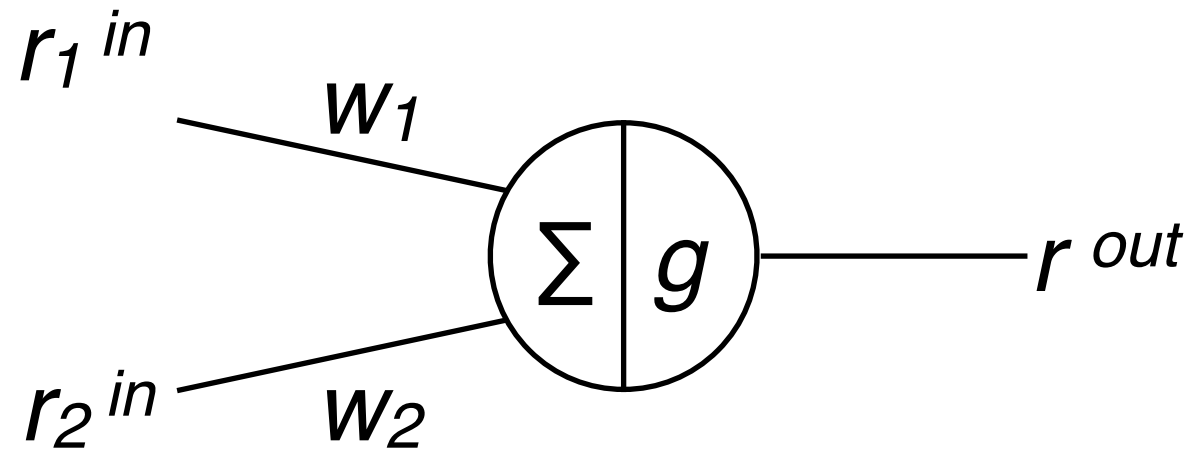


## Boolean OR

| $x_1$ | $x_2$ | $y$ |
|-------|-------|-----|
| 0     | 0     | 0   |

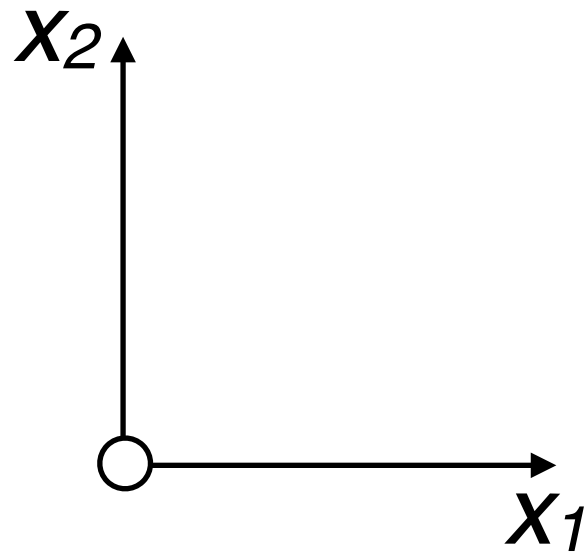


# The XOR problem

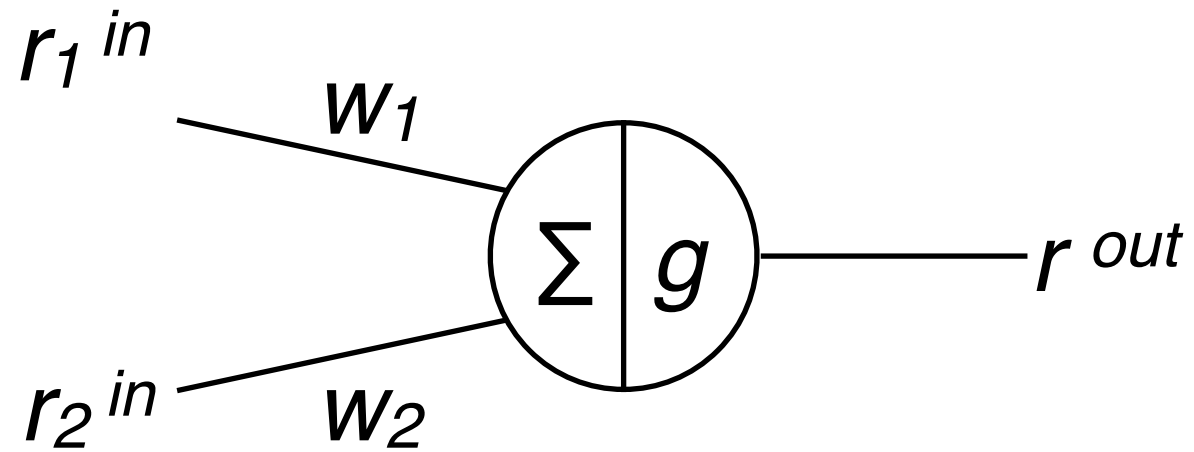


## Boolean OR

| $x_1$ | $x_2$ | $y$ |
|-------|-------|-----|
| 0     | 0     | 0   |
| 0     | 1     | 1   |

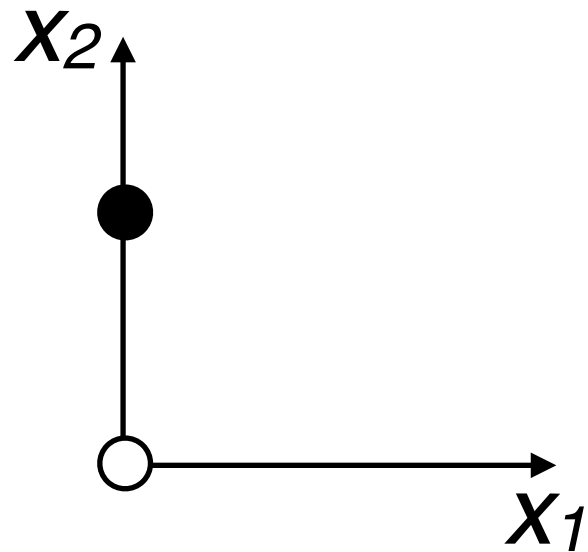


# The XOR problem



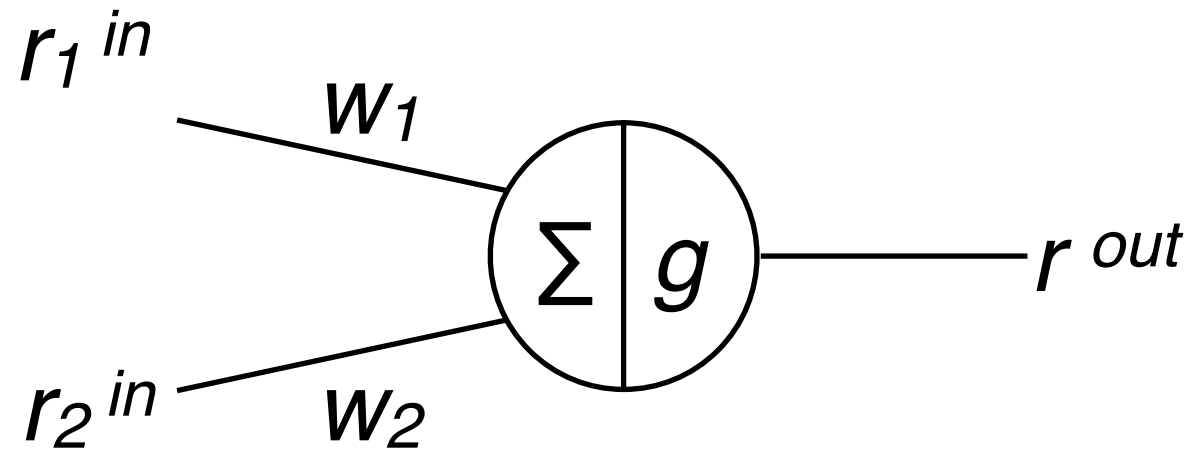
## Boolean OR

| $x_1$ | $x_2$ | $y$ |
|-------|-------|-----|
| 0     | 0     | 0   |
| 0     | 1     | 1   |



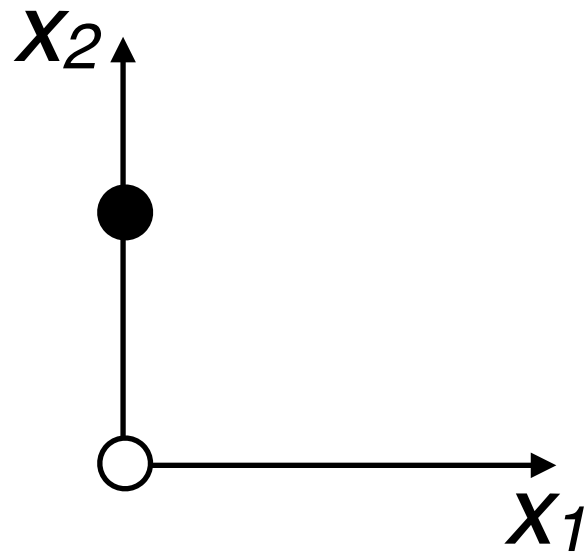


# The XOR problem

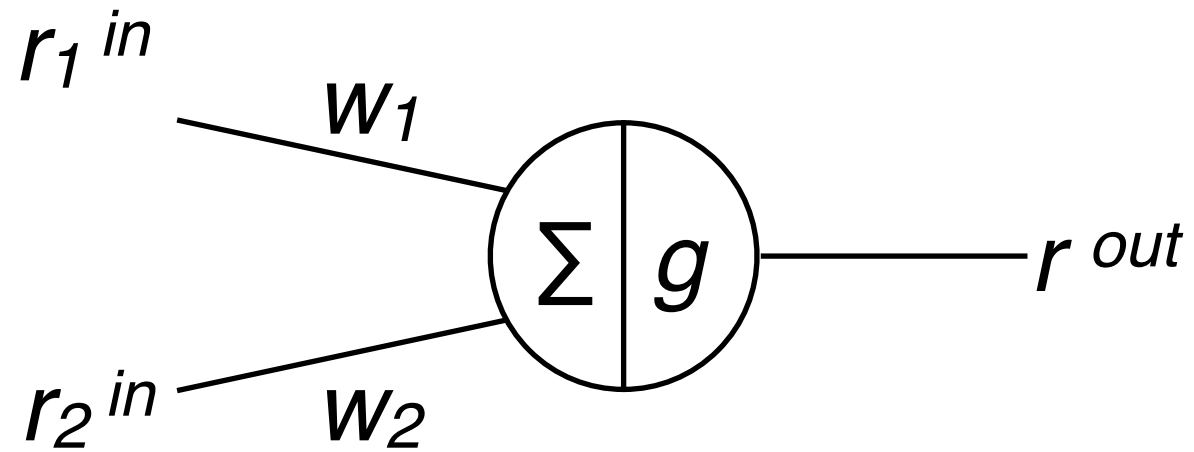


## Boolean OR

| $x_1$ | $x_2$ | $y$ |
|-------|-------|-----|
| 0     | 0     | 0   |
| 0     | 1     | 1   |
| 1     | 0     | 1   |

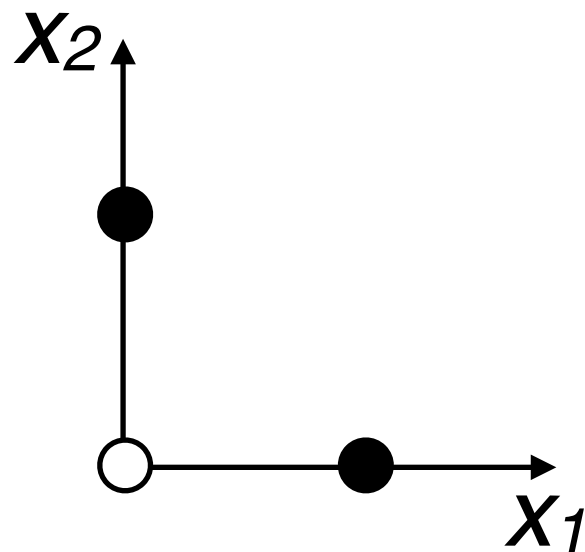


# The XOR problem

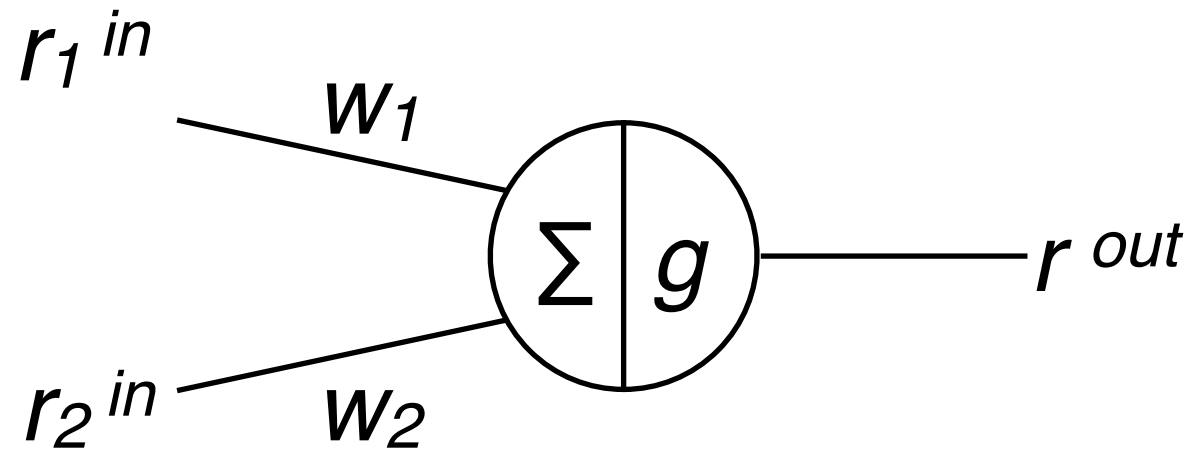


## Boolean OR

| $x_1$ | $x_2$ | $y$ |
|-------|-------|-----|
| 0     | 0     | 0   |
| 0     | 1     | 1   |
| 1     | 0     | 1   |

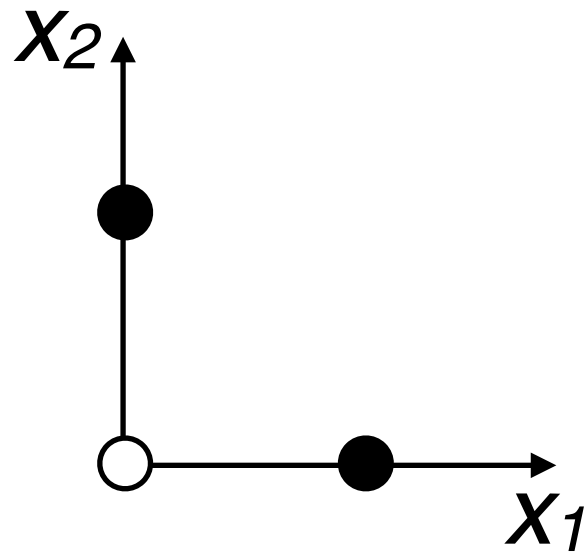


# The XOR problem

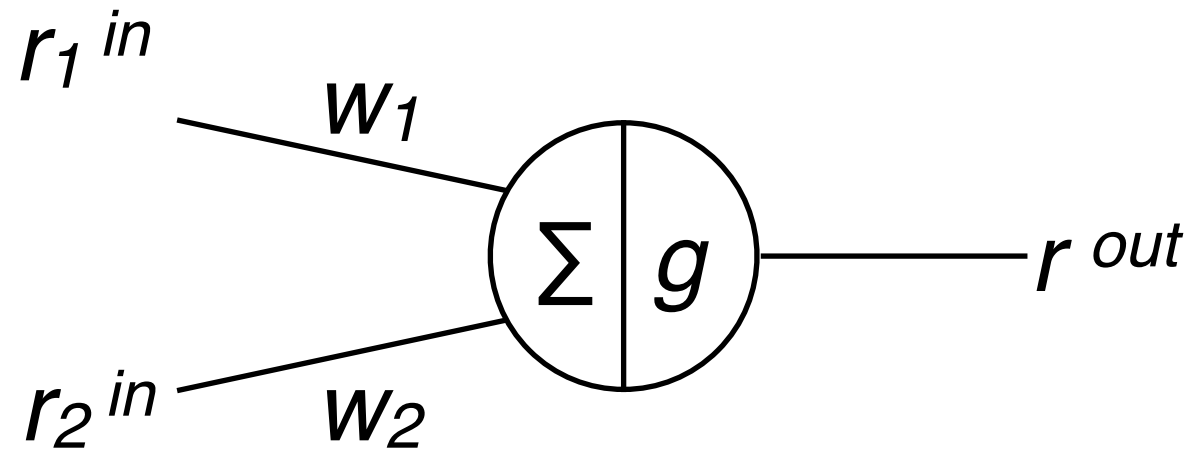


## Boolean OR

| $x_1$ | $x_2$ | $y$ |
|-------|-------|-----|
| 0     | 0     | 0   |
| 0     | 1     | 1   |
| 1     | 0     | 1   |
| 1     | 1     | 1   |

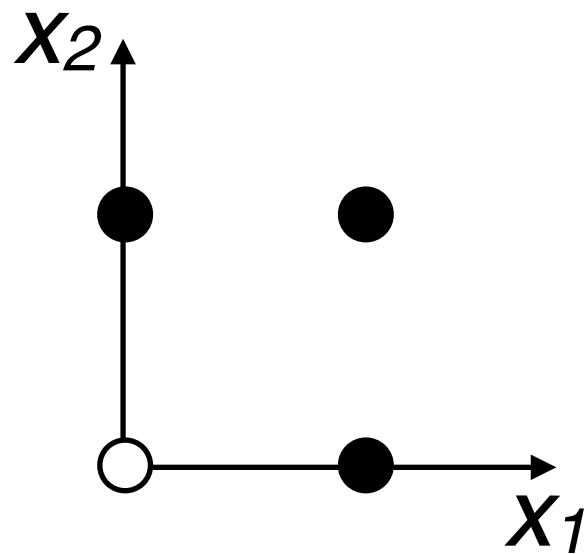


# The XOR problem

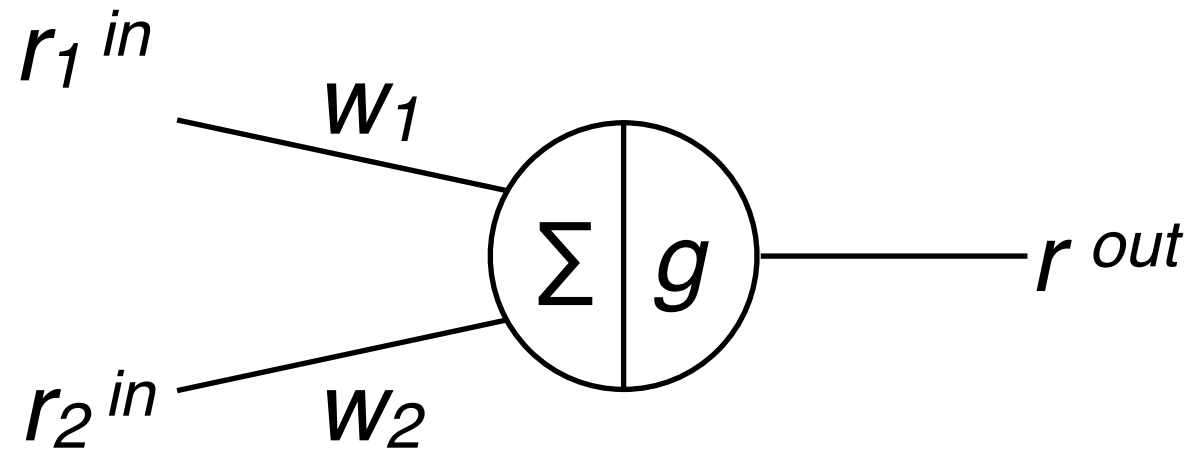


## Boolean OR

| $x_1$ | $x_2$ | $y$ |
|-------|-------|-----|
| 0     | 0     | 0   |
| 0     | 1     | 1   |
| 1     | 0     | 1   |
| 1     | 1     | 1   |

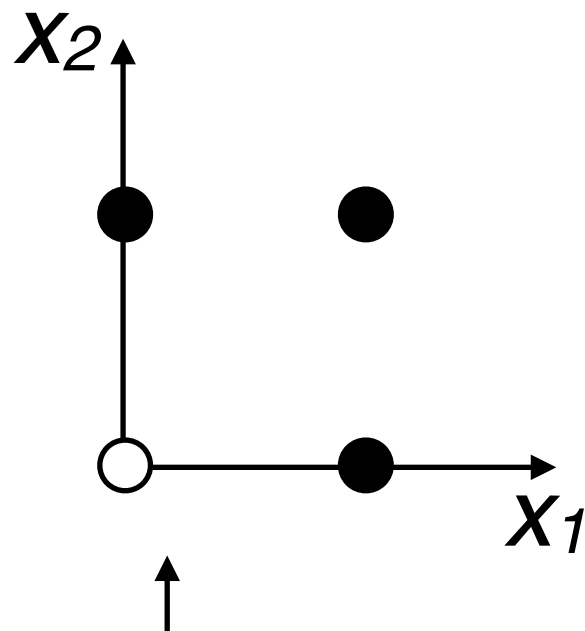


# The XOR problem



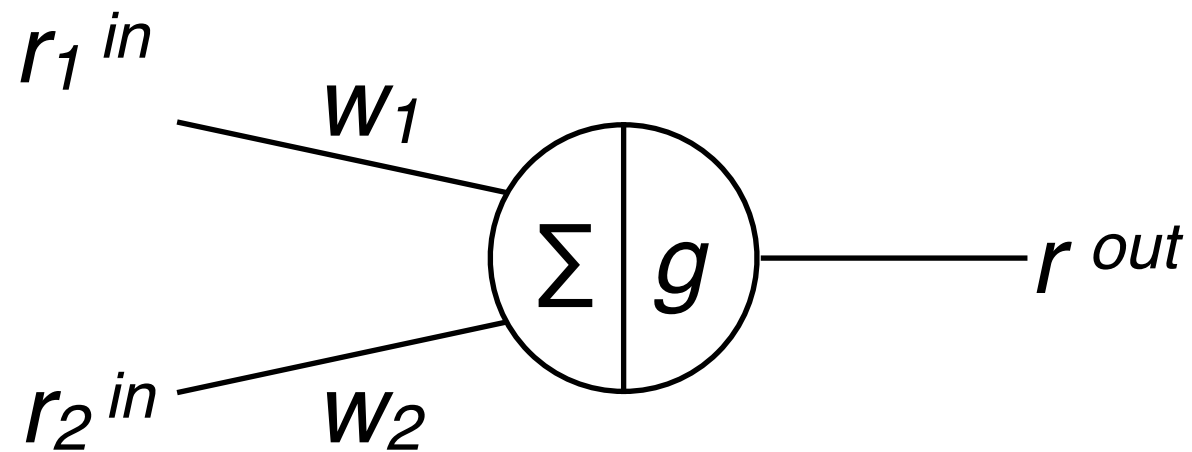
## Boolean OR

| $x_1$ | $x_2$ | $y$ |
|-------|-------|-----|
| 0     | 0     | 0   |
| 0     | 1     | 1   |
| 1     | 0     | 1   |
| 1     | 1     | 1   |



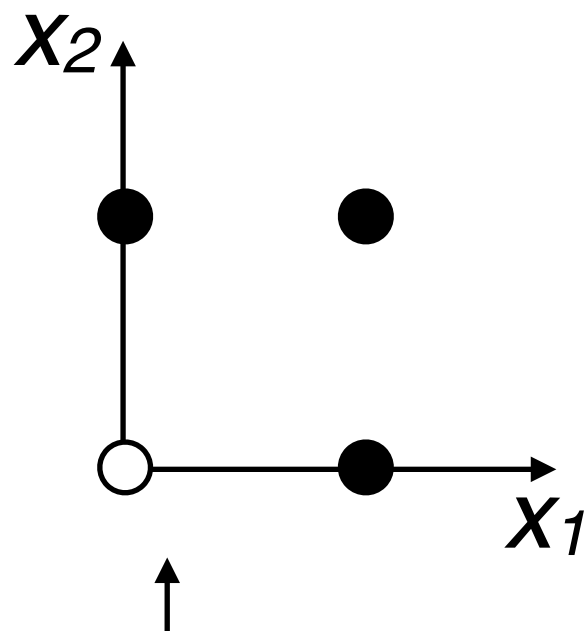
linearly separable: can be easily separated with single threshold!

# The XOR problem



## Boolean OR

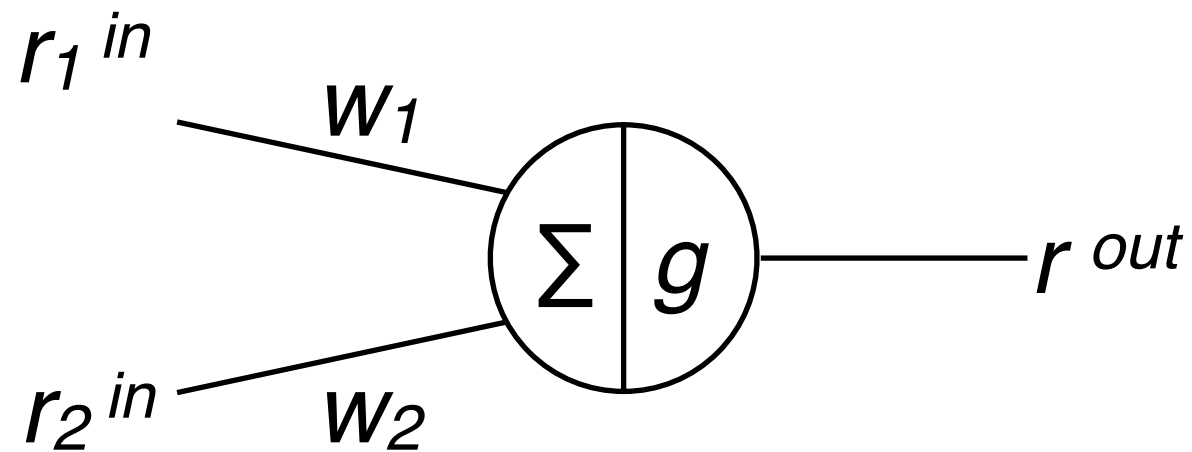
| $x_1$ | $x_2$ | $y$ |
|-------|-------|-----|
| 0     | 0     | 0   |
| 0     | 1     | 1   |
| 1     | 0     | 1   |
| 1     | 1     | 1   |



linearly separable: can be easily separated with single threshold!

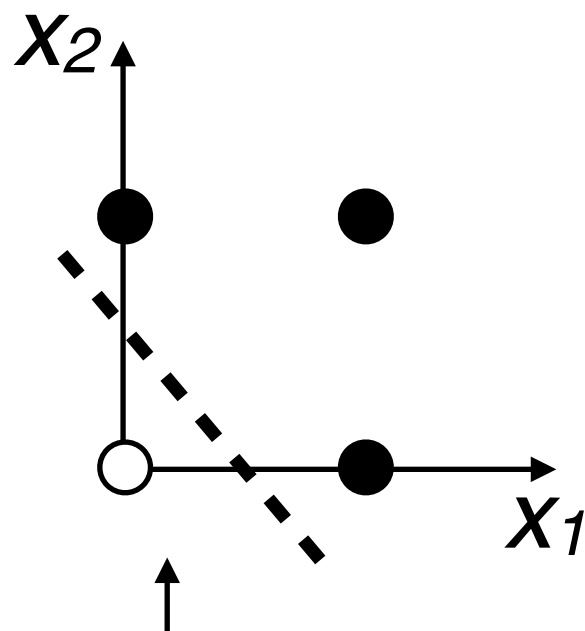
$$g(x) = \begin{cases} 1 & \text{if } x > \Theta \\ 0 & \text{else} \end{cases}$$

# The XOR problem



## Boolean OR

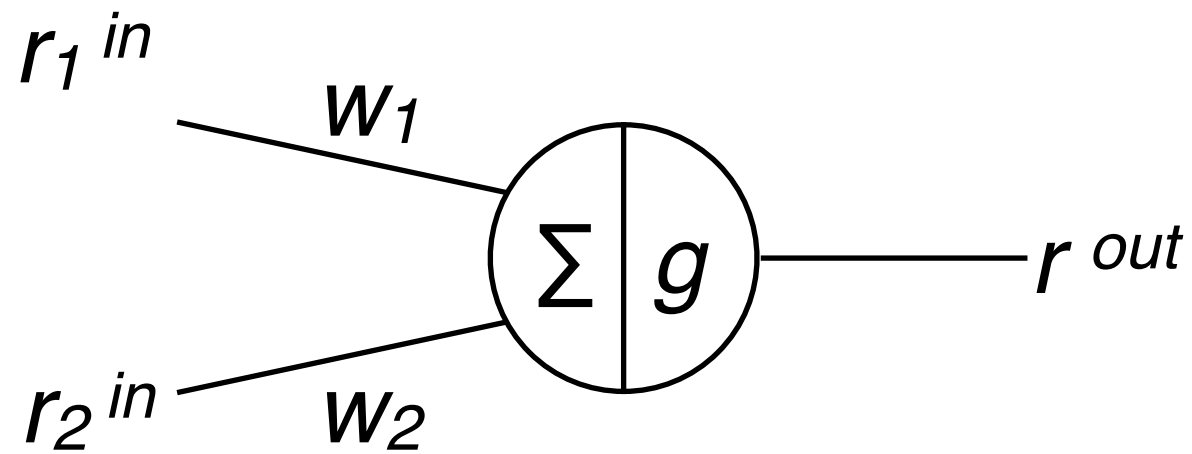
| $x_1$ | $x_2$ | $y$ |
|-------|-------|-----|
| 0     | 0     | 0   |
| 0     | 1     | 1   |
| 1     | 0     | 1   |
| 1     | 1     | 1   |



linearly separable: can be easily separated with single threshold!

$$g(x) = \begin{cases} 1 & \text{if } x > \Theta \\ 0 & \text{else} \end{cases}$$

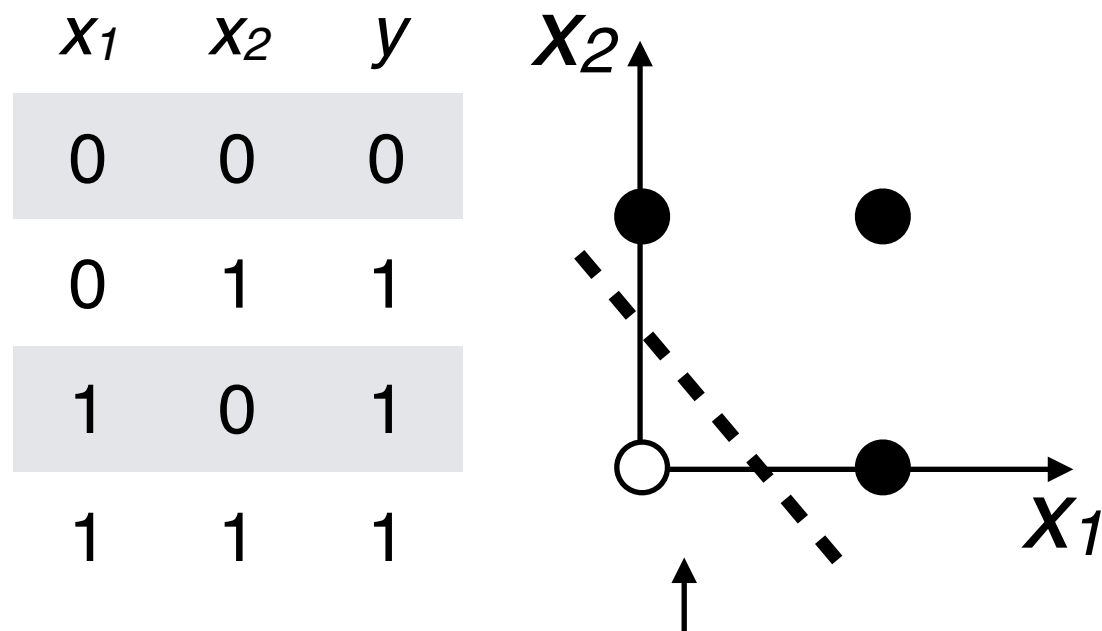
# The XOR problem



## XOR

| $x_1$ | $x_2$ | $y$ |
|-------|-------|-----|
| 0     | 0     | 0   |
| 0     | 1     | 1   |
| 1     | 0     | 1   |
| 1     | 1     | 0   |

## Boolean OR

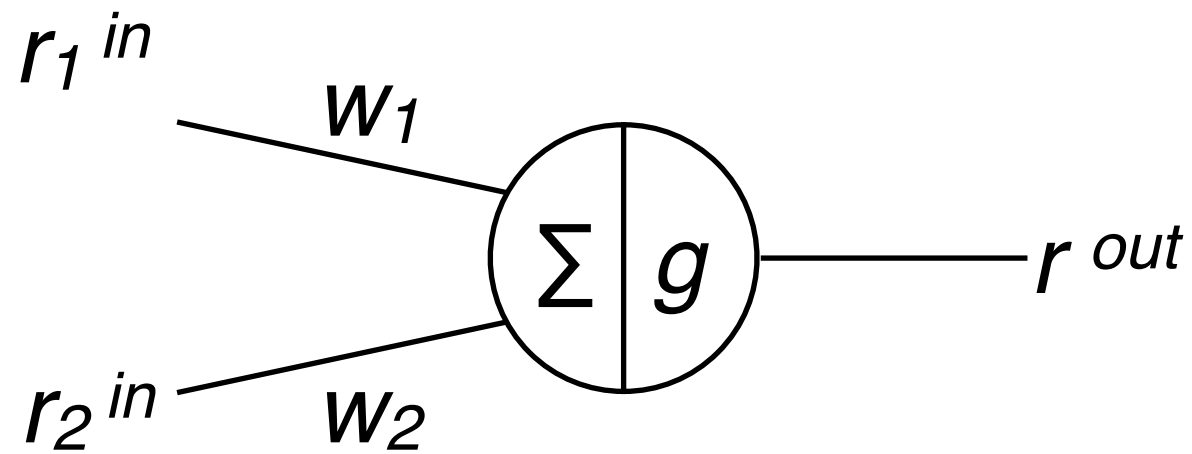


linearly separable: can be easily separated with single threshold!

$$g(x) = \begin{cases} 1 & \text{if } x > \Theta \\ 0 & \text{else} \end{cases}$$

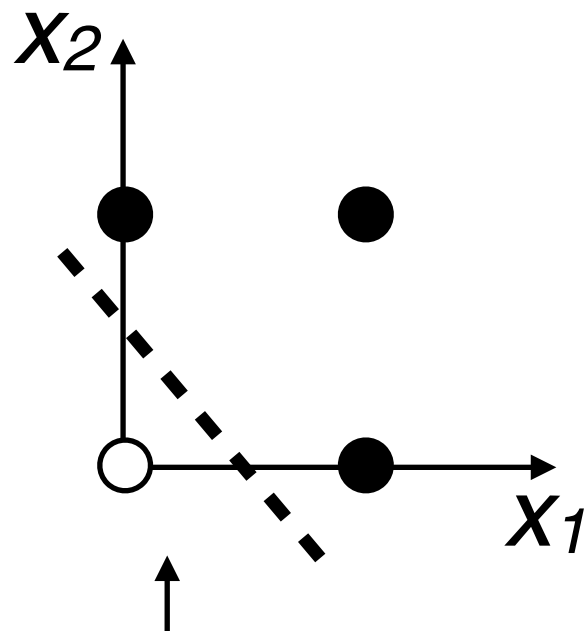


# The XOR problem



## Boolean OR

| $x_1$ | $x_2$ | $y$ |
|-------|-------|-----|
| 0     | 0     | 0   |
| 0     | 1     | 1   |
| 1     | 0     | 1   |
| 1     | 1     | 1   |

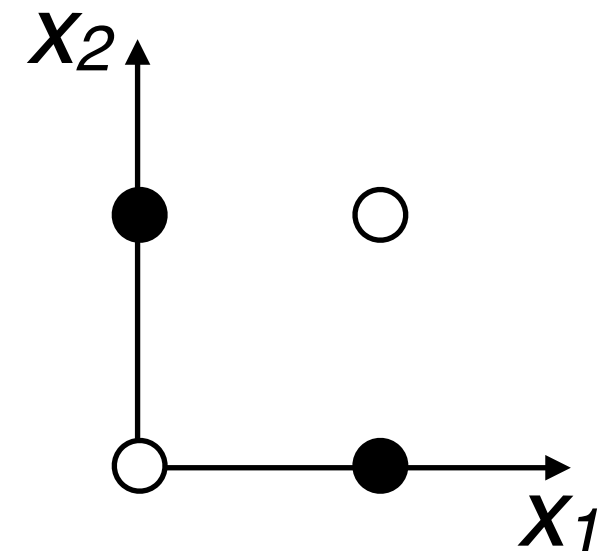


linearly separable: can be easily separated with single threshold!

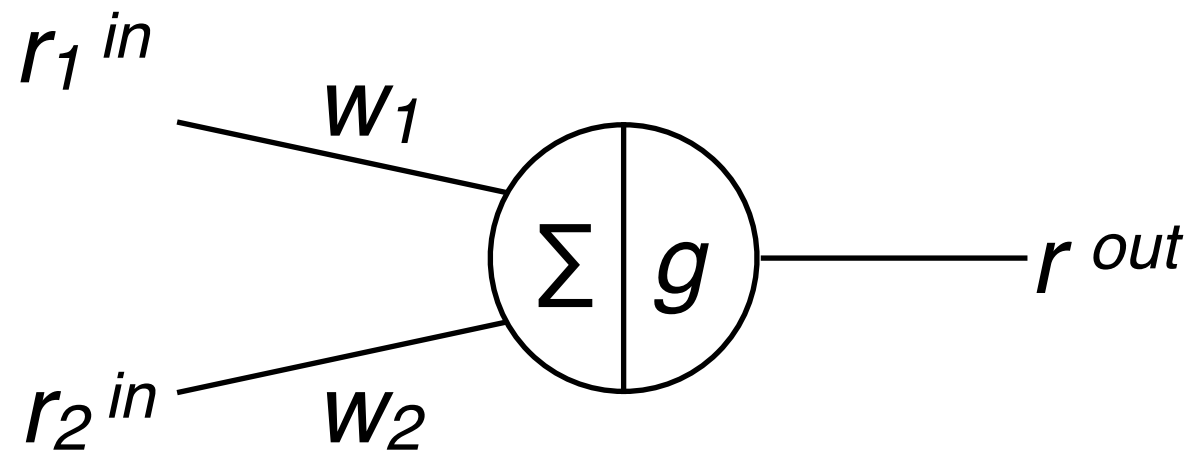
$$g(x) = \begin{cases} 1 & \text{if } x > \Theta \\ 0 & \text{else} \end{cases}$$

## XOR

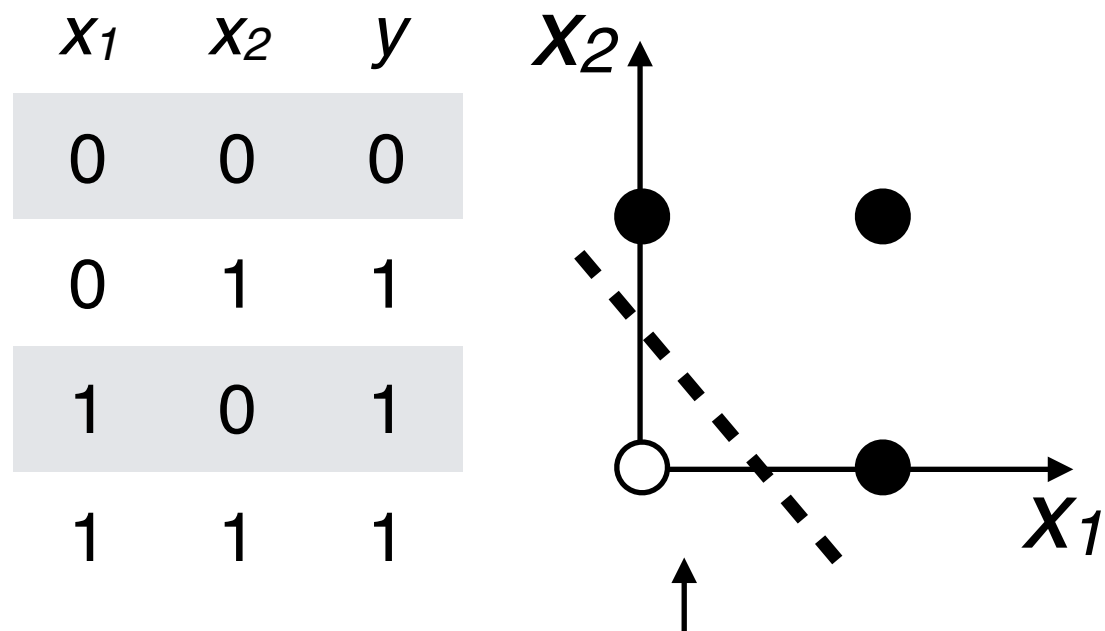
| $x_1$ | $x_2$ | $y$ |
|-------|-------|-----|
| 0     | 0     | 0   |
| 0     | 1     | 1   |
| 1     | 0     | 1   |
| 1     | 1     | 0   |



# The XOR problem



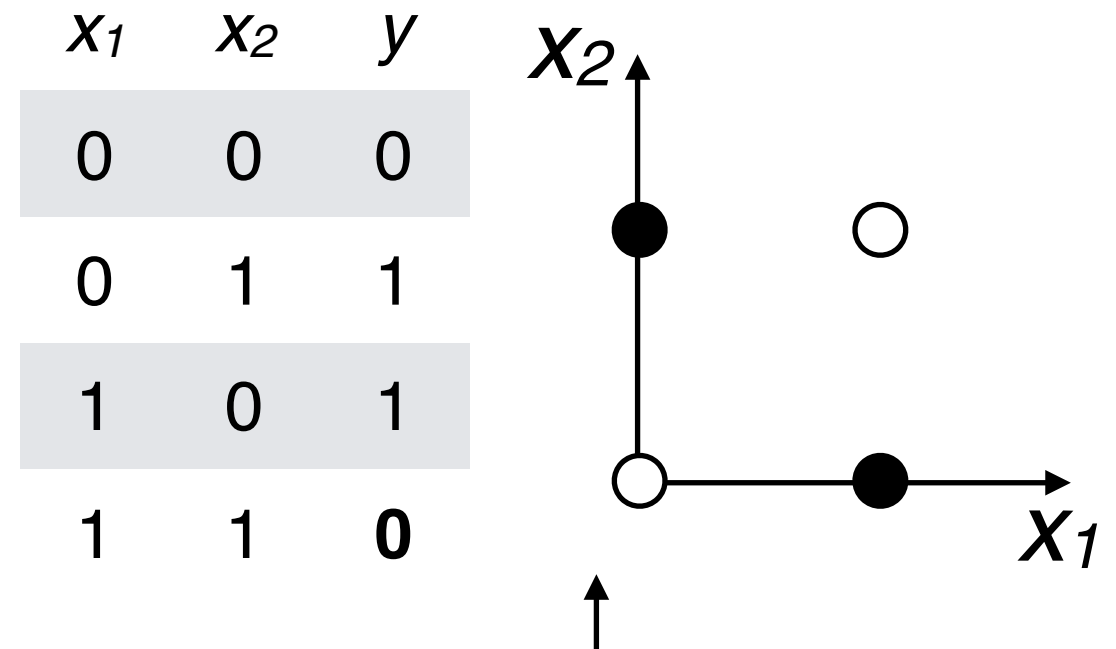
## Boolean OR



linearly separable: can be easily separated with single threshold!

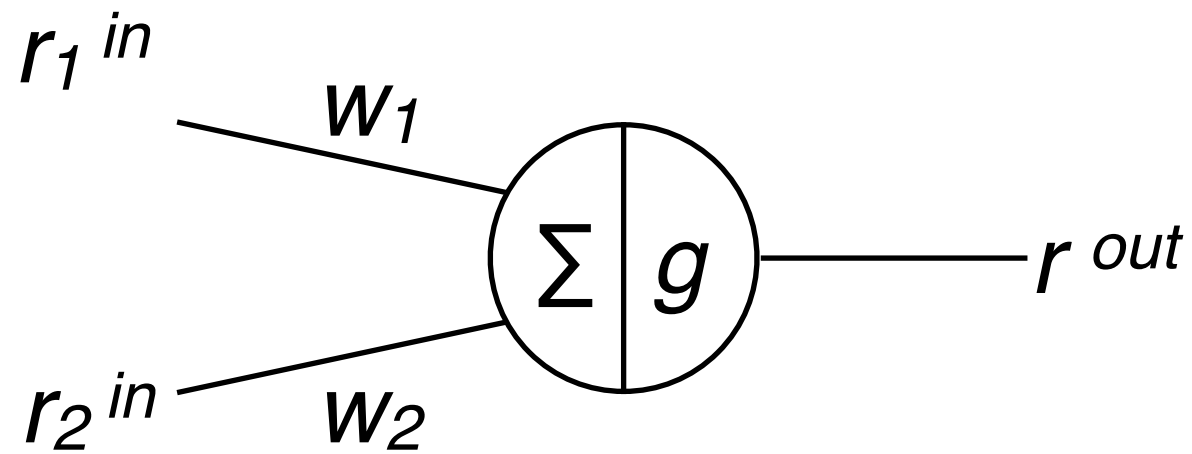
$$g(x) = \begin{cases} 1 & \text{if } x > \Theta \\ 0 & \text{else} \end{cases}$$

## XOR

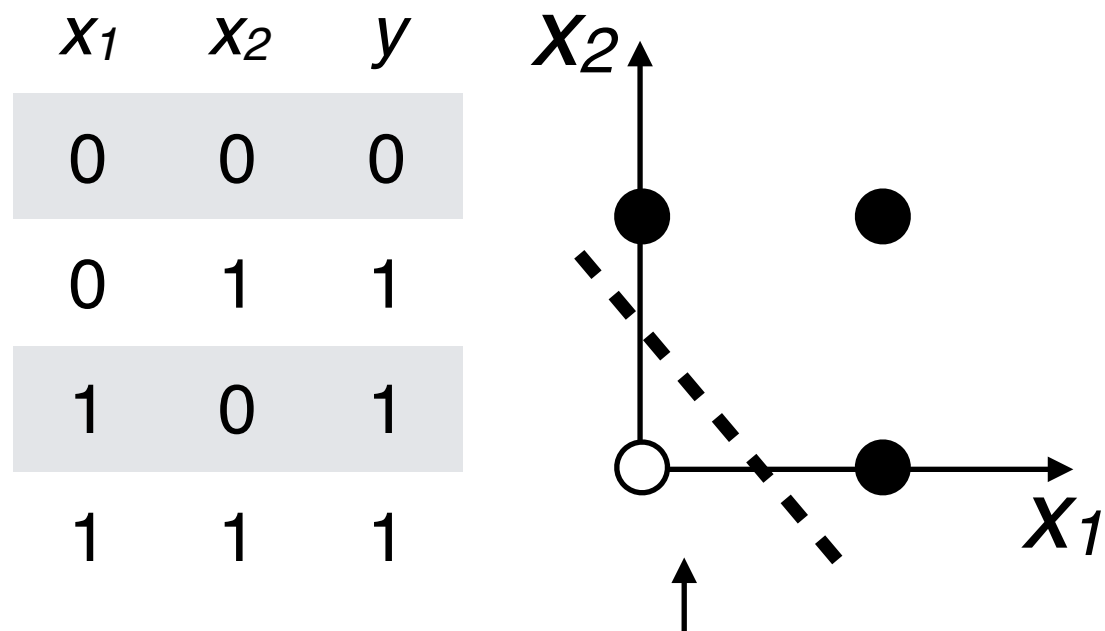


not linearly separable: more complicated!

# The XOR problem

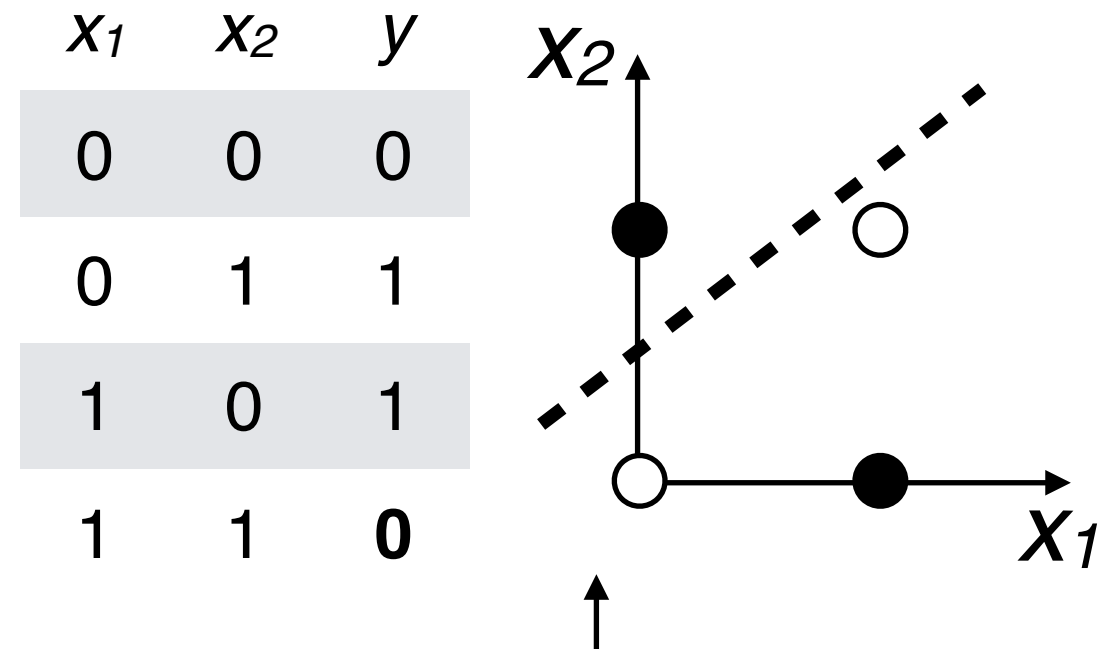


## Boolean OR



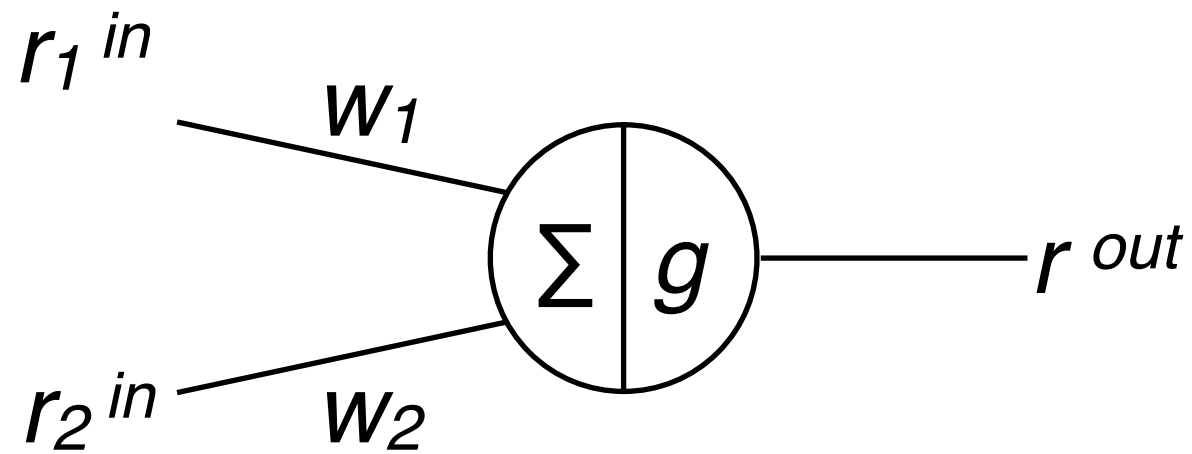
$$g(x) = \begin{cases} 1 & \text{if } x > \Theta \\ 0 & \text{else} \end{cases}$$

## XOR

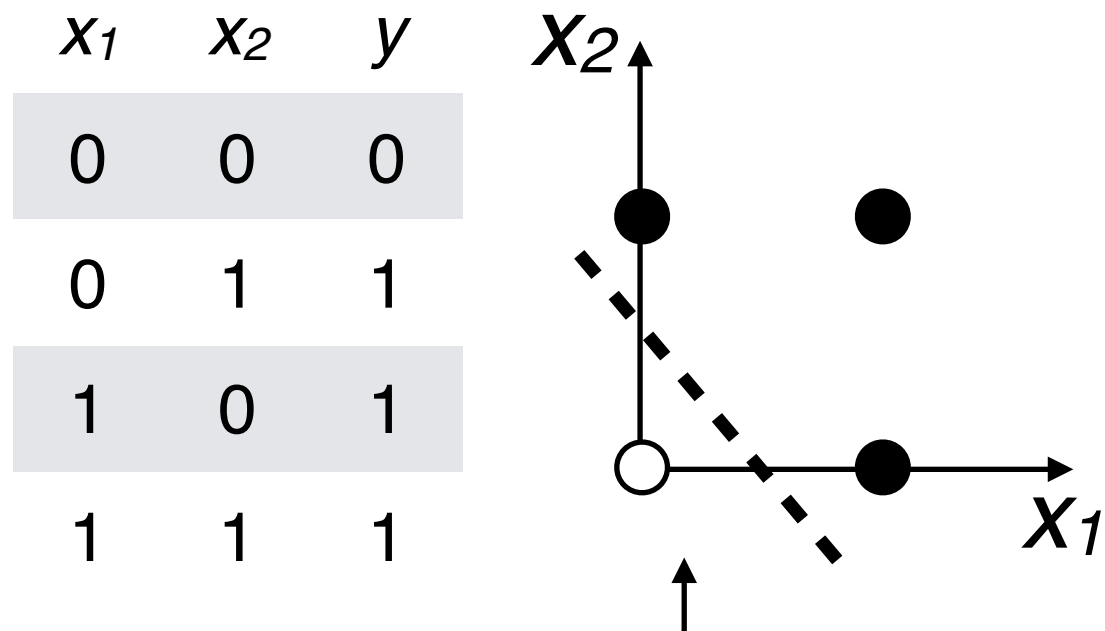


not linearly separable: more complicated!

# The XOR problem



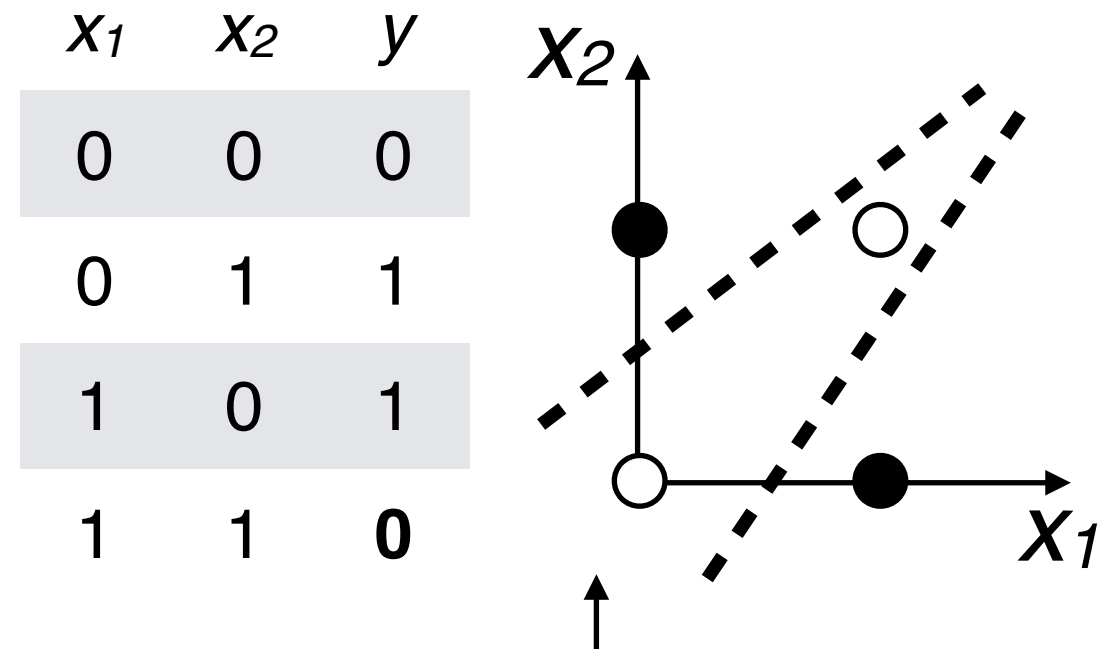
## Boolean OR



linearly separable: can be easily separated with single threshold!

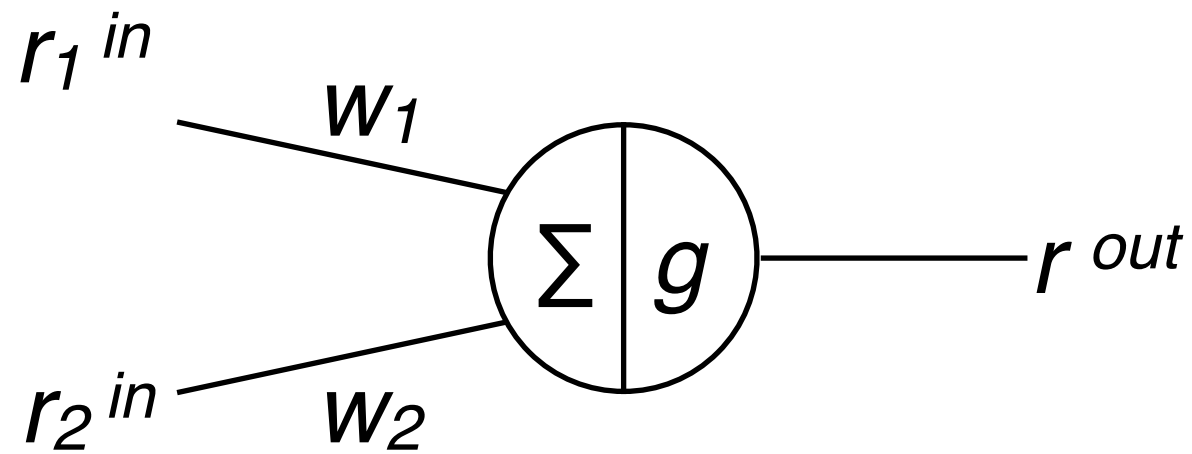
$$g(x) = \begin{cases} 1 & \text{if } x > \Theta \\ 0 & \text{else} \end{cases}$$

## XOR

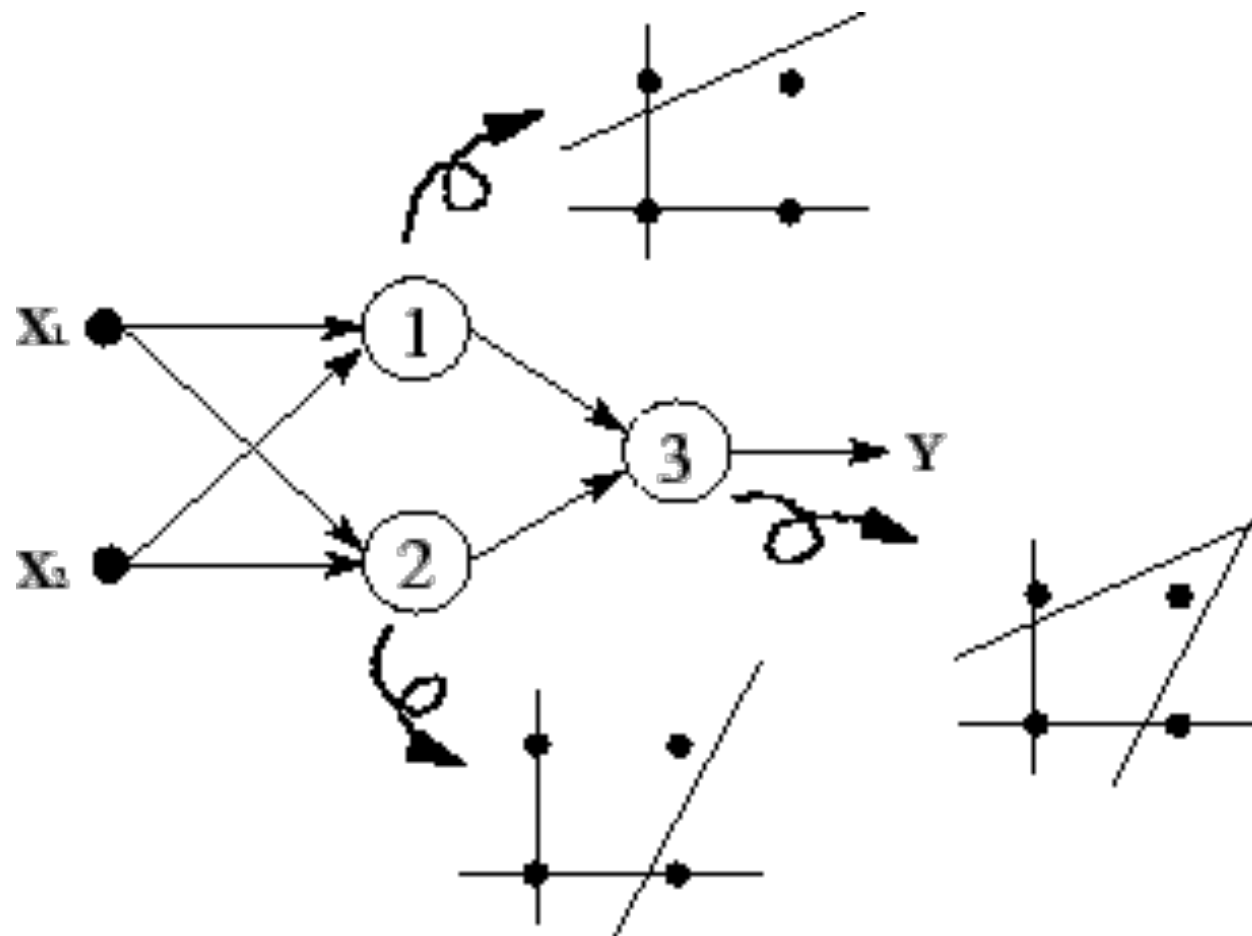


not linearly separable: more complicated!

# The XOR problem

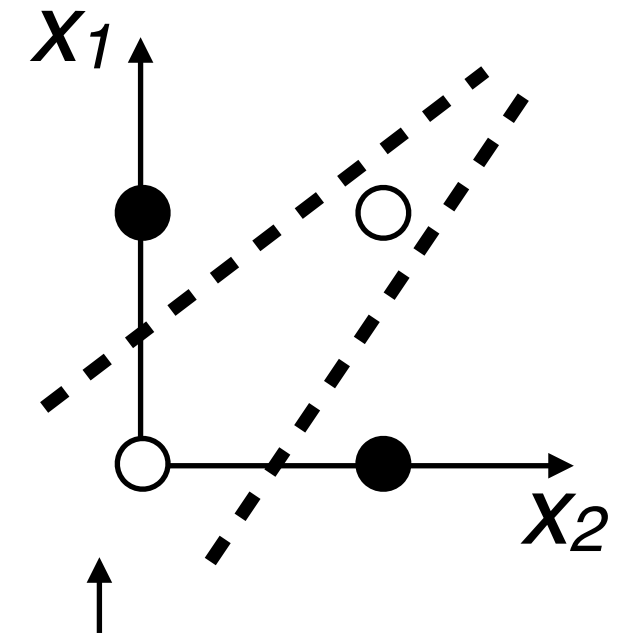


becomes



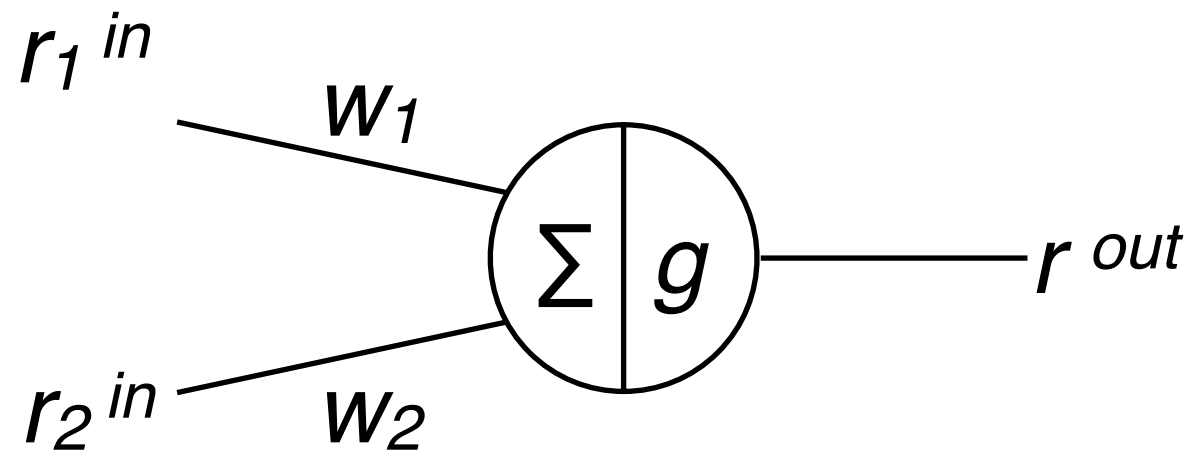
XOR

| $x_1$ | $x_2$ | $y$ |
|-------|-------|-----|
| 0     | 0     | 0   |
| 0     | 1     | 1   |
| 1     | 0     | 1   |
| 1     | 1     | 0   |

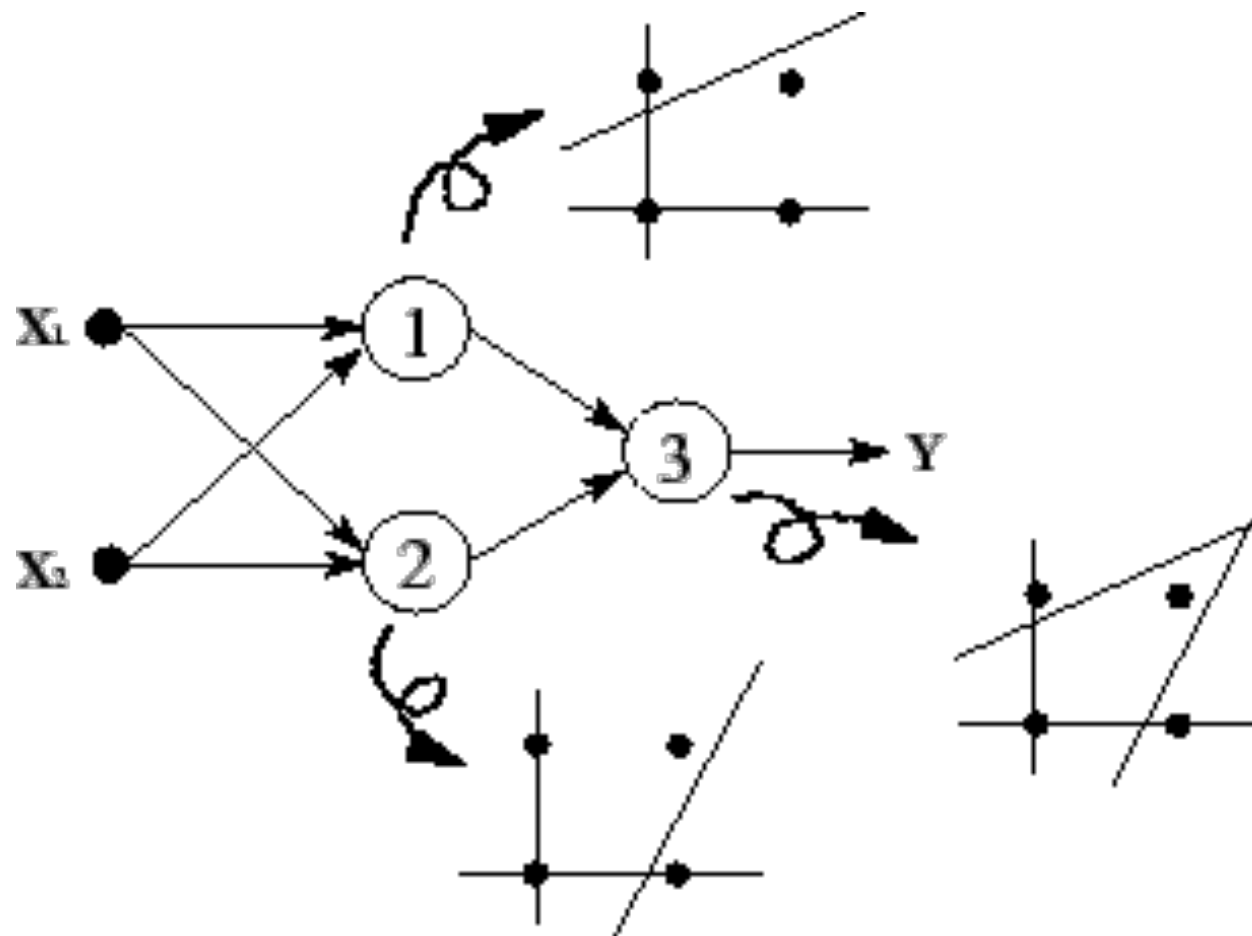


not linearly separable: more complicated!

# The XOR problem

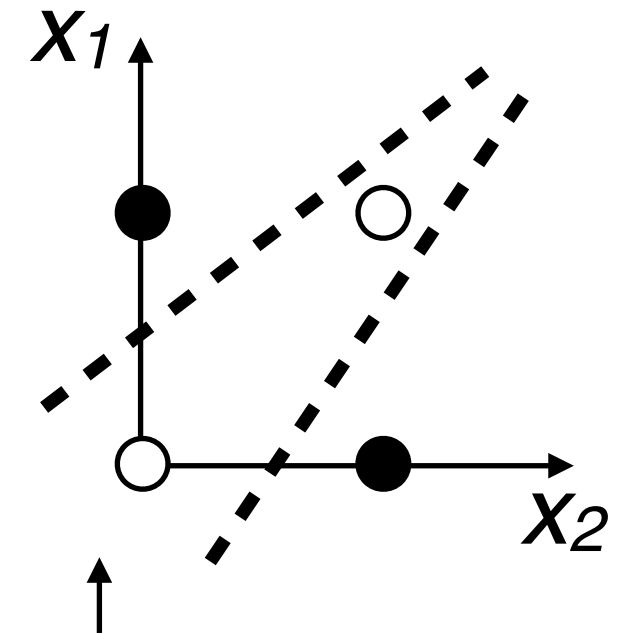


**becomes**



**XOR**

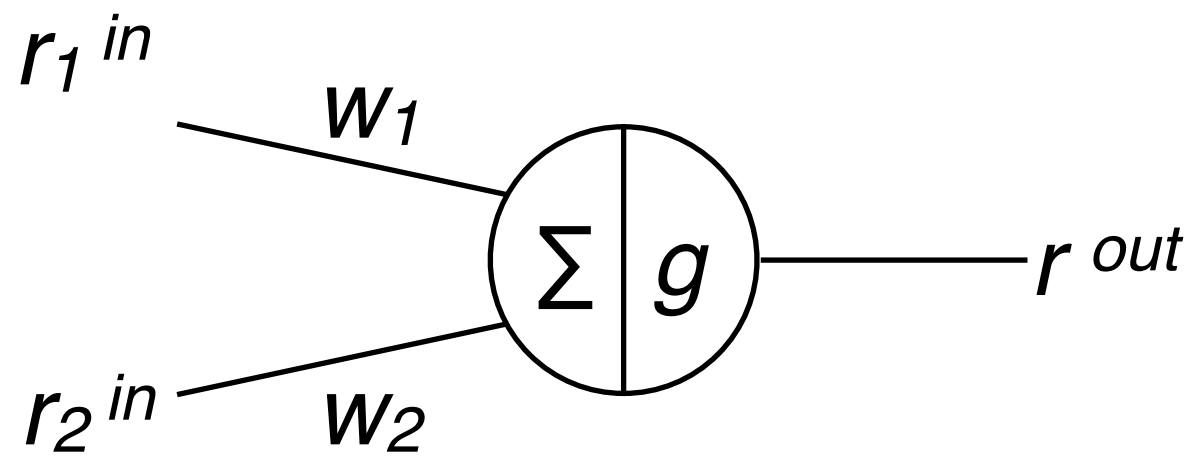
| $x_1$ | $x_2$ | $y$ |
|-------|-------|-----|
| 0     | 0     | 0   |
| 0     | 1     | 1   |
| 1     | 0     | 1   |
| 1     | 1     | 0   |



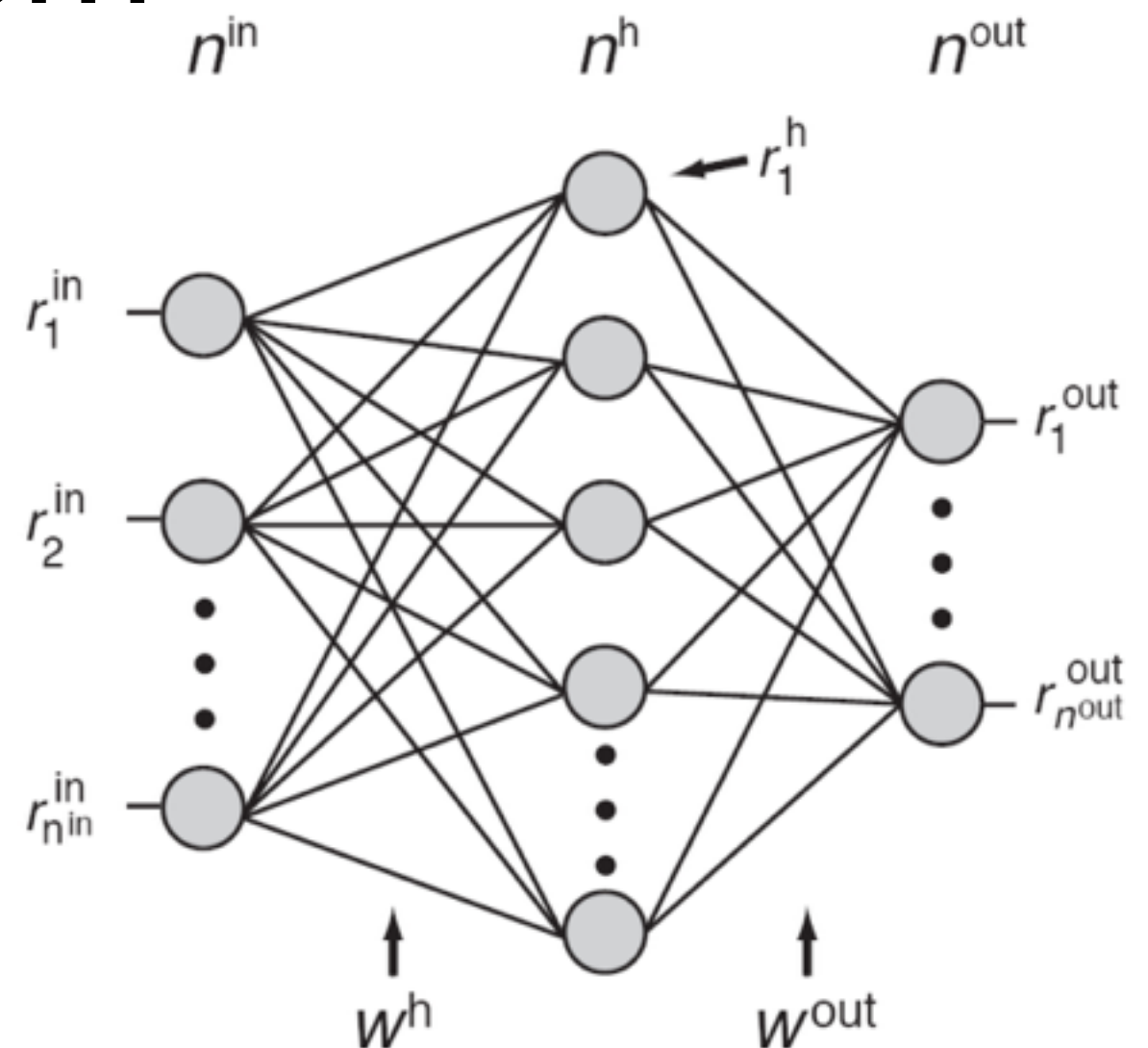
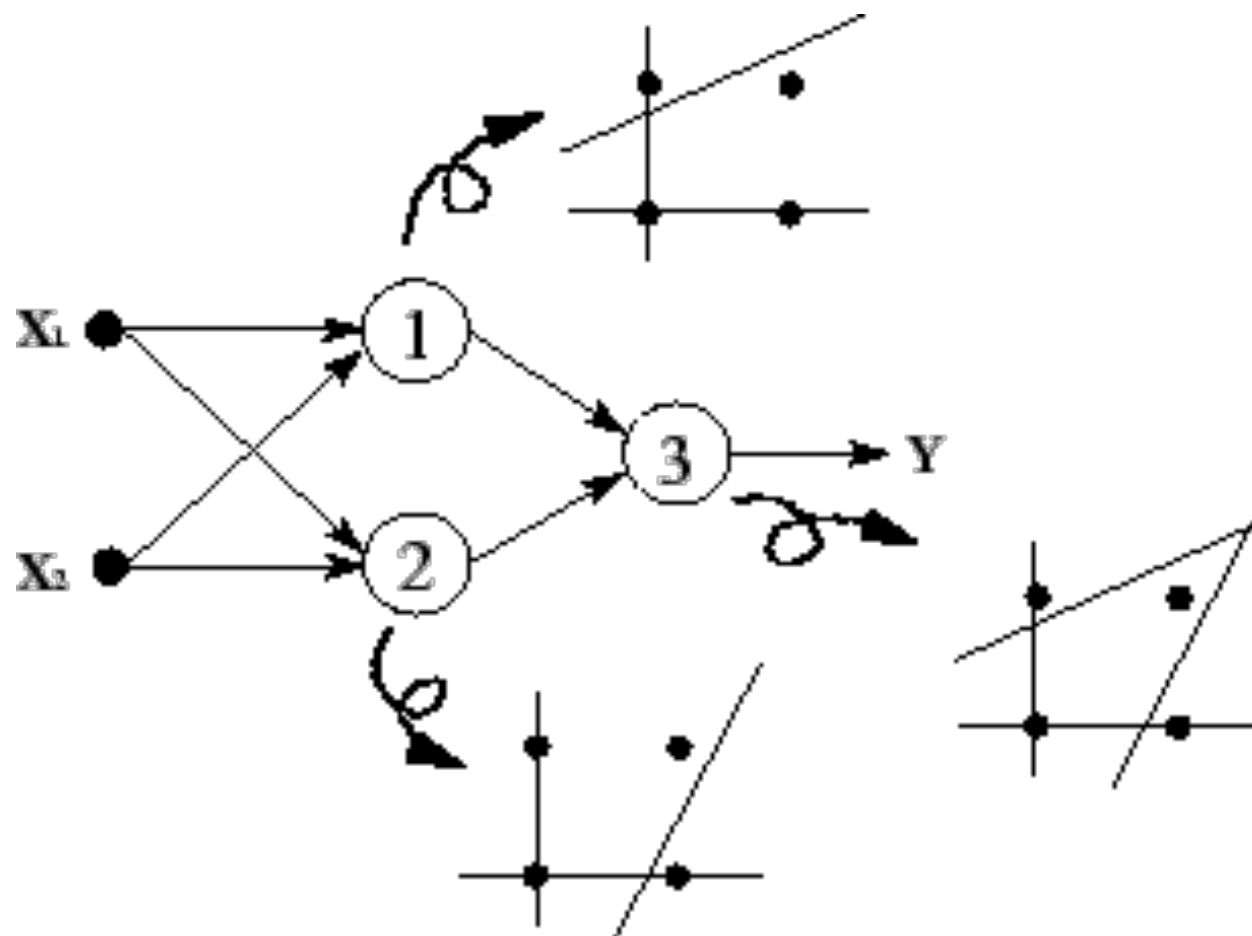
not linearly separable: more complicated!

**With enough hidden units, multilayer perceptron is the universal function approximator!**

# The XOR problem



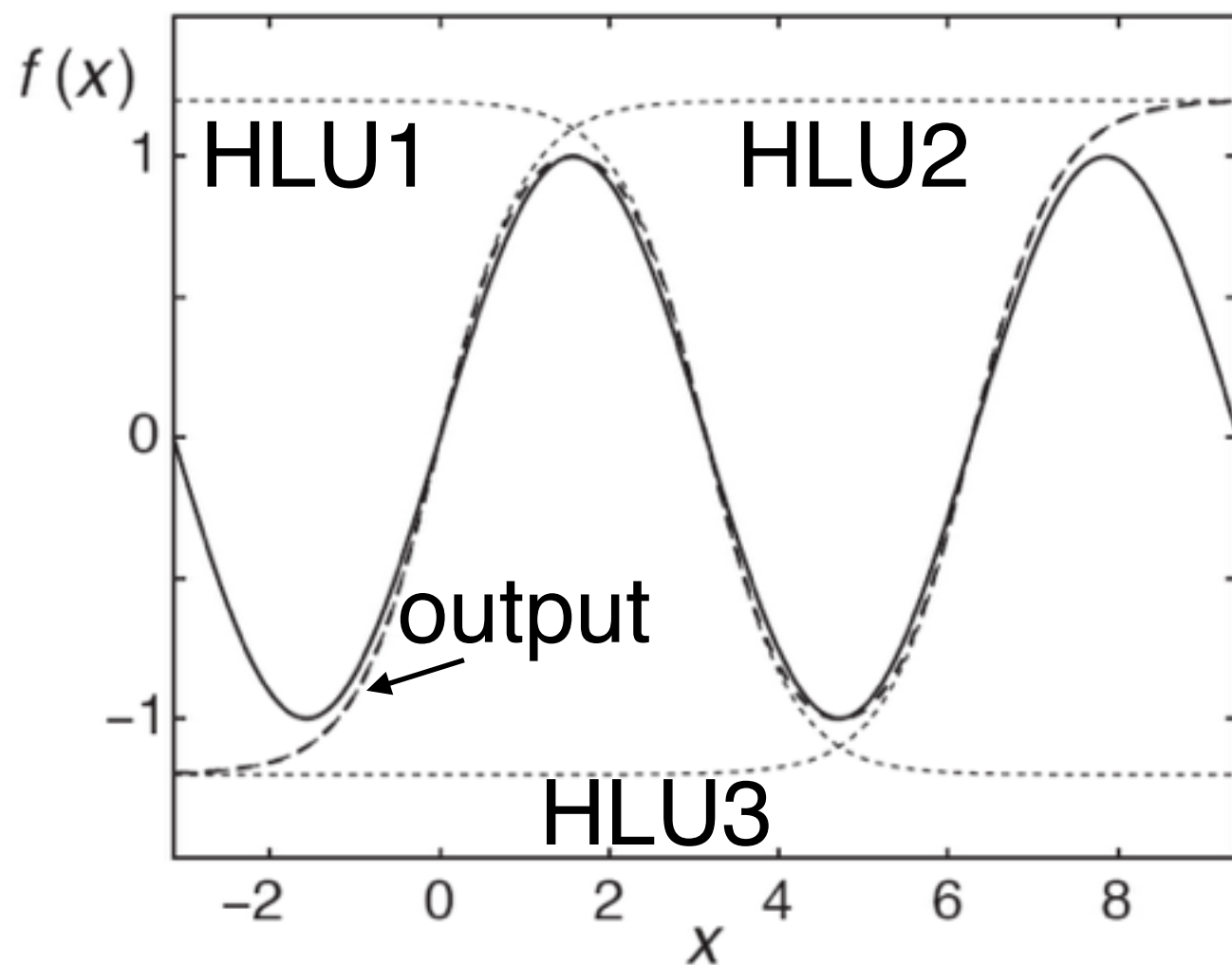
becomes



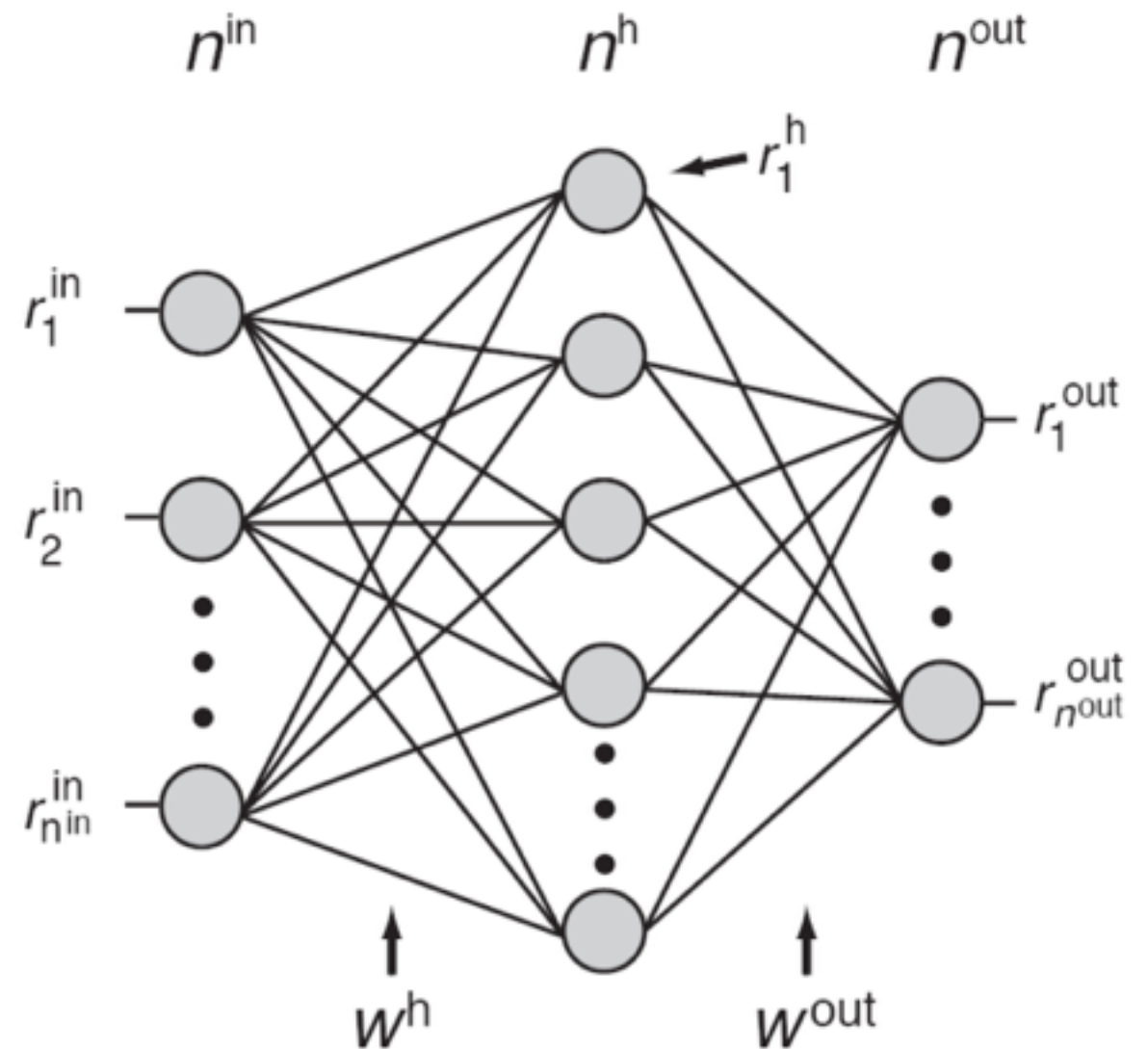
Trappenberg 2010

**With enough hidden units,  
multilayer perceptron is the  
universal function approximator!**

# Multi-layer perceptron



**E.g. sine wave approximation  
using logistic transfer function**

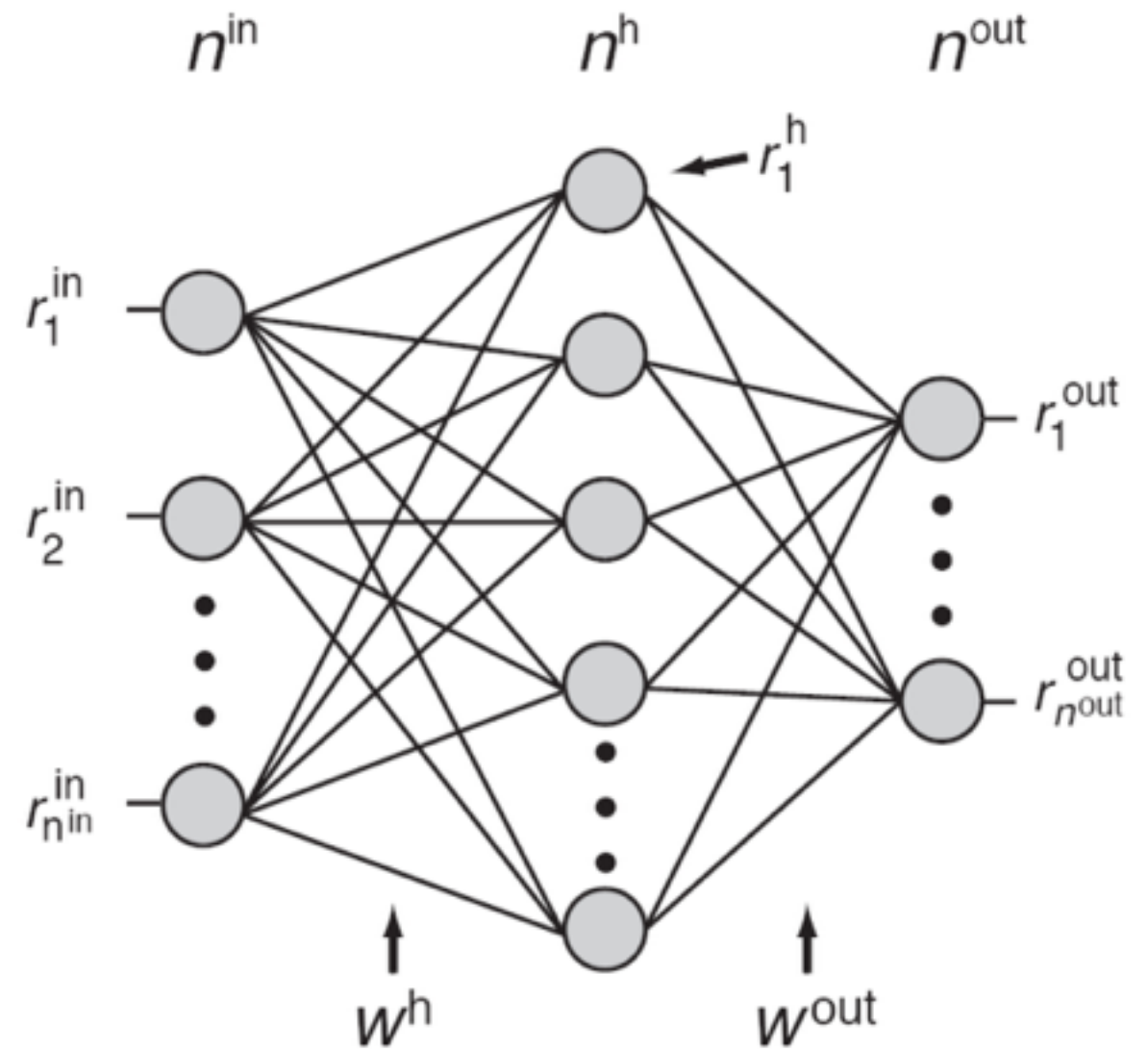


Trappenberg 2010

**With enough hidden units,  
multilayer perceptron is the  
universal function approximator!**



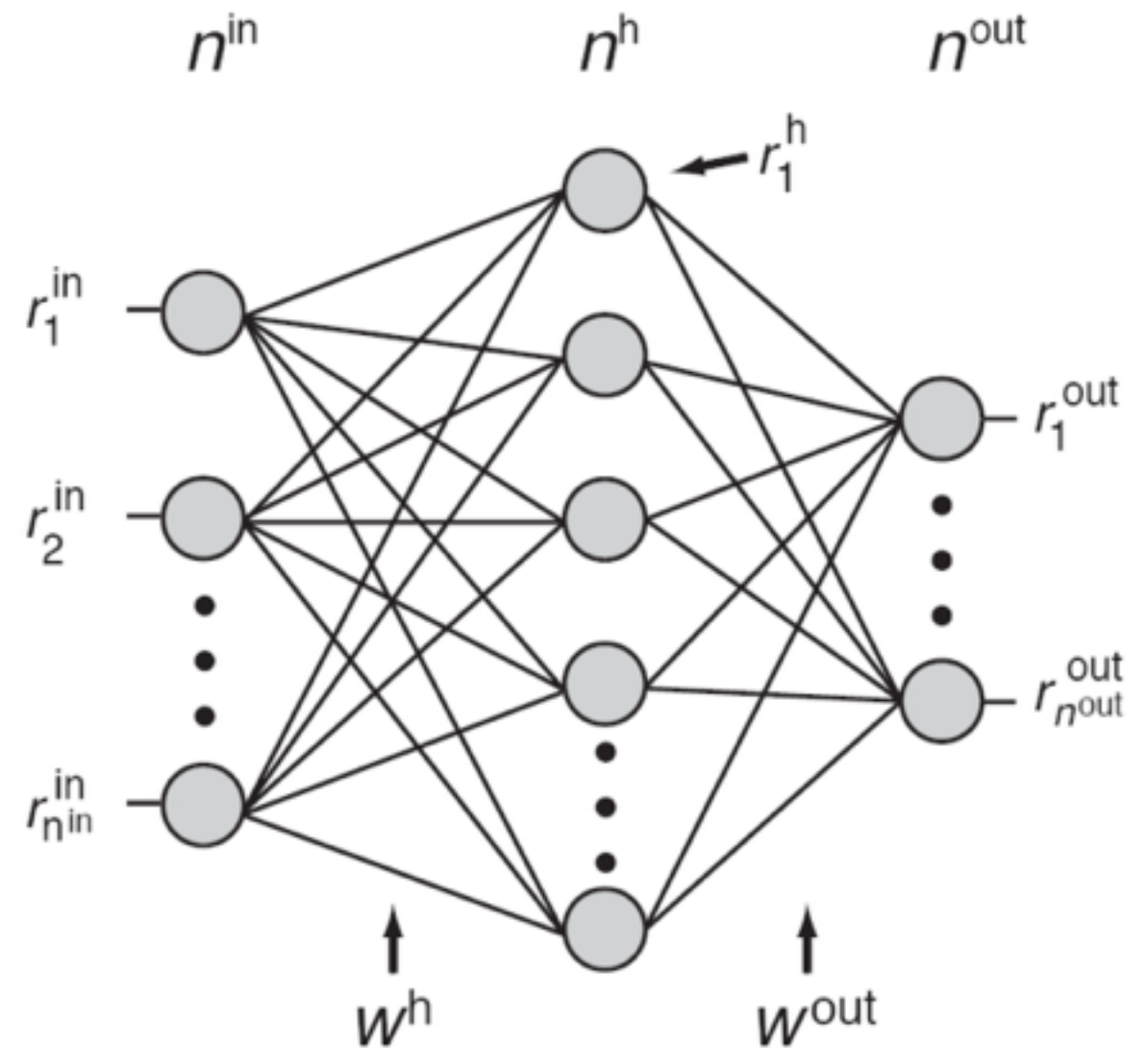
# Multi-layer perceptron



Trappenberg 2010

# Multi-layer perceptron

## Limitations

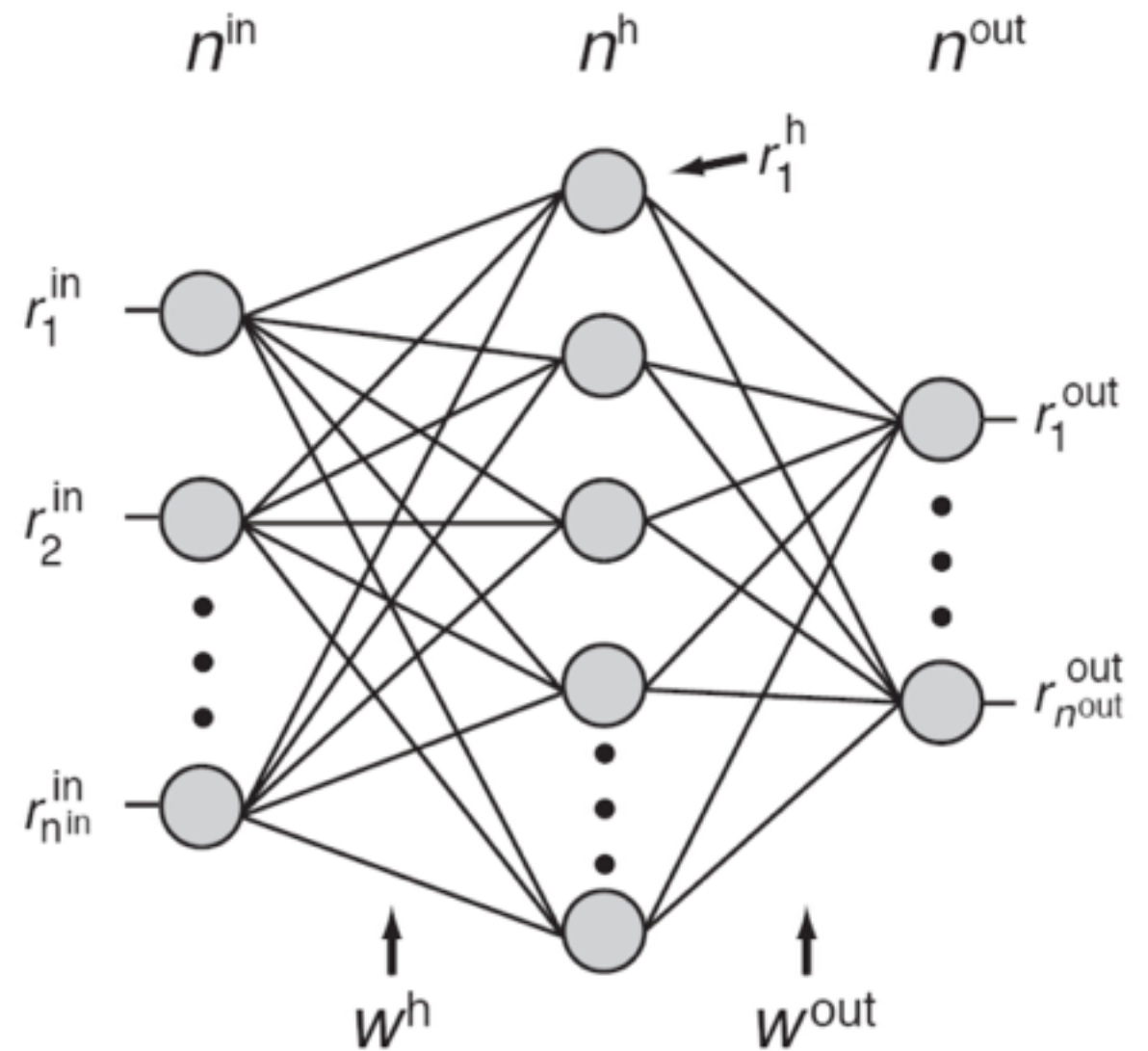


Trappenberg 2010

# Multi-layer perceptron

## Limitations

- Brain-like performance doesn't equate with actual performance

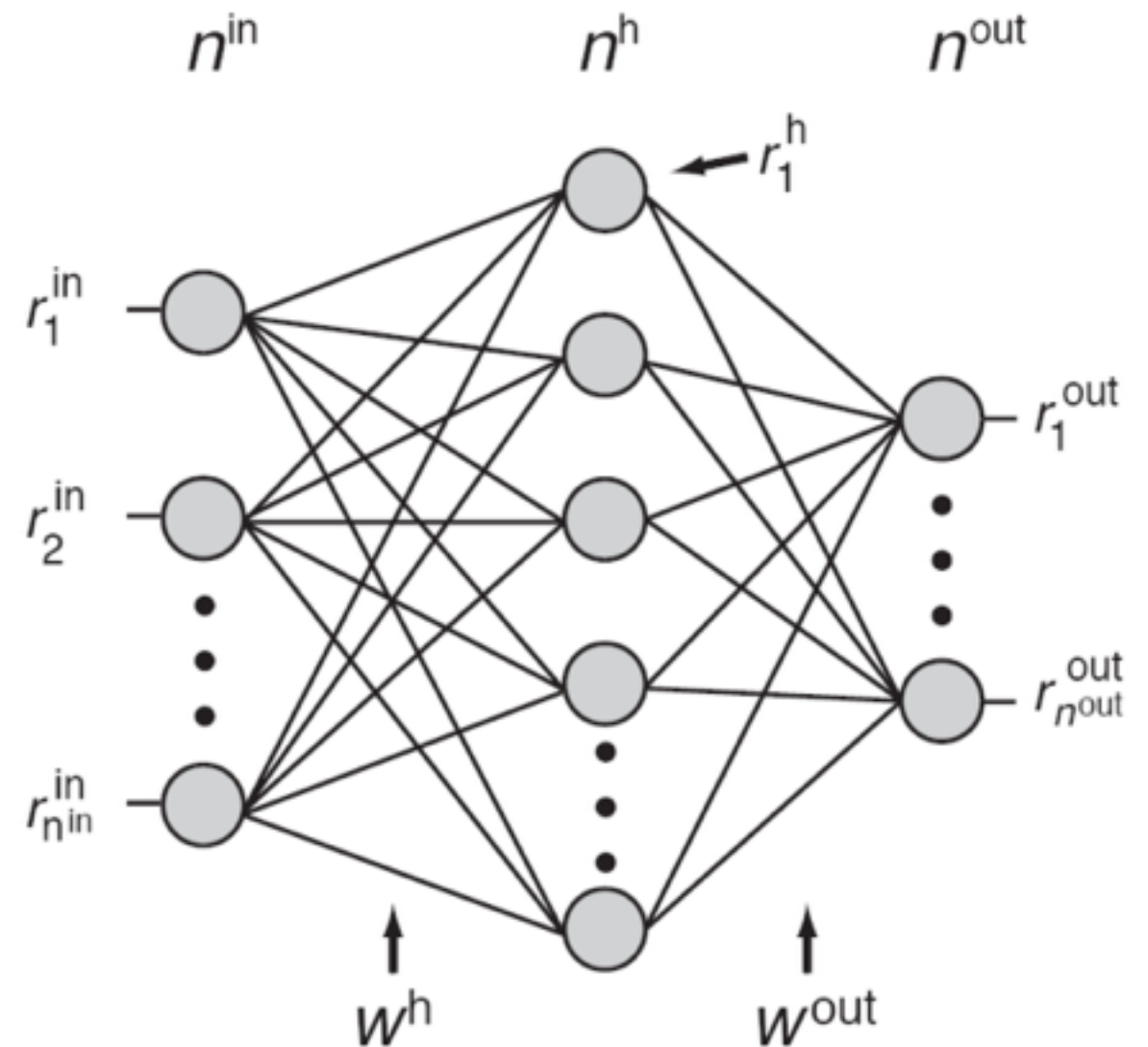


Trappenberg 2010

# Multi-layer perceptron

## Limitations

- Brain-like performance doesn't equate with actual performance
- Training rules are non-physiologic

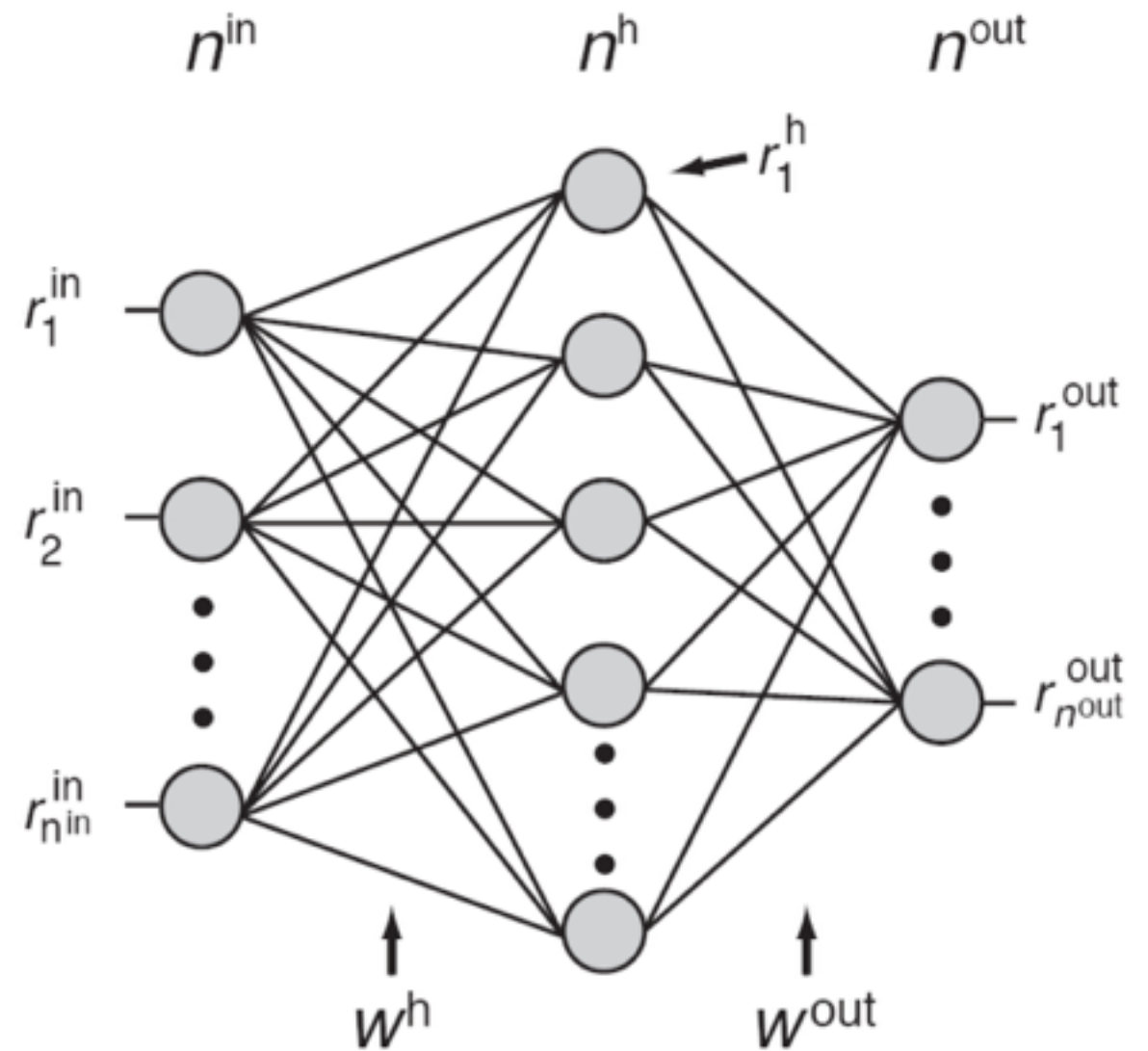


Trappenberg 2010

# Multi-layer perceptron

## Limitations

- Brain-like performance doesn't equate with actual performance
- Training rules are non-physiologic



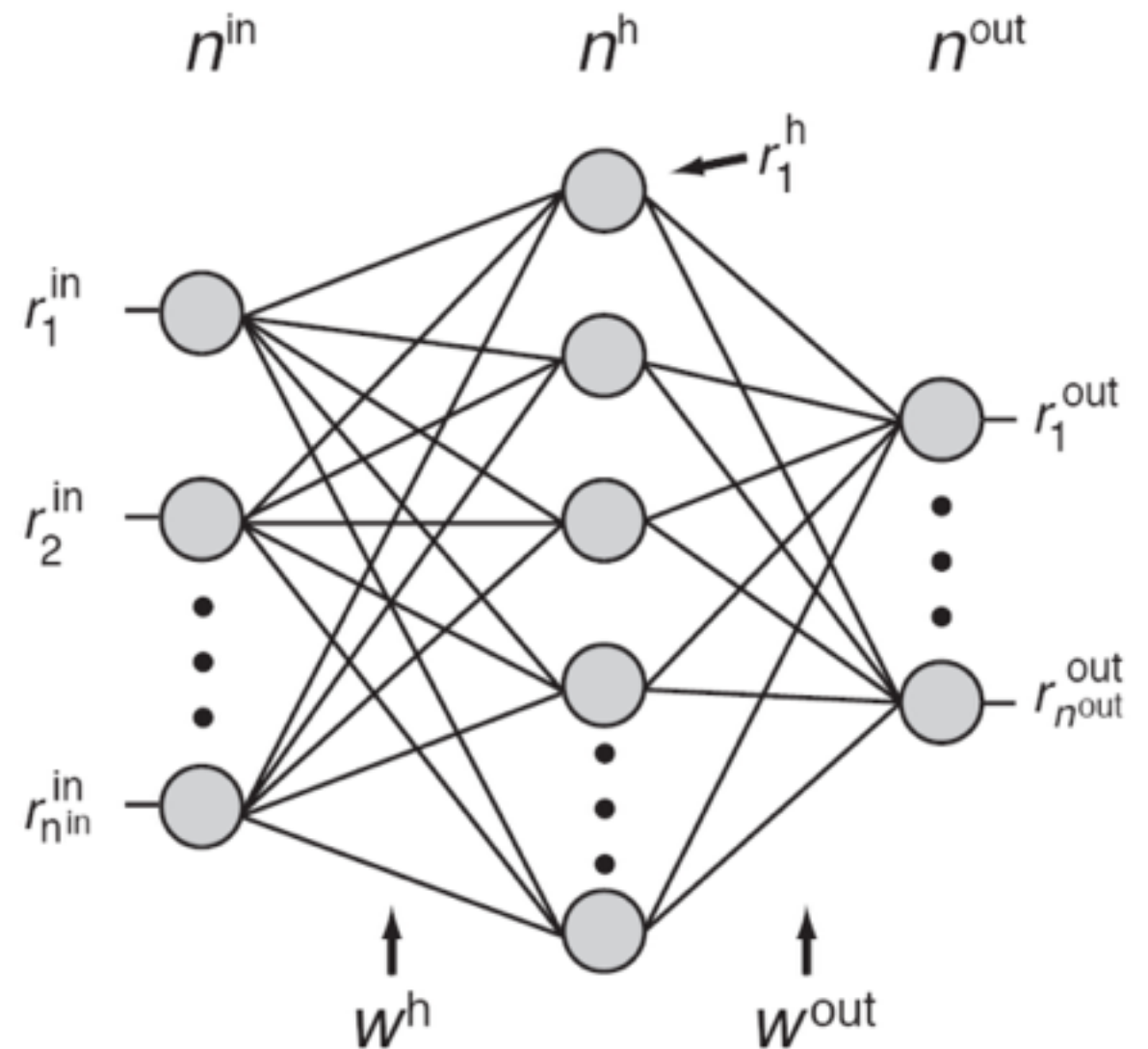
Trappenberg 2010

# Multi-layer perceptron

## Limitations

- Brain-like performance doesn't equate with actual performance
- Training rules are non-physiologic

## Strengths



Trappenberg 2010



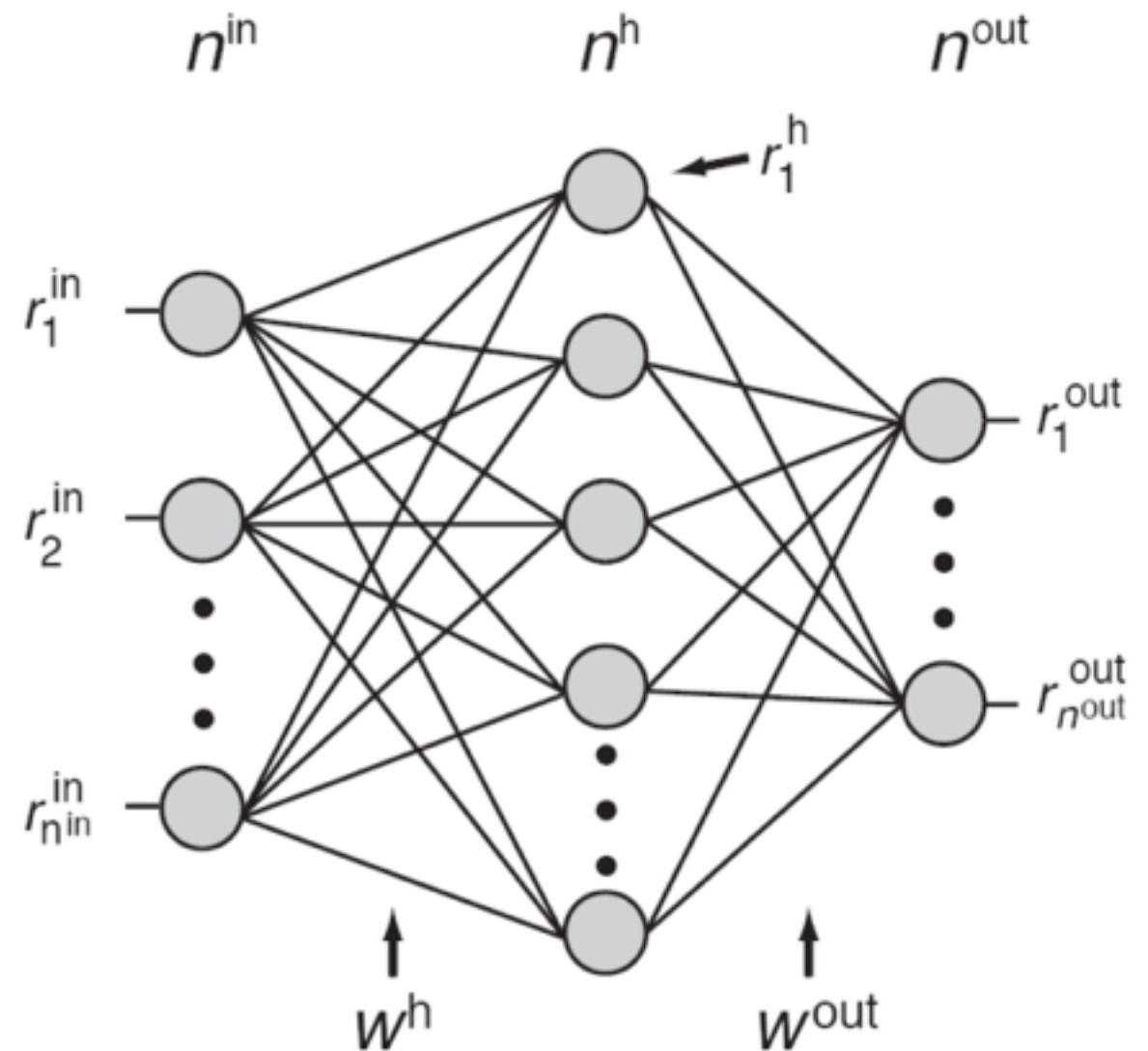
# Multi-layer perceptron

## Limitations

- Brain-like performance doesn't equate with actual performance
- Training rules are non-physiologic

## Strengths

- Hidden layer activity might resemble brain function (with appropriate inputs/outputs)



Trappenberg 2010

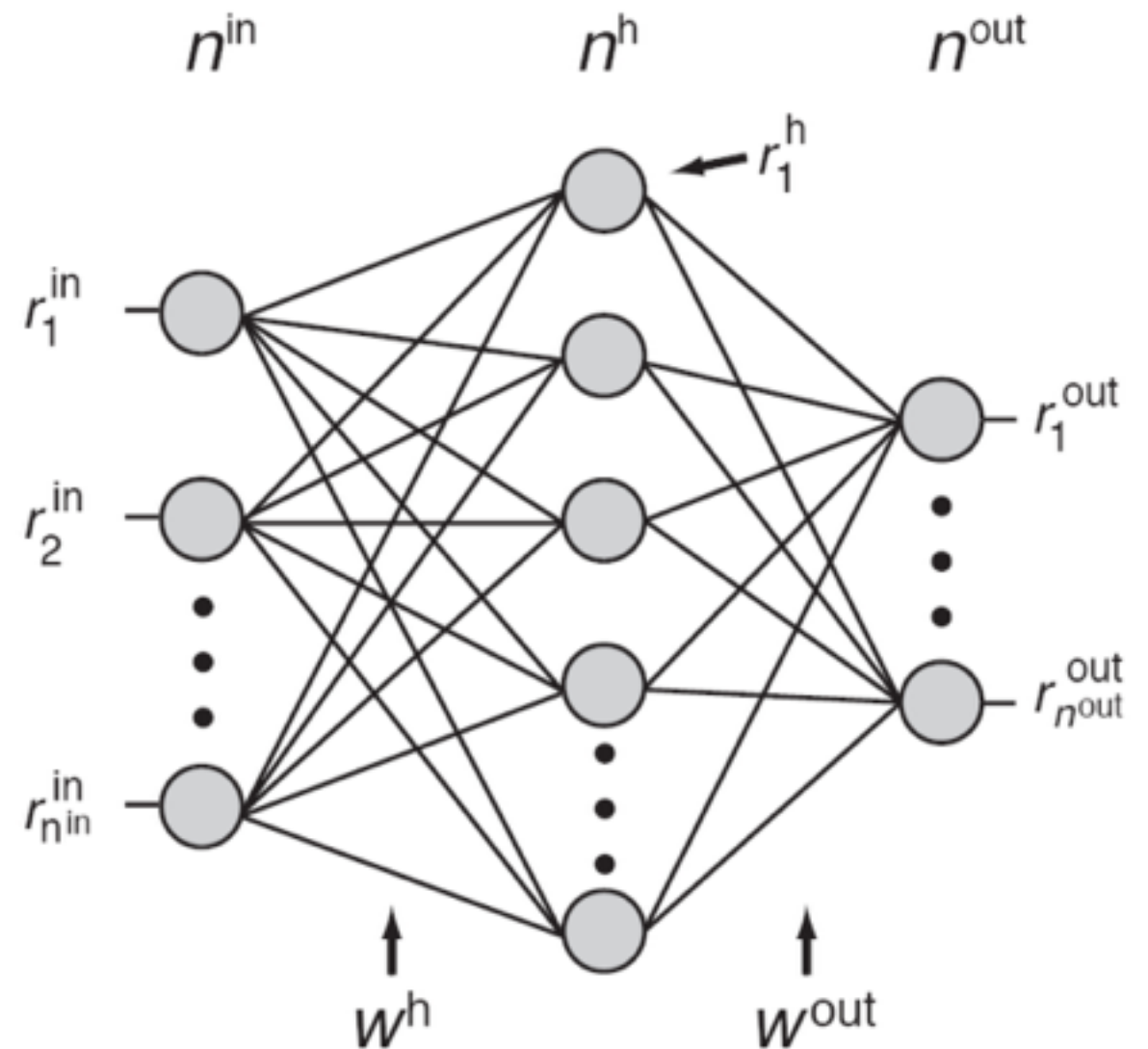
# Multi-layer perceptron

## Limitations

- Brain-like performance doesn't equate with actual performance
- Training rules are non-physiologic

## Strengths

- Hidden layer activity might resemble brain function (with appropriate inputs/outputs)
- Brain = mapping network



Trappenberg 2010



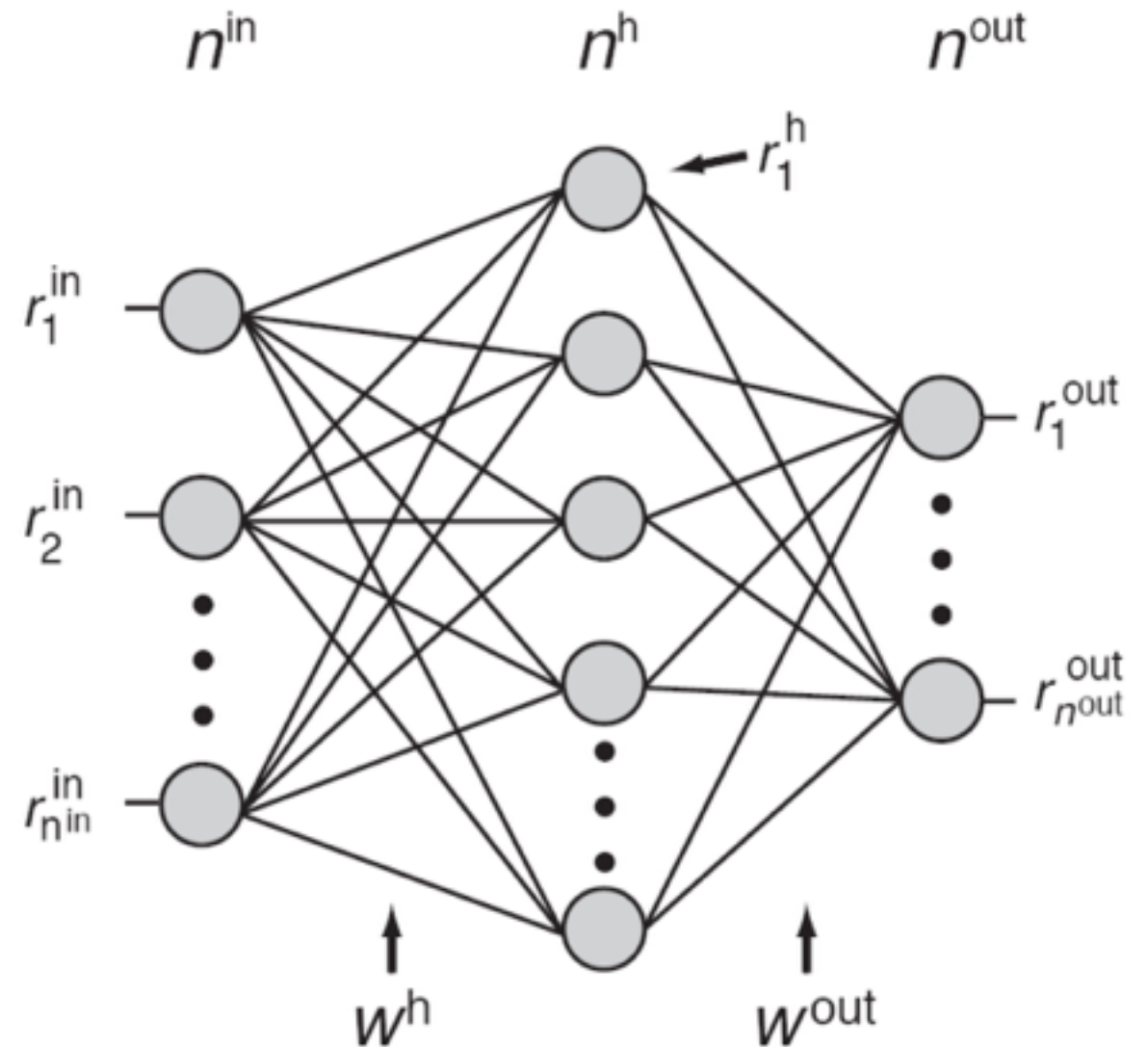
# Multi-layer perceptron

## Limitations

- Brain-like performance doesn't equate with actual performance
- Training rules are non-physiologic

## Strengths

- Hidden layer activity might resemble brain function (with appropriate inputs/outputs)
- Brain = mapping network
- Self-Organization, like the brain



Trappenberg 2010

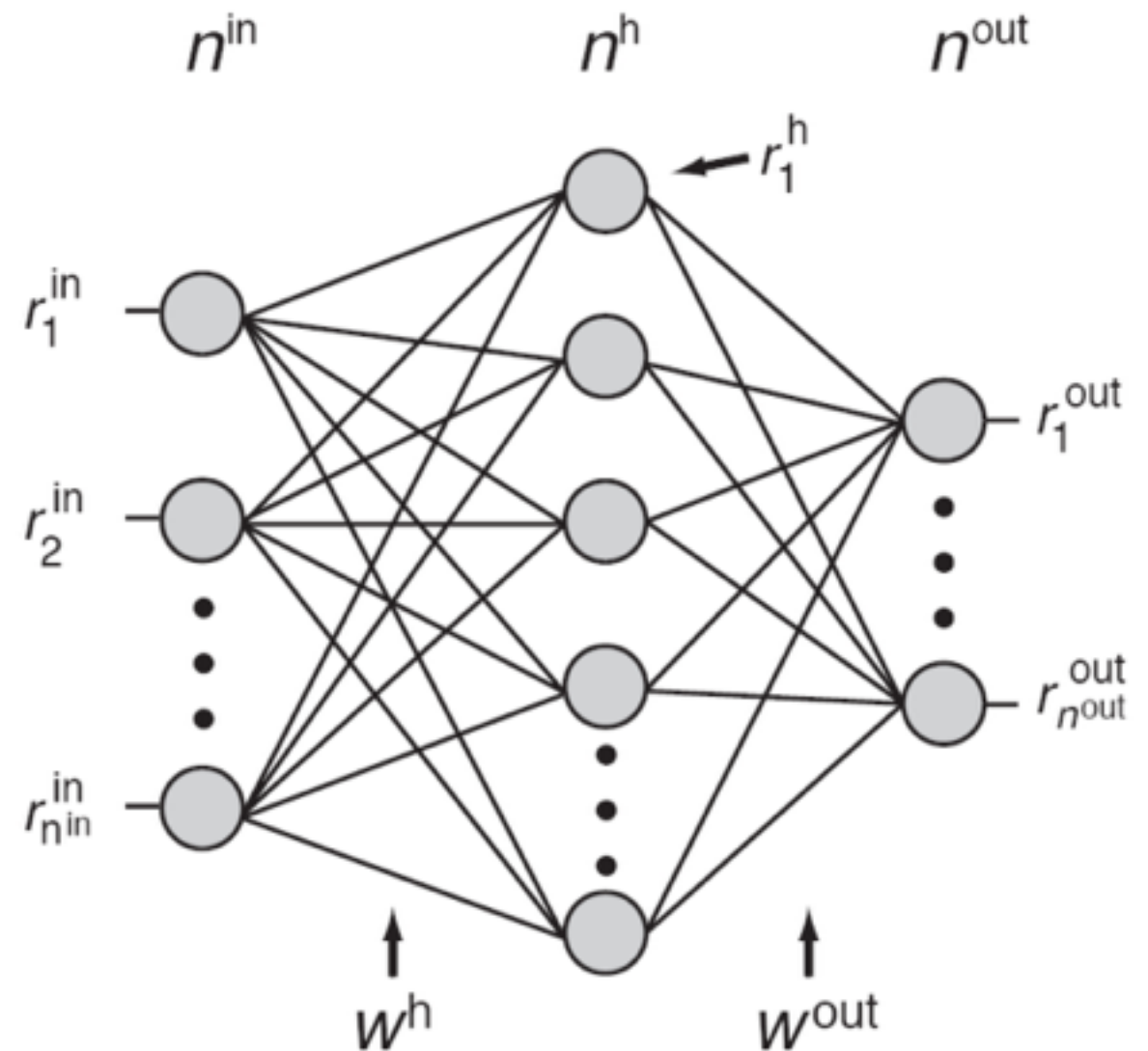
# Multi-layer perceptron

## Limitations

- Brain-like performance doesn't equate with actual performance
- Training rules are non-physiologic

## Strengths

- Hidden layer activity might resemble brain function (with appropriate inputs/outputs)
- Brain = mapping network
- Self-Organization, like the brain
- High flexibility in possible computations



Trappenberg 2010

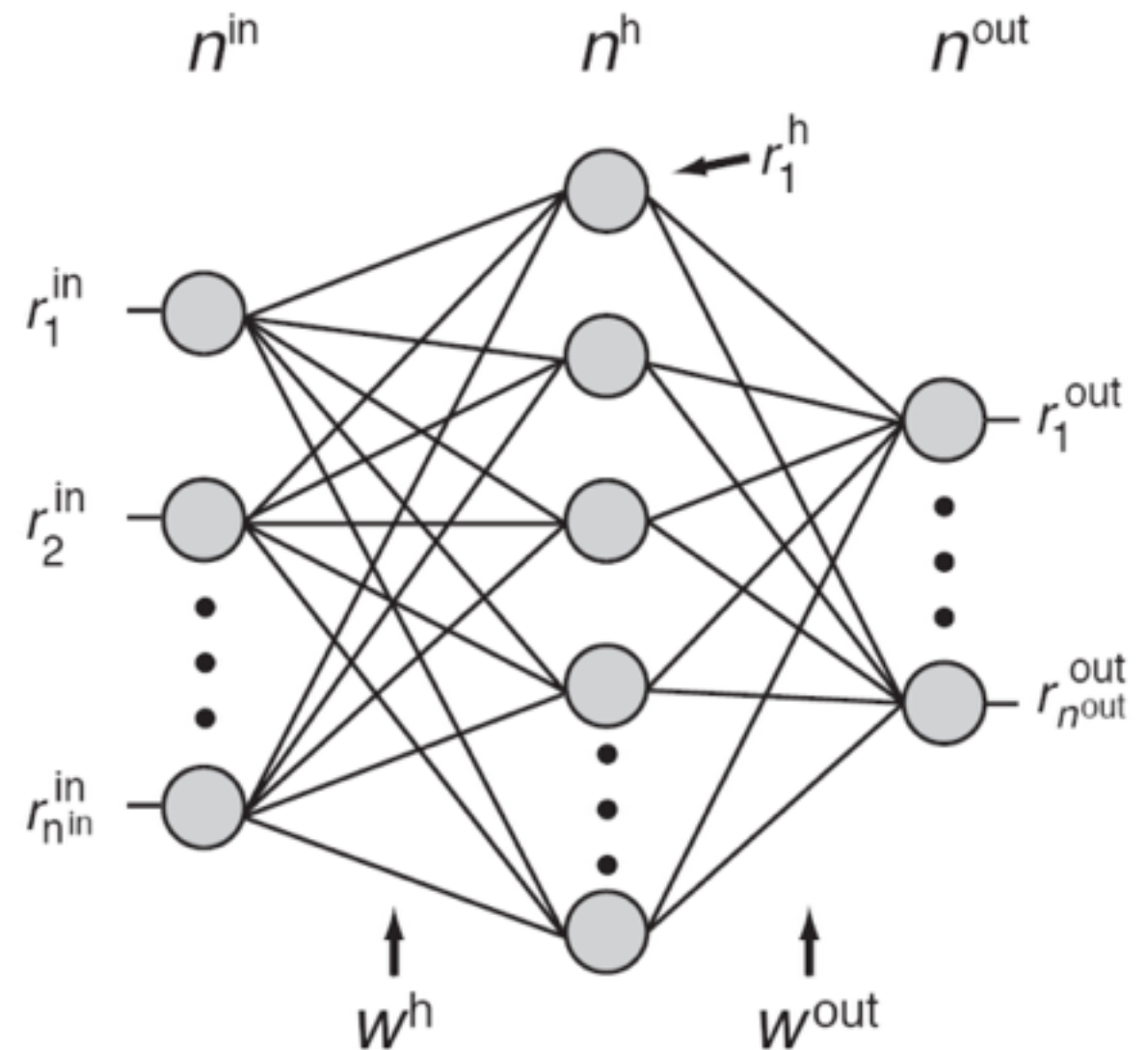
# Multi-layer perceptron

## Limitations

- Brain-like performance doesn't equate with actual performance
- Training rules are non-physiologic

## Strengths

- Hidden layer activity might resemble brain function (with appropriate inputs/outputs)
- Brain = mapping network
- Self-Organization, like the brain
- High flexibility in possible computations



Trappenberg 2010

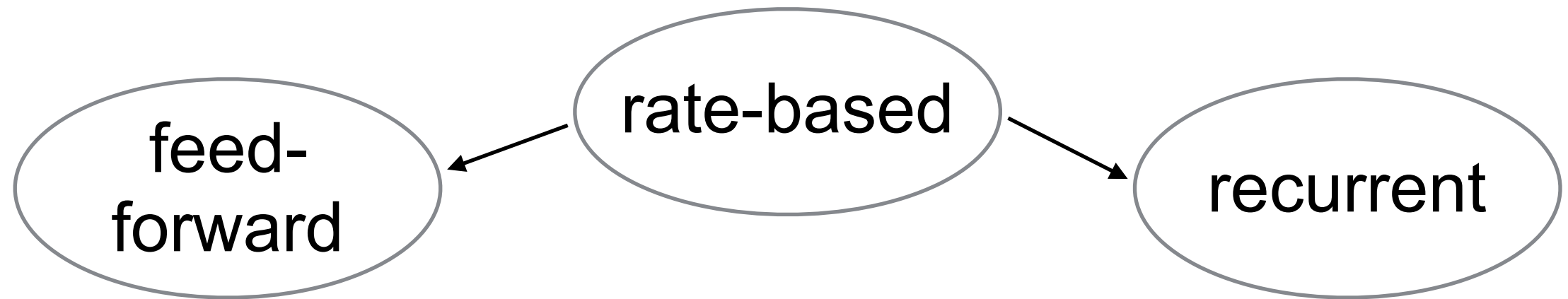
**Point: MLP is usually good for machine learning purposes, but is not necessarily good for neuroscience theory all the time!**

# Two types of NNets

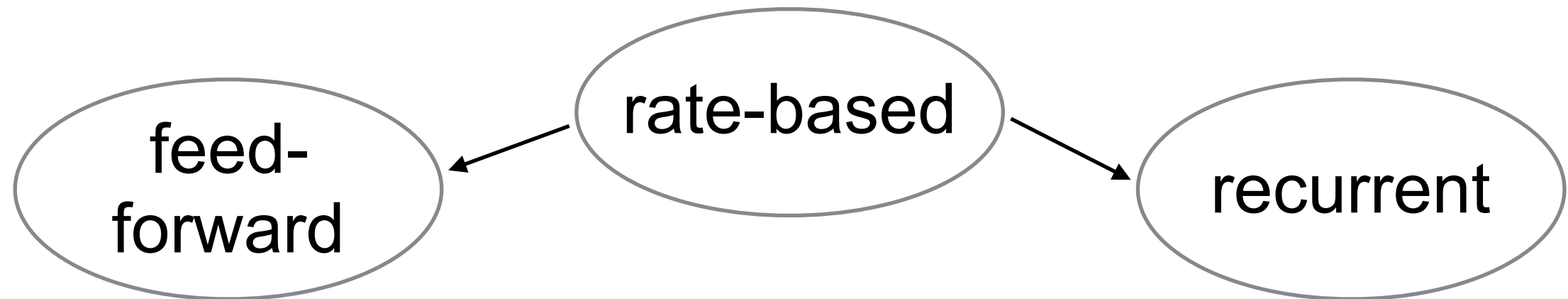


rate-based

# Two types of NNets

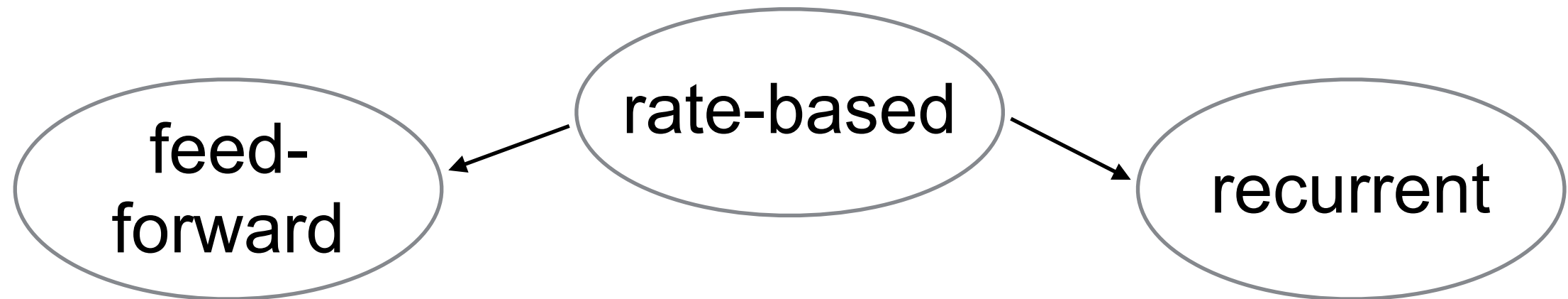


# Two types of NNets



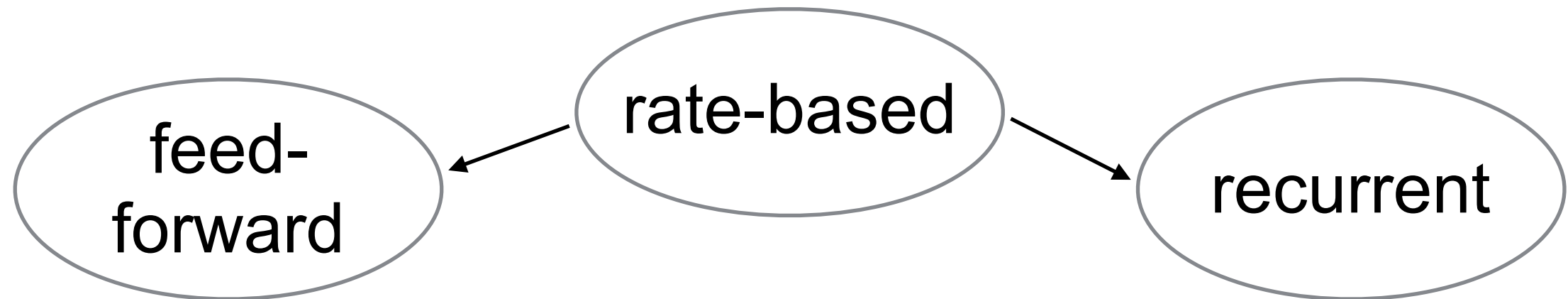
**Feed-forward:** information *only* flows forward, simplest connectivity

# Two types of NNets



**Feed-forward:** information *only* flows forward, simplest connectivity

# Two types of NNets

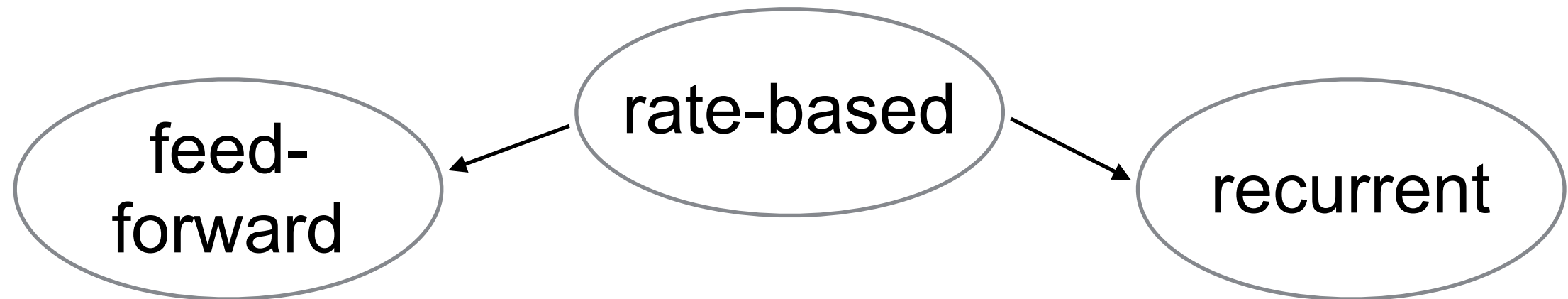


**Feed-forward:** information *only* flows forward, simplest connectivity

- Kernel machines



# Two types of NNets



**Feed-forward:** information *only* flows forward, simplest connectivity

- Kernel machines
- Radial basis function networks

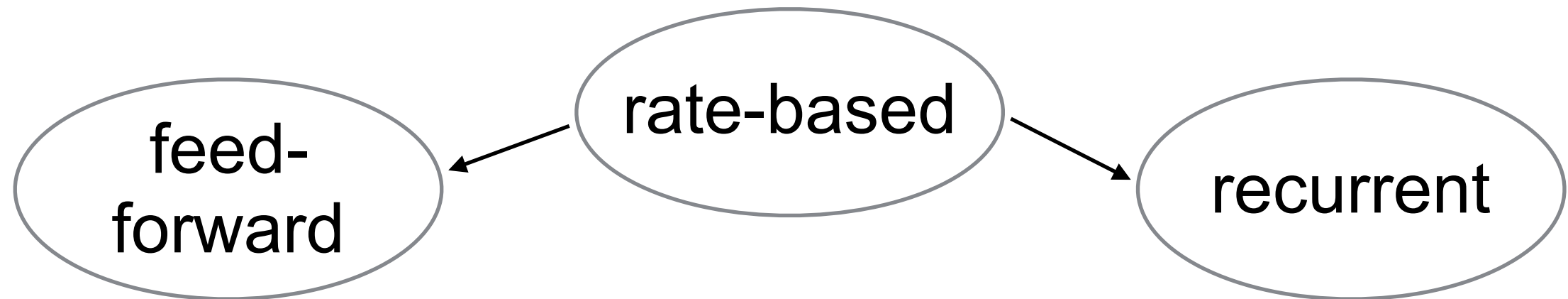
# Two types of NNets



**Feed-forward:** information *only* flows forward, simplest connectivity

- Kernel machines
- Radial basis function networks
- Probabilistic mapping networks

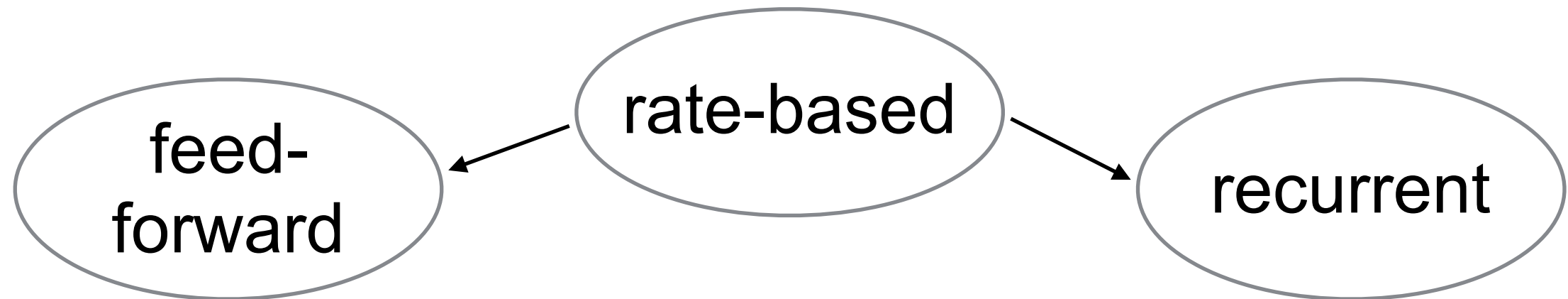
# Two types of NNets



**Feed-forward:** information *only* flows forward, simplest connectivity

- Kernel machines
- Radial basis function networks
- Probabilistic mapping networks
- Bayesian nets

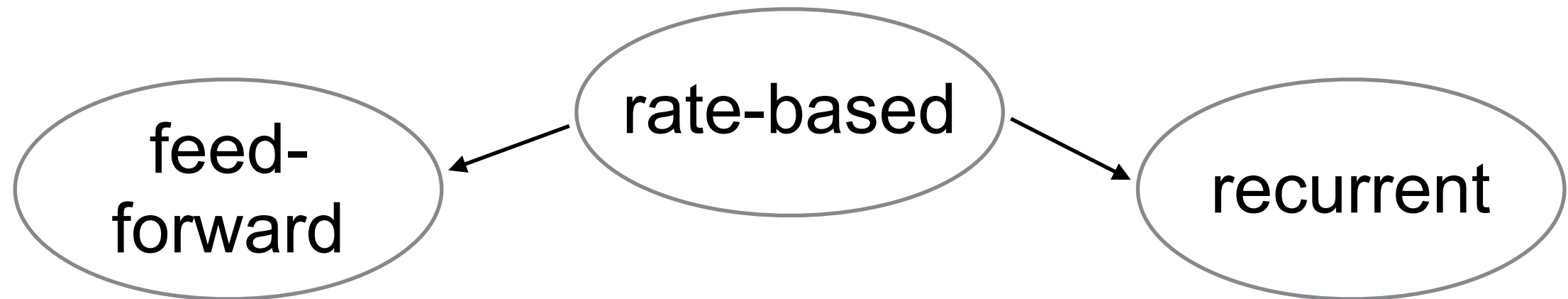
# Two types of NNets



**Feed-forward:** information *only* flows forward, simplest connectivity

- Kernel machines
- Radial basis function networks
- Probabilistic mapping networks
- Bayesian nets
- Stochastic nets

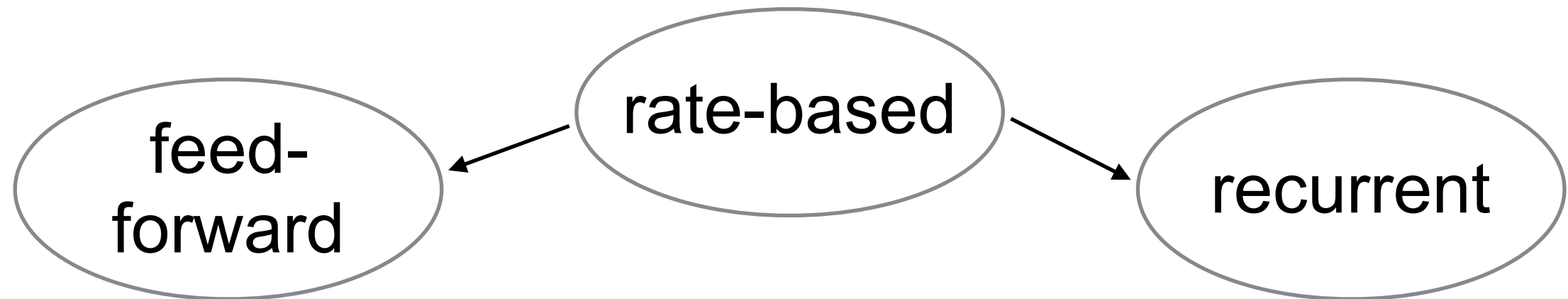
# Two types of NNets



**Feed-forward:** information *only* flows forward, simplest connectivity

- Kernel machines
- Radial basis function networks
- Probabilistic mapping networks
- Bayesian nets
- Stochastic nets
- Modular nets

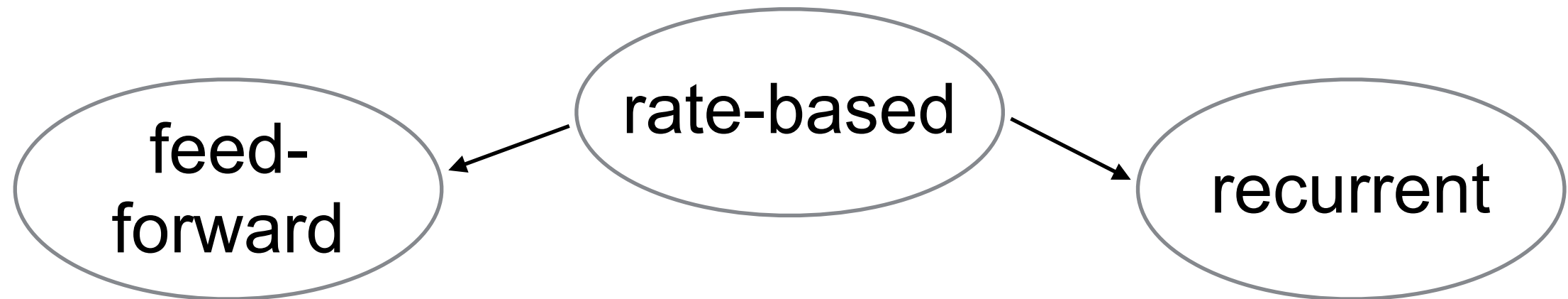
# Two types of NNets



**Feed-forward:** information *only* flows forward, simplest connectivity

- Kernel machines
- Radial basis function networks
- Probabilistic mapping networks
- Bayesian nets
- Stochastic nets
- Modular nets
- Committee machines

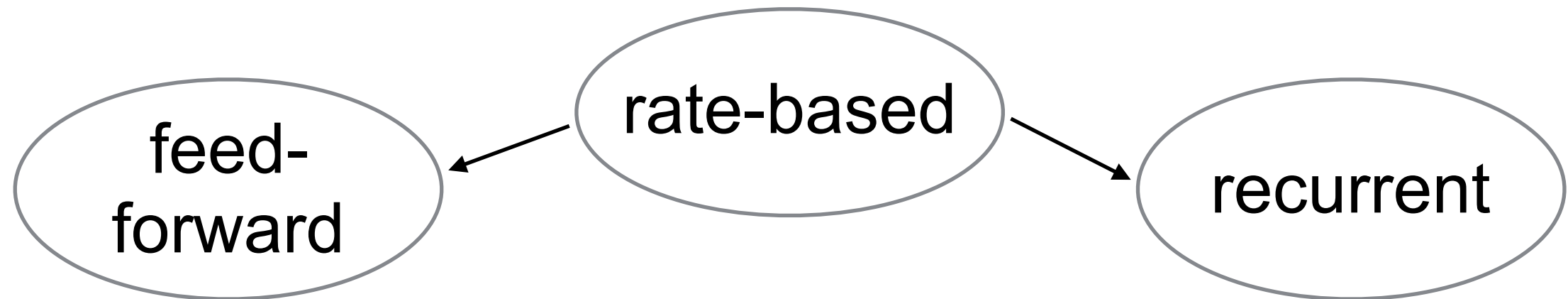
# Two types of NNets



**Feed-forward:** information *only* flows forward, simplest connectivity

- Kernel machines
- Radial basis function networks
- Probabilistic mapping networks
- Bayesian nets
- Stochastic nets
- Modular nets
- Committee machines
- Associative neural nets

# Two types of NNets

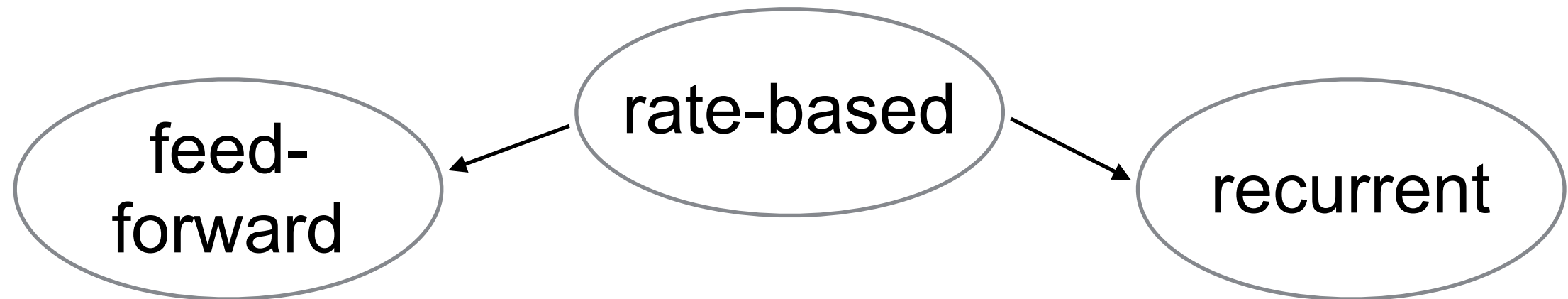


**Feed-forward:** information *only* flows forward, simplest connectivity

- Kernel machines
- Radial basis function networks
- Probabilistic mapping networks
- Bayesian nets
- Stochastic nets
- Modular nets
- Committee machines
- Associative neural nets
- Holographic associative nets (???)



# Two types of NNets



**Feed-forward:** information *only* flows forward, simplest connectivity

- Kernel machines
- Radial basis function networks
- Probabilistic mapping networks
- Bayesian nets
- Stochastic nets
- Modular nets
- Committee machines
- Associative neural nets
- Holographic associative nets (???)
- Fuzzy nets

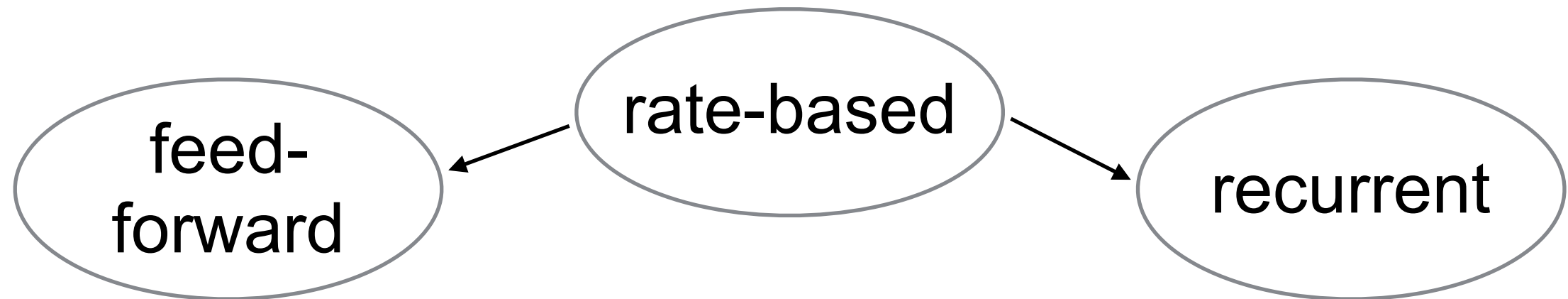
# Two types of NNets



**Feed-forward:** information *only* flows forward, simplest connectivity

- Kernel machines
- Radial basis function networks
- Probabilistic mapping networks
- Bayesian nets
- Stochastic nets
- Modular nets
- Committee machines
- Associative neural nets
- Holographic associative nets (???)
- Fuzzy nets
- Compositional pattern-producing nets

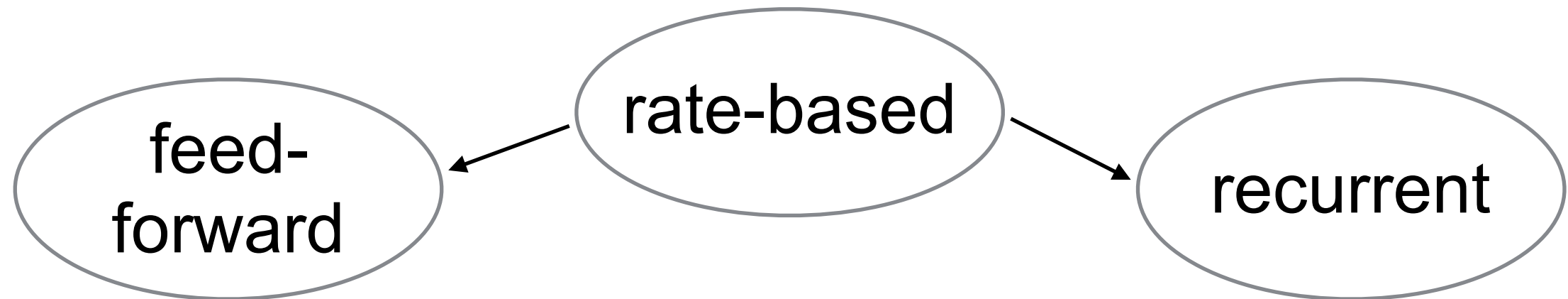
# Two types of NNets



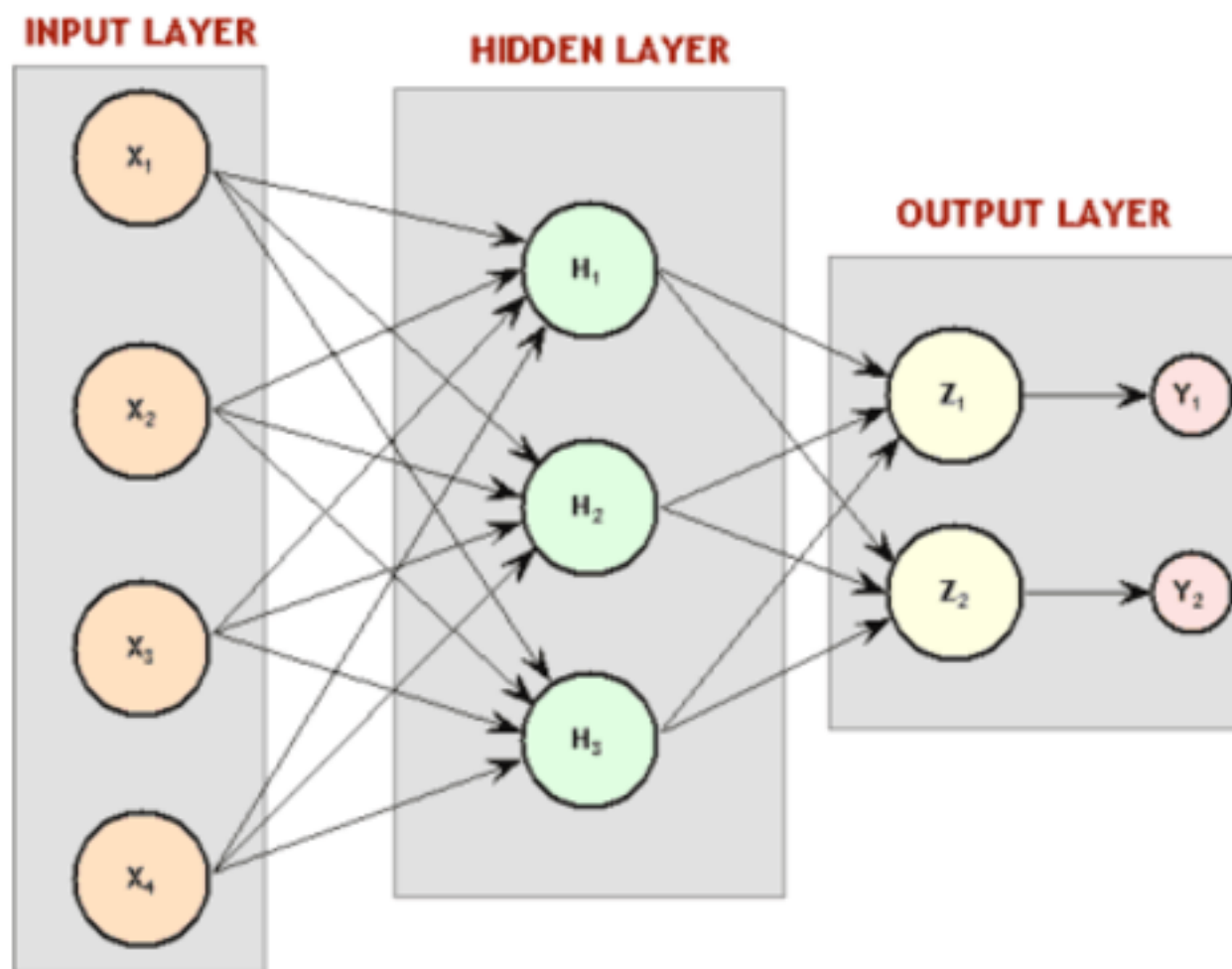
**Feed-forward:** information *only* flows forward, simplest connectivity

- Kernel machines
- Radial basis function networks
- Probabilistic mapping networks
- Bayesian nets
- Stochastic nets
- Modular nets
- Committee machines
- Associative neural nets
- Holographic associative nets (???)
- Fuzzy nets
- Compositional pattern-producing nets
- etc.

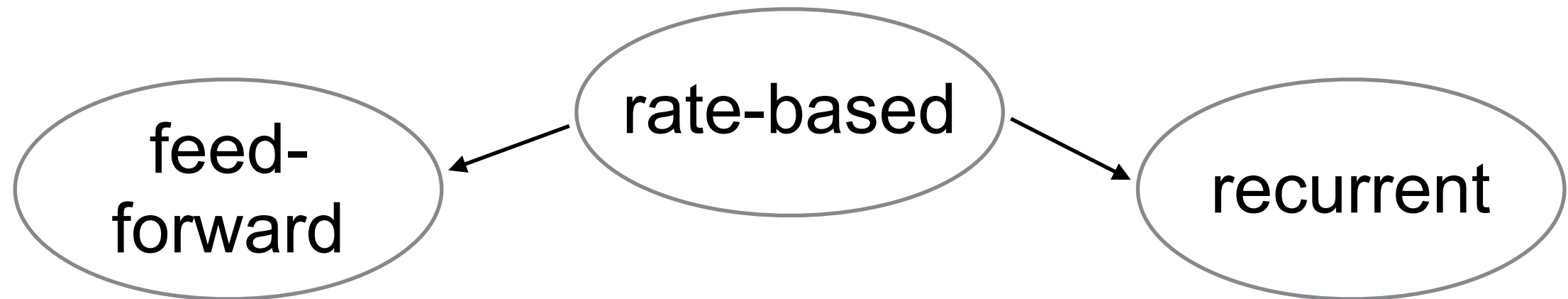
# Two types of NNets



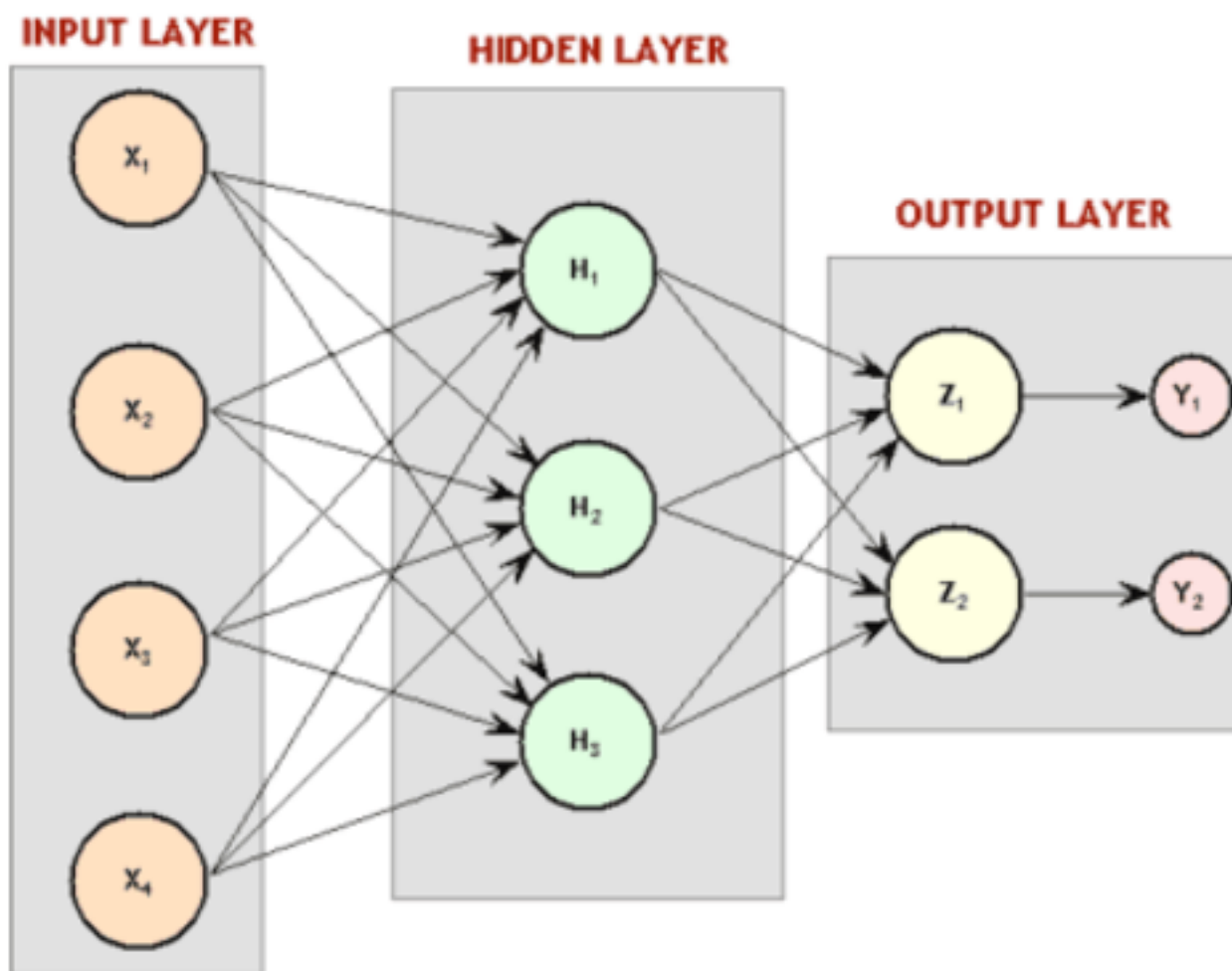
**Feed-forward:** information *only* flows forward, simplest connectivity



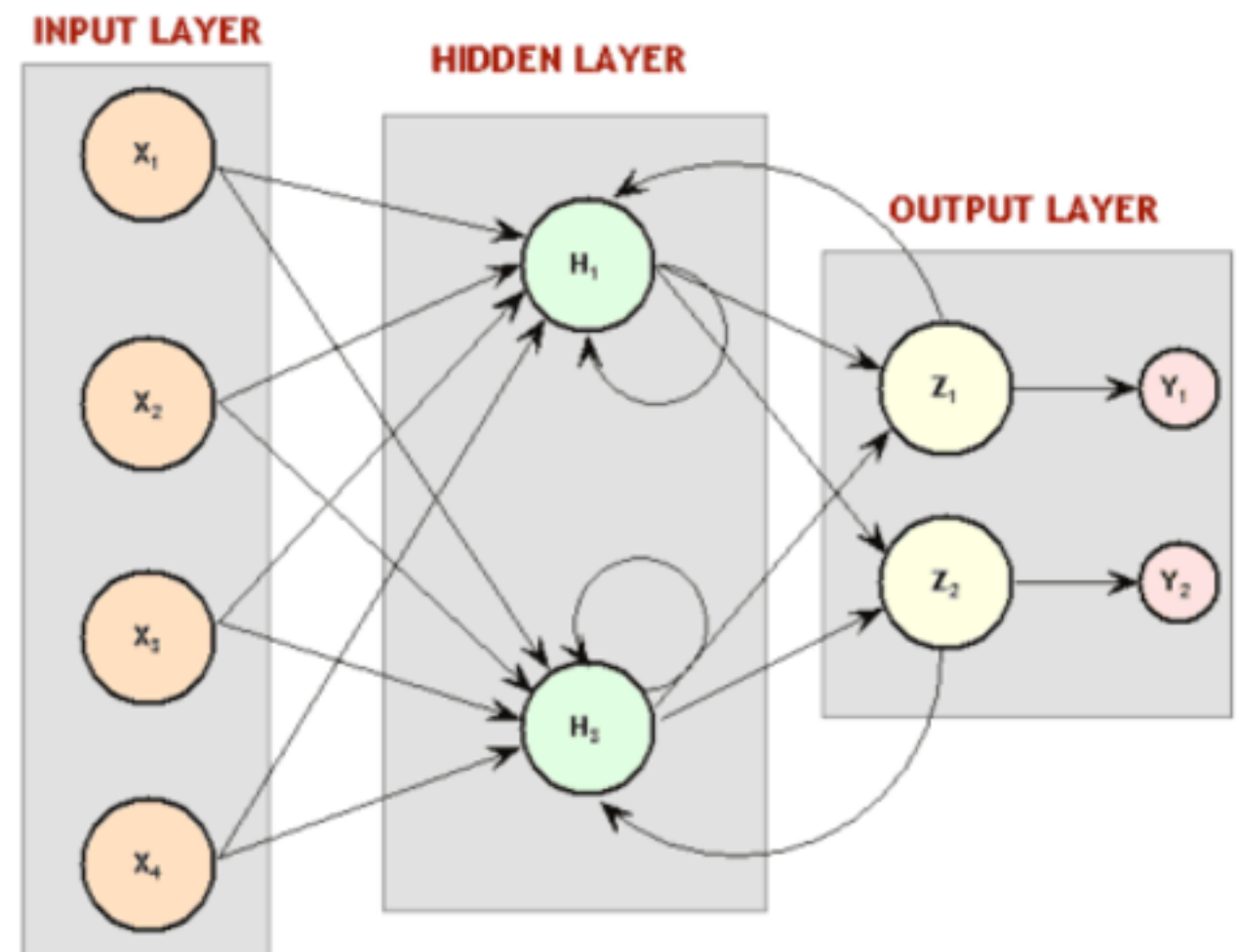
# Two types of NNets



**Feed-forward:** information *only* flows forward, simplest connectivity

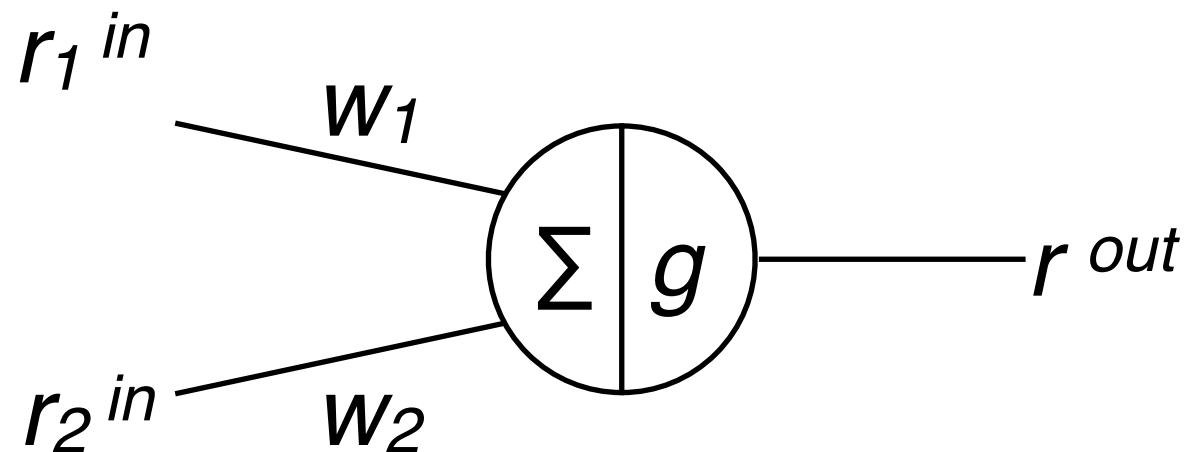


**Recurrent:** information can flow forward and backward, useful for time-resolved problems



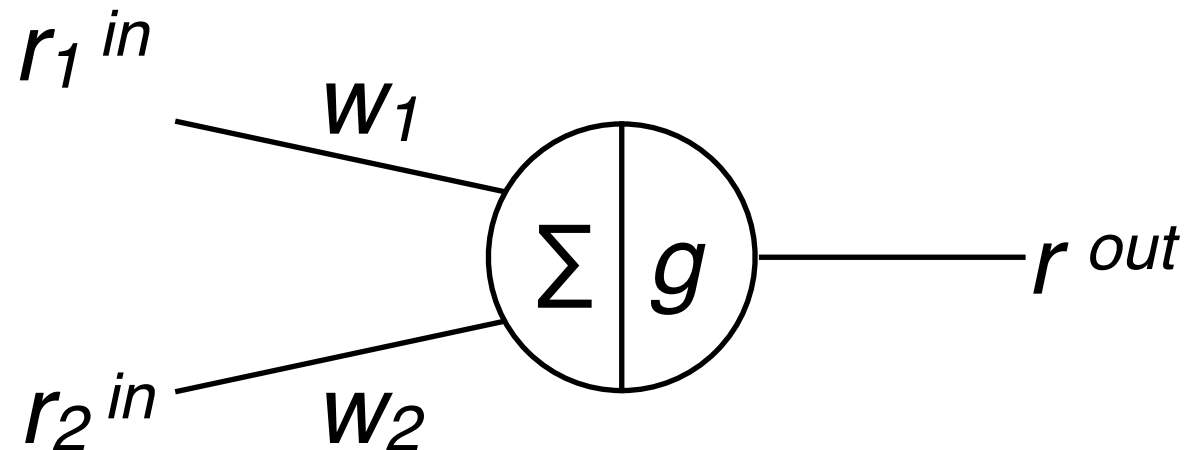
# How do we train a neural network?

**Minimize some cost function... gradient descent!**



# How do we train a neural network?

**Minimize some cost function... gradient descent!**

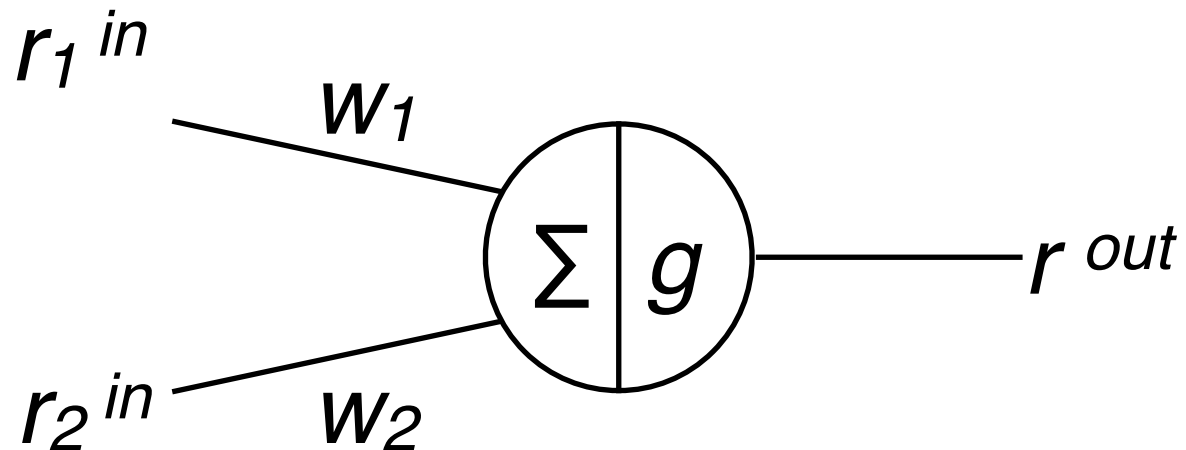


**MSE**

$$E = \frac{1}{2} \sum_i (r_i^{out} - y_i)^2$$

# How do we train a neural network?

**Minimize some cost function... gradient descent!**



**MSE**

$$E = \frac{1}{2} \sum_i (r_i^{out} - y_i)^2$$

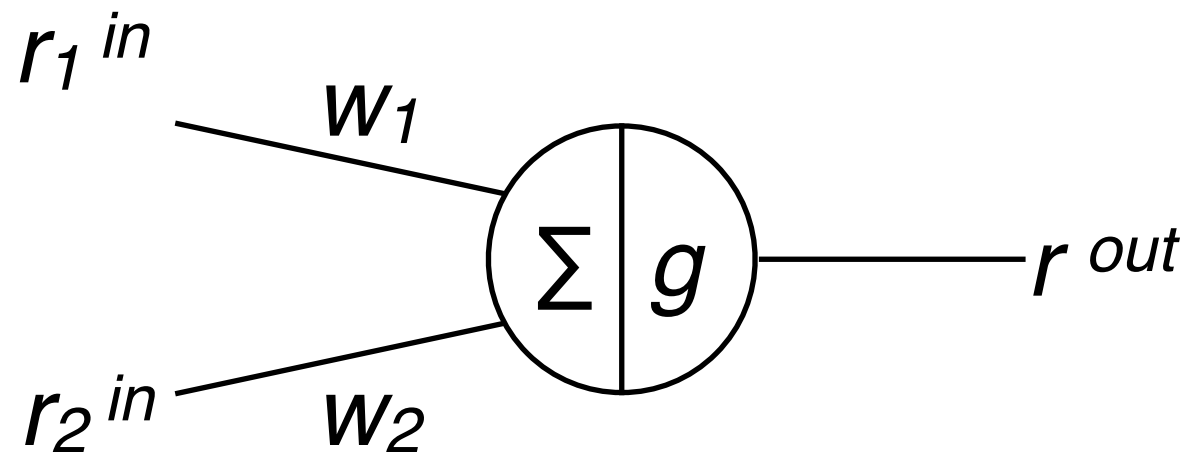
desired output  
(data, training set)

An arrow points from the text "desired output (data, training set)" to the  $y_i$  term in the equation above.



# How do we train a neural network?

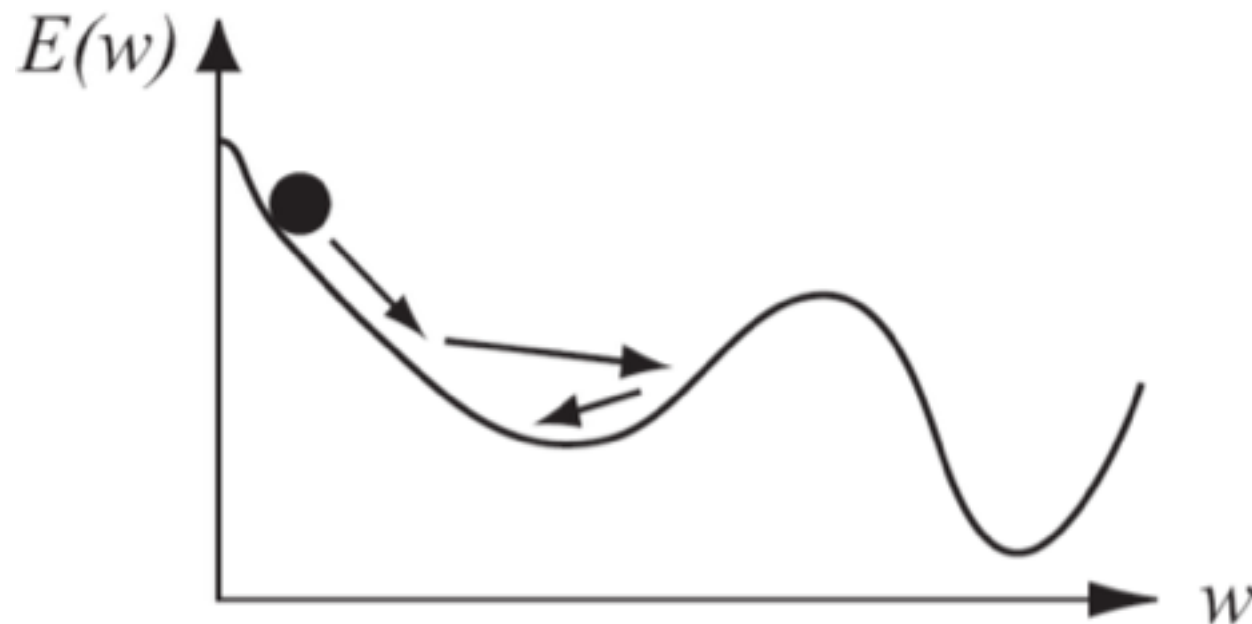
**Minimize some cost function... gradient descent!**



**MSE**

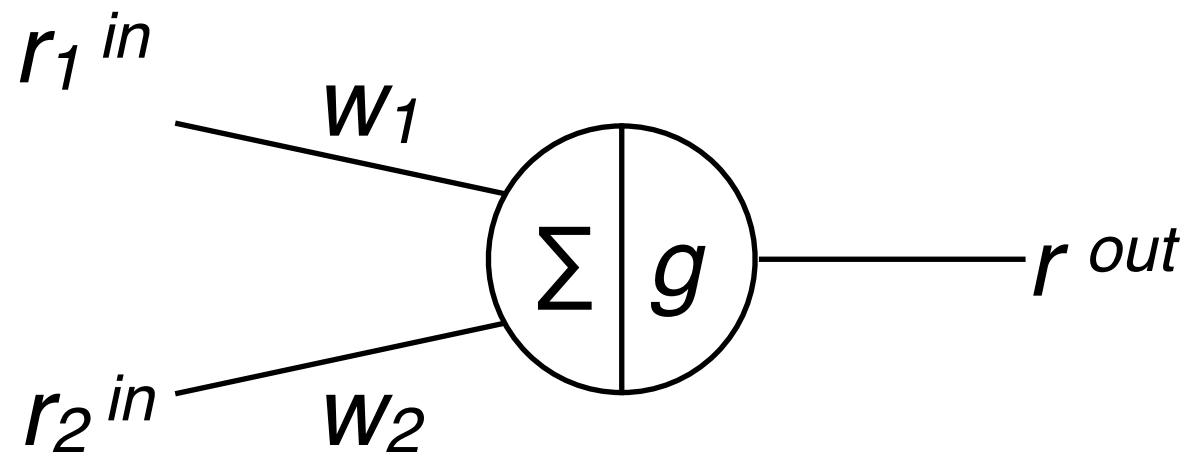
$$E = \frac{1}{2} \sum_i (r_i^{out} - y_i)^2$$

desired output  
(data, training set)



# How do we train a neural network?

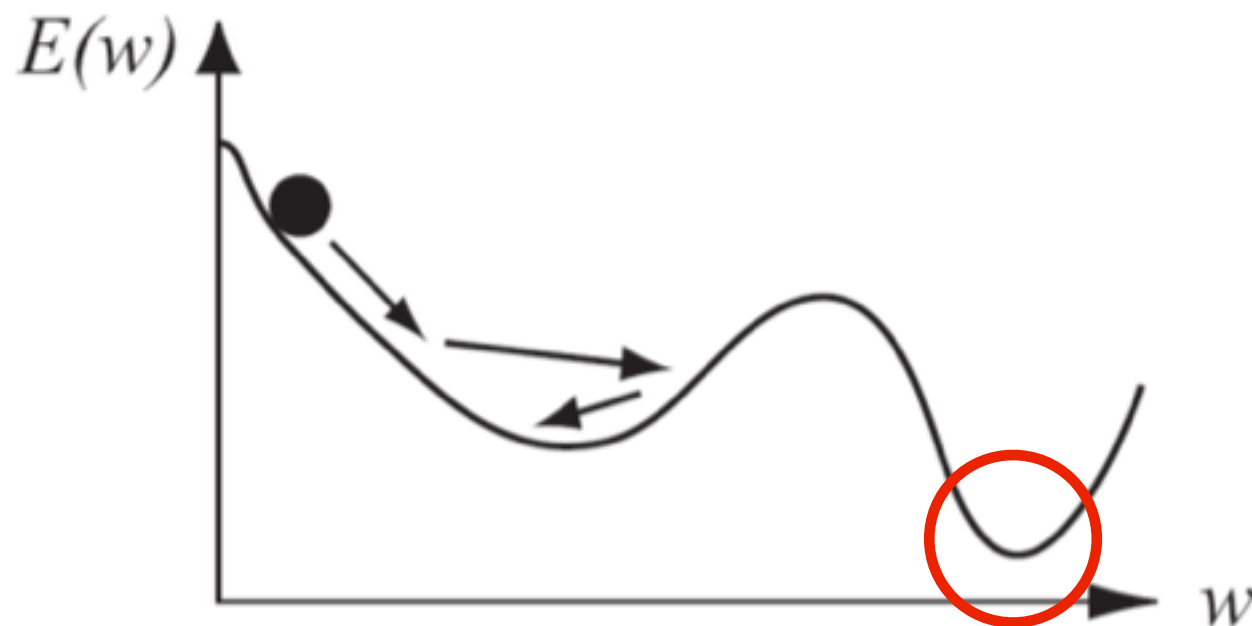
**Minimize some cost function... gradient descent!**



**MSE**

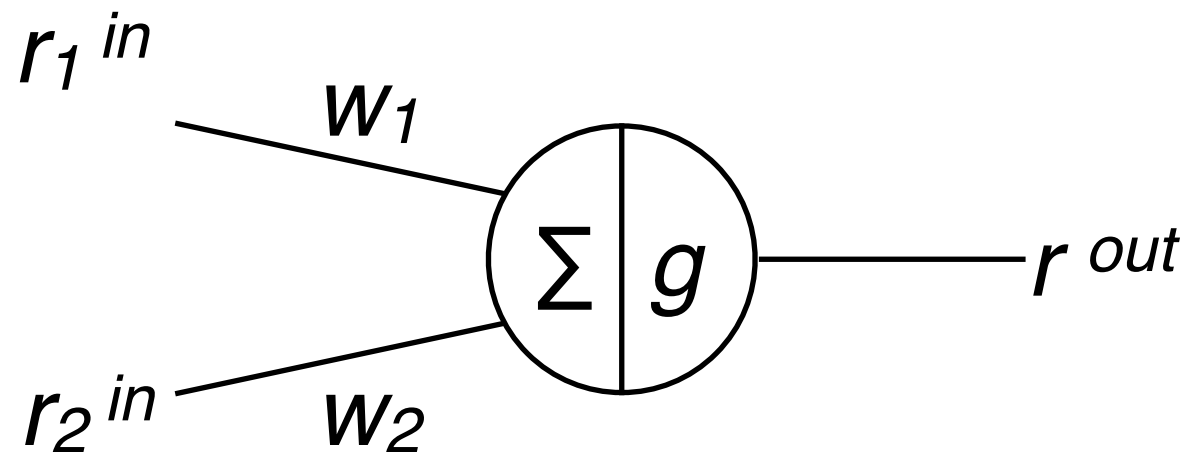
$$E = \frac{1}{2} \sum_i (r_i^{out} - y_i)^2$$

desired output  
(data, training set)



# How do we train a neural network?

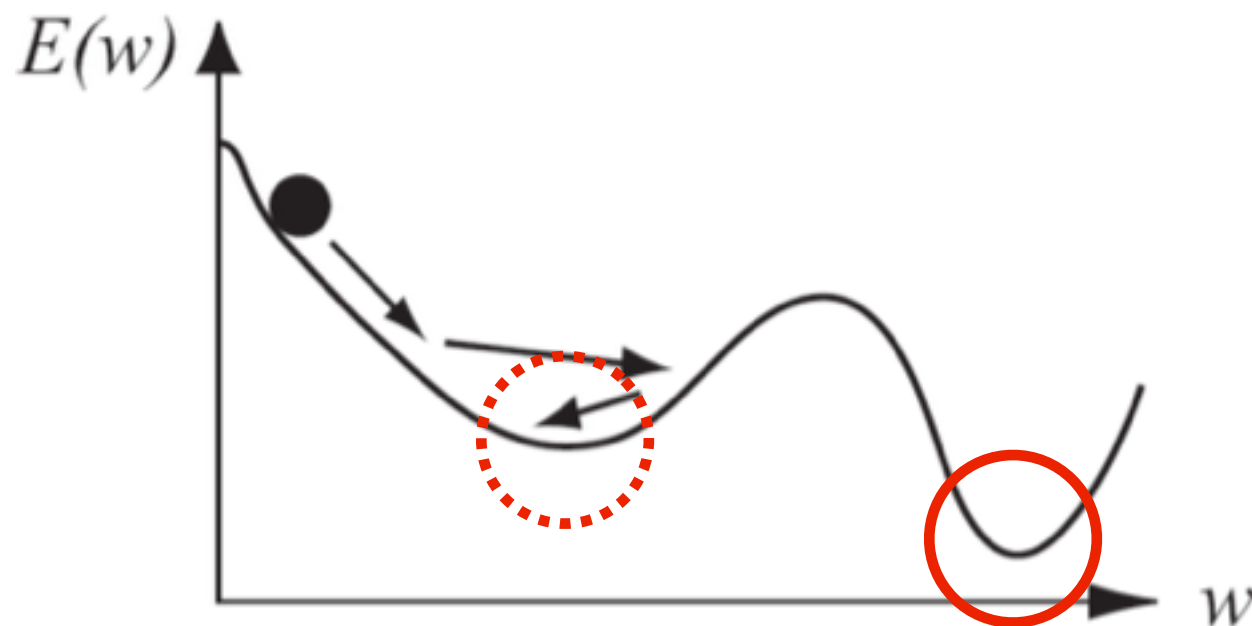
**Minimize some cost function... gradient descent!**



**MSE**

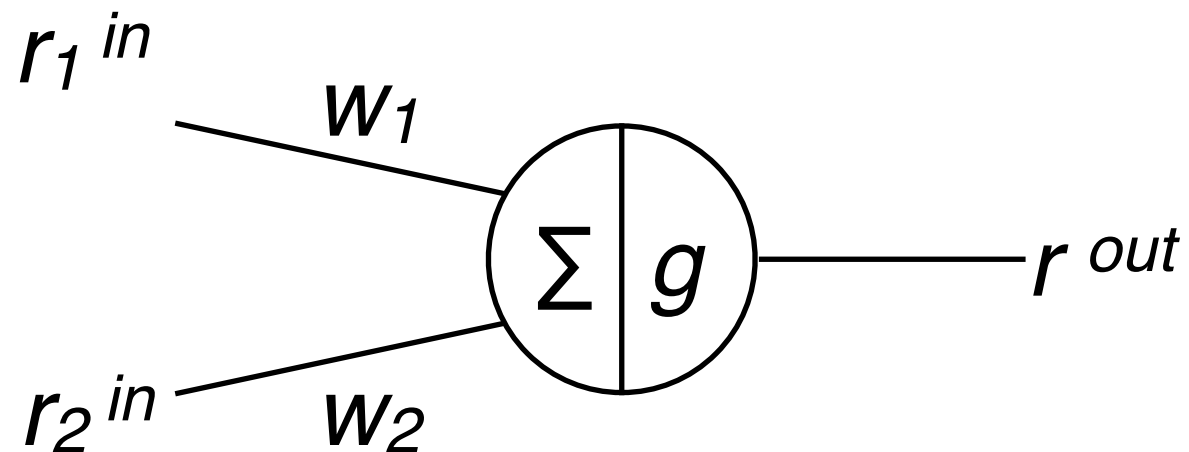
$$E = \frac{1}{2} \sum_i (r_i^{out} - y_i)^2$$

desired output  
(data, training set)



# How do we train a neural network?

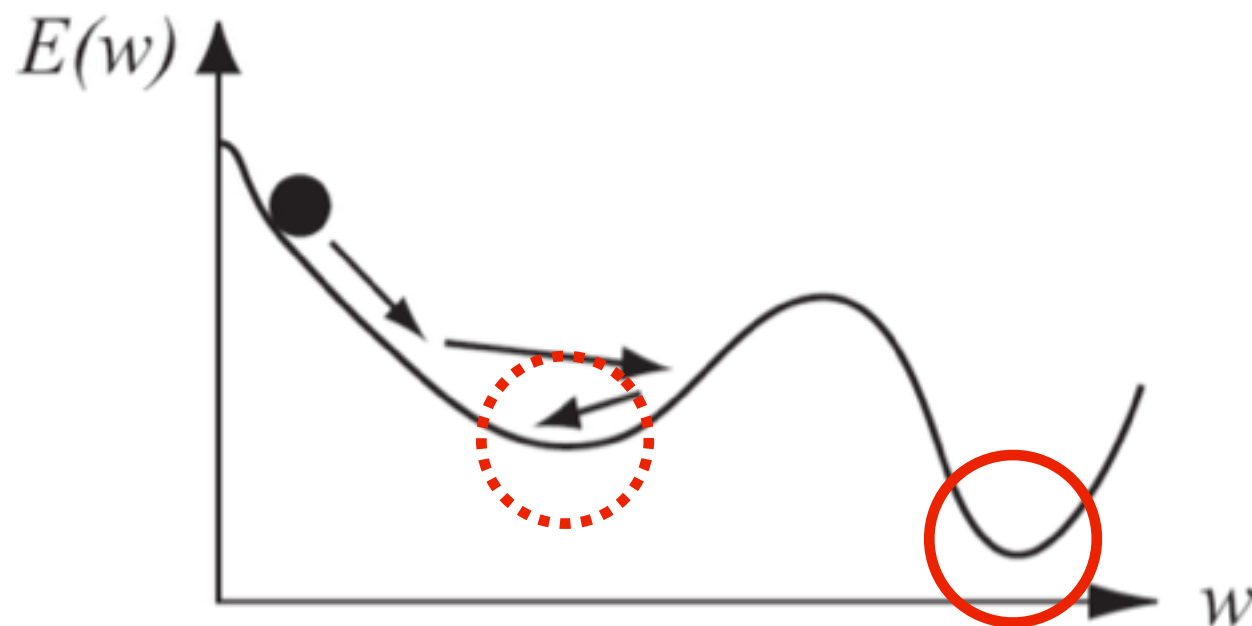
**Minimize some cost function... gradient descent!**



**MSE**

$$E = \frac{1}{2} \sum_i (r_i^{out} - y_i)^2$$

desired output  
(data, training set)

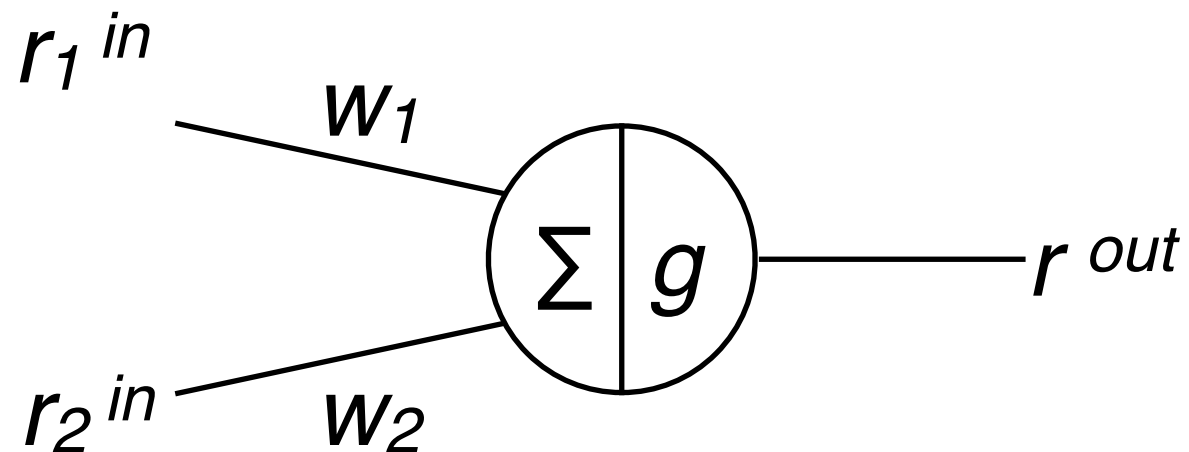


$$w_{ij} \leftarrow w_{ij} + \Delta w_{ij}$$

$$\Delta w_{ij} = -\epsilon \left( \frac{\partial E}{\partial w_{ij}} \right)$$

# How do we train a neural network?

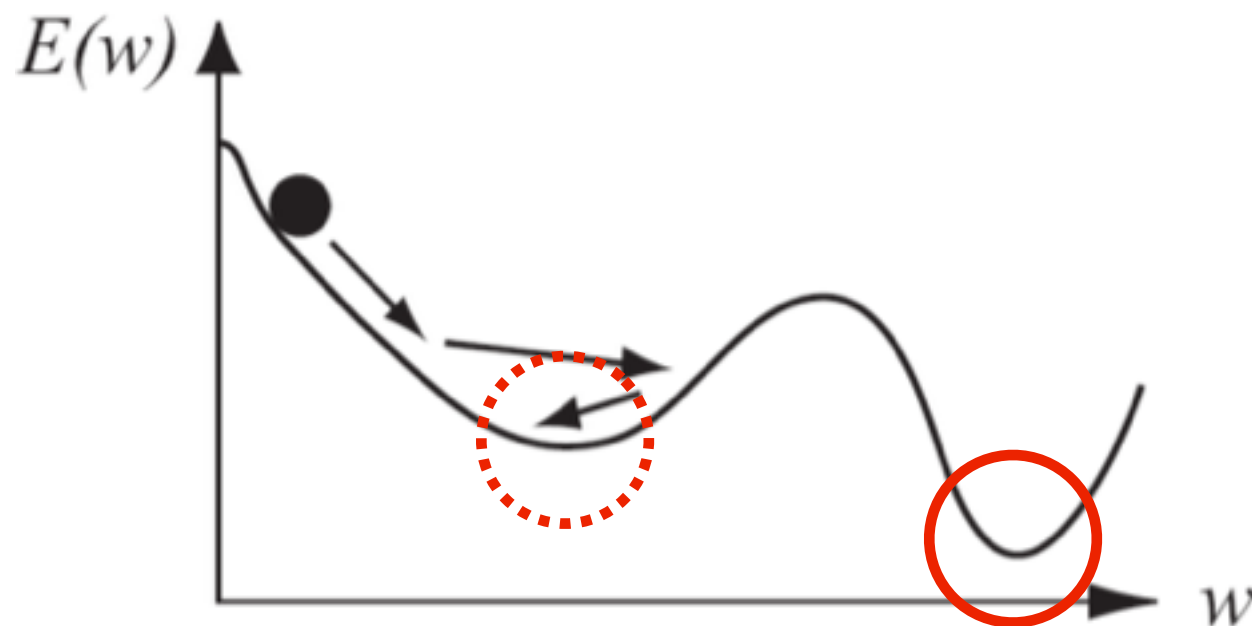
**Minimize some cost function... gradient descent!**



**MSE**

$$E = \frac{1}{2} \sum_i (r_i^{out} - y_i)^2$$

desired output  
(data, training set)



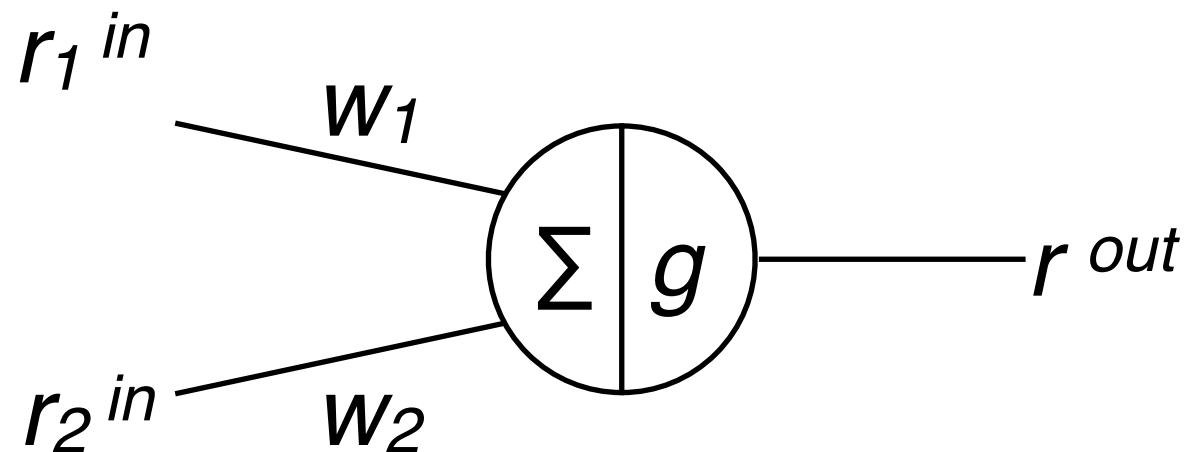
$$w_{ij} \leftarrow w_{ij} + \Delta w_{ij}$$

$$\Delta w_{ij} = -\epsilon \left( \frac{\partial E}{\partial w_{ij}} \right)$$

learning  
rate

# How do we train a neural network?

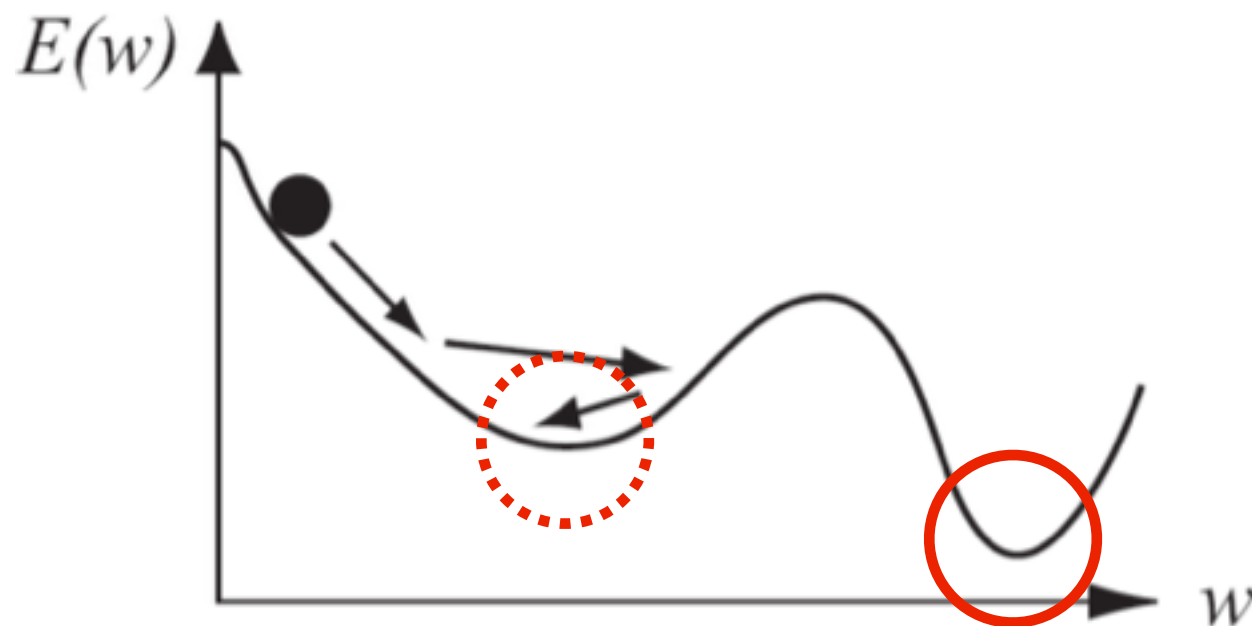
**Minimize some cost function... gradient descent!**



**MSE**

$$E = \frac{1}{2} \sum_i (r_i^{out} - y_i)^2$$

desired output  
(data, training set)



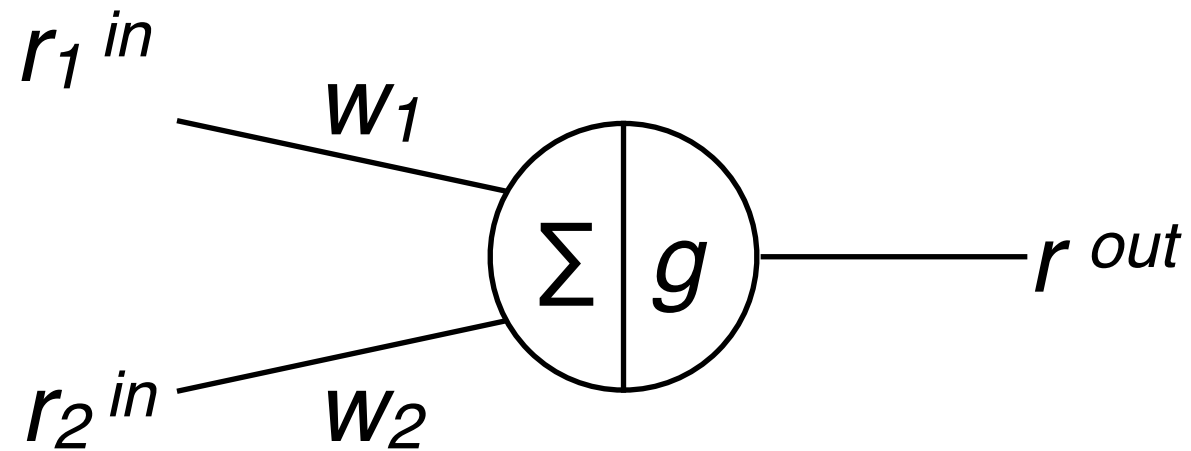
$$w_{ij} \leftarrow w_{ij} + \Delta w_{ij}$$

$$\Delta w_{ij} = -\epsilon \left( \frac{\partial E}{\partial w_{ij}} \right)$$

learning  
rate

gradient of MSE  
wrt weights

# How do we train a neural network?

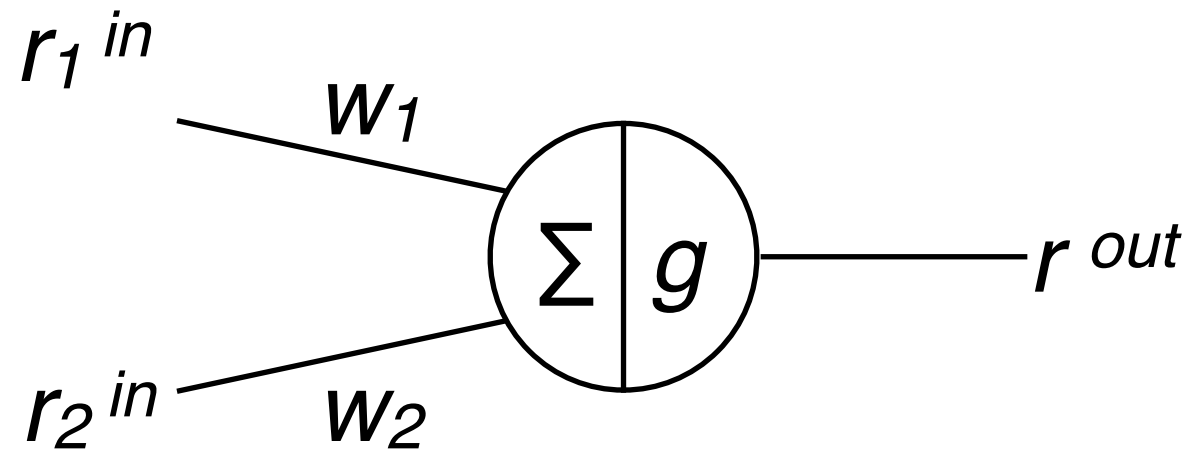


$$E = \frac{1}{2} \sum_i (r_i^{out} - y_i)^2$$

$$w_{ij} \leftarrow w_{ij} + \Delta w_{ij}$$

$$\Delta w_{ij} = -\epsilon \left( \frac{\partial E}{\partial w_{ij}} \right)$$

# How do we train a neural network?



$$E = \frac{1}{2} \sum_i (r_i^{out} - y_i)^2$$

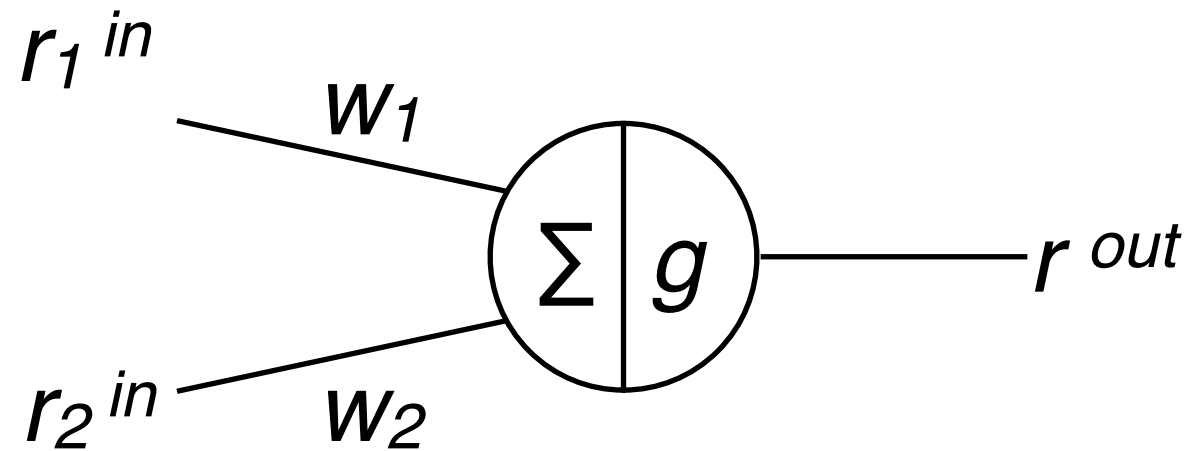
$$w_{ij} \leftarrow w_{ij} + \Delta w_{ij}$$

$$\Delta w_{ij} = -\epsilon \left( \frac{\partial E}{\partial w_{ij}} \right)$$

change in weight i,j depends  
on learning rate and dependence of  
error on weight change at i,j



# How do we train a neural network?



$$E = \frac{1}{2} \sum_i (r_i^{out} - y_i)^2$$

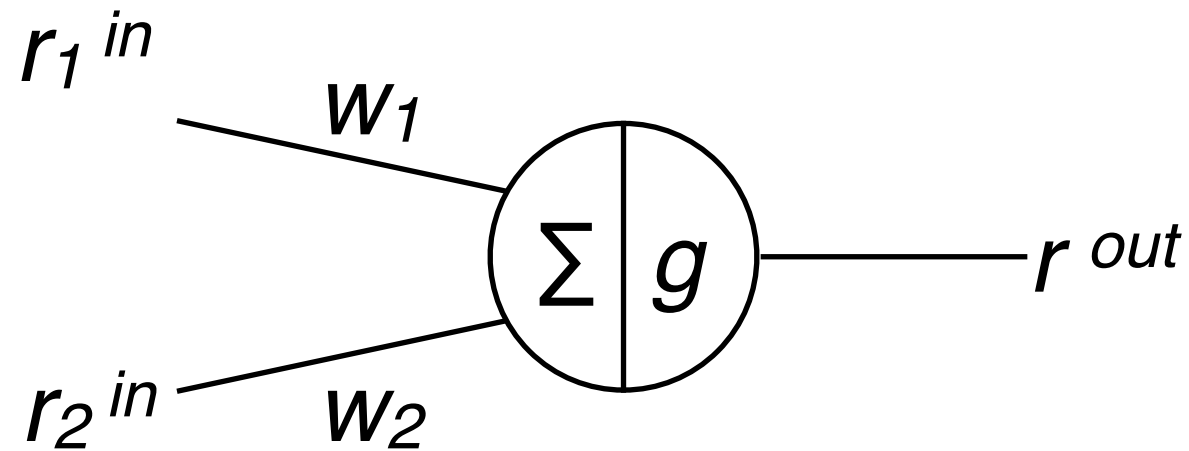
$$w_{ij} \leftarrow w_{ij} + \Delta w_{ij}$$

$$\Delta w_{ij} = -\epsilon \left( \frac{\partial E}{\partial w_{ij}} \right)$$

$$\frac{\partial E}{\partial w_{ij}} = \frac{1}{2} \frac{\partial}{\partial w_{ij}} \sum_i \underbrace{\left( g \left( \underbrace{\sum_j w_{ij} r_j^{in}}_{h_i} \right) - y_i \right)^2}_f$$

change in weight i,j depends  
on learning rate and dependence of  
error on weight change at i,j

# How do we train a neural network?



$$E = \frac{1}{2} \sum_i (r_i^{out} - y_i)^2$$

$$w_{ij} \leftarrow w_{ij} + \Delta w_{ij}$$

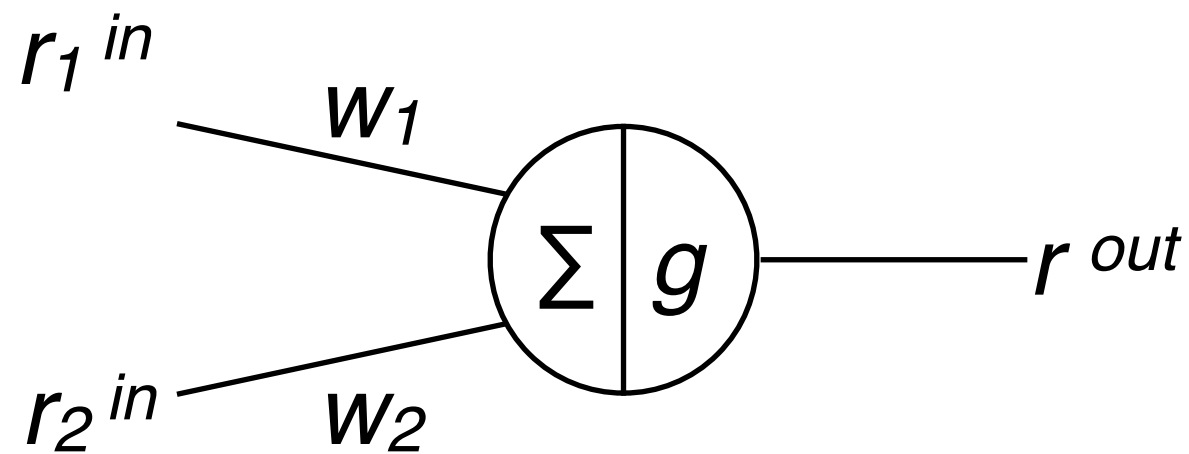
$$\Delta w_{ij} = -\epsilon \left( \frac{\partial E}{\partial w_{ij}} \right)$$

$$\frac{\partial E}{\partial w_{ij}} = \frac{1}{2} \frac{\partial}{\partial w_{ij}} \sum_i \underbrace{\left( g \left( \underbrace{\sum_j w_{ij} r_j^{in}}_{h_i} \right) - y_i \right)^2}_f$$

change in weight i,j depends  
on learning rate and dependence of  
error on weight change at i,j

error's dependence on weight i,j, rewritten  
using MSE equation

# How do we train a neural network?



$$E = \frac{1}{2} \sum_i (r_i^{out} - y_i)^2$$

$$w_{ij} \leftarrow w_{ij} + \Delta w_{ij}$$

$$\Delta w_{ij} = -\epsilon \left( \frac{\partial E}{\partial w_{ij}} \right)$$

$$\frac{\partial E}{\partial w_{ij}} = \frac{1}{2} \frac{\partial}{\partial w_{ij}} \sum_i \underbrace{\left( g \left( \underbrace{\sum_j w_{ij} r_j^{in}}_{h_i} \right) - y_i \right)^2}_f$$

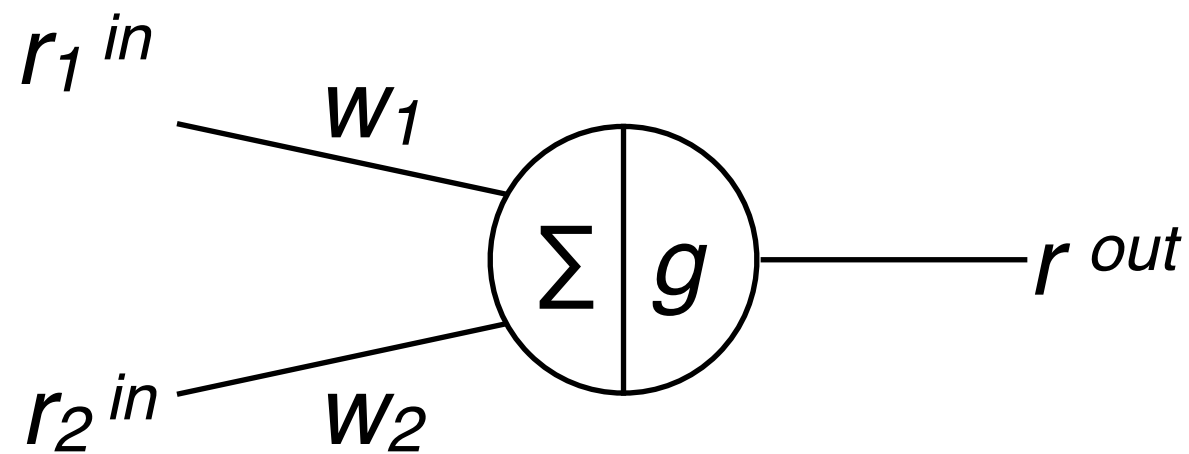
change in weight i,j depends  
on learning rate and dependence of  
error on weight change at i,j

$$\frac{\partial f}{\partial w_{ij}} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial h} \frac{\partial h}{\partial w_{ij}}$$

(using chain rule)

error's dependence on weight i,j, rewritten  
using MSE equation

# How do we train a neural network?



$$E = \frac{1}{2} \sum_i (r_i^{out} - y_i)^2$$

$$w_{ij} \leftarrow w_{ij} + \Delta w_{ij}$$

$$\Delta w_{ij} = -\epsilon \left( \frac{\partial E}{\partial w_{ij}} \right)$$

$$\frac{\partial E}{\partial w_{ij}} = \frac{1}{2} \frac{\partial}{\partial w_{ij}} \sum_i \underbrace{\left( g \left( \underbrace{\sum_j w_{ij} r_j^{in}}_{h_i} \right) - y_i \right)^2}_f$$

change in weight i,j depends on learning rate and dependence of error on weight change at i,j

$$\frac{\partial f}{\partial w_{ij}} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial h} \frac{\partial h}{\partial w_{ij}}$$

(using chain rule)

error's dependence on weight i,j, rewritten using MSE equation

$$\Delta w_{ij} = \epsilon (g'(h_i) (y_i - r_i^{out}) r_j^{in}) \quad \textbf{learning rule (derivation?)}$$

# Adding layers...

Start by finding the output rates...

**2-layer (1 hidden) perceptron**

$$E = \frac{1}{2} \sum_i (r_i^{out} - y_i)^2$$

$$\Delta w_{ij} = \varepsilon (g'(h_i) (y_i - r_i^{out}) r_j^{in})$$

# Adding layers...

Start by finding the output rates...

**2-layer (1 hidden) perceptron**

$$\mathbf{r}^{out} = g(\mathbf{w}^{out} \mathbf{r}^h)$$

$$r_i^{out} = g\left(\sum_j w_{ij}^{out} r_j^h\right)$$

$$E = \frac{1}{2} \sum_i (r_i^{out} - y_i)^2$$

$$\Delta w_{ij} = \varepsilon (g'(h_i) (y_i - r_i^{out}) r_j^{in})$$

# Adding layers...

Start by finding the output rates...

**2-layer (1 hidden) perceptron**

$$\mathbf{r}^{out} = g(\mathbf{w}^{out} \mathbf{r}^h)$$

$$r_i^{out} = g\left(\sum_j w_{ij}^{out} r_j^h\right)$$

$$\mathbf{r}^{out} = g^{out}\left(\mathbf{w}^{out} g^h(\mathbf{w}^h \mathbf{r}^{in})\right)$$

$$E = \frac{1}{2} \sum_i (r_i^{out} - y_i)^2$$

$$\Delta w_{ij} = \varepsilon (g'(h_i) (y_i - r_i^{out}) r_j^{in})$$

# Adding layers...

Start by finding the output rates...

**2-layer (1 hidden) perceptron**

$$\mathbf{r}^{out} = g(\mathbf{w}^{out} \mathbf{r}^h)$$

$$r_i^{out} = g\left(\sum_j w_{ij}^{out} r_j^h\right)$$

$$\mathbf{r}^{out} = g^{out}\left(\mathbf{w}^{out} g^h\left(\mathbf{w}^h \mathbf{r}^{in}\right)\right)$$

**3-layer (2 hidden) perceptron**

$$\mathbf{r}^{out} = g^{out}\left(\mathbf{w}^{out} g^{h_{out-1}}\left(\mathbf{w}^{h_{out-1}} g^{h_{out-2}}\left(\mathbf{w}^{h_{out-2}} \mathbf{r}^{in}\right)\right)\right)$$

$$E = \frac{1}{2} \sum_i (r_i^{out} - y_i)^2$$

$$\Delta w_{ij} = \varepsilon (g'(h_i) (y_i - r_i^{out}) r_j^{in})$$



# Adding layers...

Start by finding the output rates...

**2-layer (1 hidden) perceptron**

$$\mathbf{r}^{out} = g(\mathbf{w}^{out} \mathbf{r}^h)$$

$$r_i^{out} = g\left(\sum_j w_{ij}^{out} r_j^h\right)$$

$$E = \frac{1}{2} \sum_i (r_i^{out} - y_i)^2$$

$$\Delta w_{ij} = \varepsilon (g'(h_i) (y_i - r_i^{out}) r_j^{in})$$

$$\mathbf{r}^{out} = g^{out}(\mathbf{w}^{out} g^h(\mathbf{w}^h \mathbf{r}^{in}))$$

**3-layer (2 hidden) perceptron**

$$\mathbf{r}^{out} = g^{out}(\mathbf{w}^{out} g^{h_{out-1}}(\mathbf{w}^{h_{out-1}} g^{h_{out-2}}(\mathbf{w}^{h_{out-2}} \mathbf{r}^{in})))$$

**n-layer (n-1 hidden) perceptron**

$$\mathbf{r}^{out} = g^{out}(\mathbf{w}^{out} g^{h_{out-1}}(\mathbf{w}^{h_{out-1}} g^{h_{out-2}}(\mathbf{w}^{h_{out-2}} \dots g^{h_{out-n+1}}(\mathbf{w}^{h_{out-n+1}} g^{h_{out-n}}(\mathbf{w}^{h_{out-n}} \mathbf{r}^{in}))))))$$

# Adding layers...

Start by finding the output rates...

**2-layer (1 hidden) perceptron**

$$\mathbf{r}^{out} = g(\mathbf{w}^{out} \mathbf{r}^h)$$

$$r_i^{out} = g\left(\sum_j w_{ij}^{out} r_j^h\right)$$

$$E = \frac{1}{2} \sum_i (r_i^{out} - y_i)^2$$

$$\Delta w_{ij} = \varepsilon (g'(h_i) (y_i - r_i^{out}) r_j^{in})$$

$$\mathbf{r}^{out} = g^{out}(\mathbf{w}^{out} g^h(\mathbf{w}^h \mathbf{r}^{in}))$$

**3-layer (2 hidden) perceptron**

$$\mathbf{r}^{out} = g^{out}(\mathbf{w}^{out} g^{h_{out-1}}(\mathbf{w}^{h_{out-1}} g^{h_{out-2}}(\mathbf{w}^{h_{out-2}} \mathbf{r}^{in})))$$

**n-layer (n-1 hidden) perceptron**

$$\mathbf{r}^{out} = g^{out}(\mathbf{w}^{out} g^{h_{out-1}}(\mathbf{w}^{h_{out-1}} g^{h_{out-2}}(\mathbf{w}^{h_{out-2}} \dots g^{h_{out-n+1}}(\mathbf{w}^{h_{out-n+1}} g^{h_{out-n}}(\mathbf{w}^{h_{out-n}} \mathbf{r}^{in}))))))$$

**Idea is to nest each layer's output rates within the next...**

# Training through the layers...

**Generalized delta rule (output wts)**

$$\frac{\partial E}{\partial w_{ij}^{out}} = \frac{1}{2} \frac{\partial}{\partial w_{ij}^{out}} \sum_i (r_i^{out} - y_i)^2$$

$$E = \frac{1}{2} \sum_i (r_i^{out} - y_i)^2$$

$$\Delta w_{ij} = \varepsilon (g'(h_i) (y_i - r_i^{out}) r_j^{in})$$

# Training through the layers...

**Generalized delta rule (output wts)**

$$\begin{aligned}\frac{\partial E}{\partial w_{ij}^{out}} &= \frac{1}{2} \frac{\partial}{\partial w_{ij}^{out}} \sum_i (r_i^{out} - y_i)^2 \\ &= \delta_i^{out} r_j^h\end{aligned}$$

$$E = \frac{1}{2} \sum_i (r_i^{out} - y_i)^2$$

$$\Delta w_{ij} = \varepsilon (g'(h_i) (y_i - r_i^{out}) r_j^{in})$$

# Training through the layers...

**Generalized delta rule (output wts)**

$$\begin{aligned}\frac{\partial E}{\partial w_{ij}^{out}} &= \frac{1}{2} \frac{\partial}{\partial w_{ij}^{out}} \sum_i (r_i^{out} - y_i)^2 \\ &= \delta_i^{out} r_j^h\end{aligned}$$

$$E = \frac{1}{2} \sum_i (r_i^{out} - y_i)^2$$

$$\Delta w_{ij} = \varepsilon (g'(h_i) (y_i - r_i^{out}) r_j^{in})$$

with

$$\delta_i^{out} = g^{out \prime}(h_i^h) (r_i^{out} - y_i) \quad \text{delta rule for output weights}$$

# Training through the layers...

**Generalized delta rule (output wts)**

$$\begin{aligned}\frac{\partial E}{\partial w_{ij}^{out}} &= \frac{1}{2} \frac{\partial}{\partial w_{ij}^{out}} \sum_i (r_i^{out} - y_i)^2 \\ &= \delta_i^{out} r_j^h\end{aligned}$$

$$E = \frac{1}{2} \sum_i (r_i^{out} - y_i)^2$$

$$\Delta w_{ij} = \varepsilon (g'(h_i) (y_i - r_i^{out}) r_j^{in})$$

with

$$\delta_i^{out} = g^{out \prime}(h_i^h) (r_i^{out} - y_i) \quad \text{delta rule for output weights}$$

**Hidden layer weights**

# Training through the layers...

## Generalized delta rule (output wts)

$$\begin{aligned}\frac{\partial E}{\partial w_{ij}^{out}} &= \frac{1}{2} \frac{\partial}{\partial w_{ij}^{out}} \sum_i (r_i^{out} - y_i)^2 \\ &= \delta_i^{out} r_j^h\end{aligned}$$

$$E = \frac{1}{2} \sum_i (r_i^{out} - y_i)^2$$

$$\Delta w_{ij} = \varepsilon (g'(h_i) (y_i - r_i^{out}) r_j^{in})$$

with

$$\delta_i^{out} = g^{out \prime}(h_i^h) (r_i^{out} - y_i) \quad \text{delta rule for output weights}$$

## Hidden layer weights

$$\frac{\partial E}{\partial w_{ij}^h} = \frac{1}{2} \frac{\partial}{\partial w_{ij}^h} \sum_i (r_i^{out} - y_i)^2$$

# Training through the layers...

## Generalized delta rule (output wts)

$$\begin{aligned}\frac{\partial E}{\partial w_{ij}^{out}} &= \frac{1}{2} \frac{\partial}{\partial w_{ij}^{out}} \sum_i (r_i^{out} - y_i)^2 \\ &= \delta_i^{out} r_j^h\end{aligned}$$

$$E = \frac{1}{2} \sum_i (r_i^{out} - y_i)^2$$

$$\Delta w_{ij} = \varepsilon (g'(h_i) (y_i - r_i^{out}) r_j^{in})$$

with

$$\delta_i^{out} = g^{out \prime}(h_i^h) (r_i^{out} - y_i) \quad \text{delta rule for output weights}$$

## Hidden layer weights

$$\frac{\partial E}{\partial w_{ij}^h} = \frac{1}{2} \frac{\partial}{\partial w_{ij}^h} \sum_i (r_i^{out} - y_i)^2$$

$$\begin{aligned}\frac{\partial E}{\partial w_{ij}^h} &= \frac{1}{2} \frac{\partial}{\partial w_{ij}^h} \sum_i (g^{out}(\sum_j w_{ij}^{out} g^h(\sum_k w_{jk}^h r_k^{in})) - y_i)^2 \\ &= \delta_i^h r_j^{in}\end{aligned}$$



# Training through the layers...

## Generalized delta rule (output wts)

$$\begin{aligned}\frac{\partial E}{\partial w_{ij}^{out}} &= \frac{1}{2} \frac{\partial}{\partial w_{ij}^{out}} \sum_i (r_i^{out} - y_i)^2 \\ &= \delta_i^{out} r_j^h\end{aligned}$$

$$E = \frac{1}{2} \sum_i (r_i^{out} - y_i)^2$$

$$\Delta w_{ij} = \varepsilon (g'(h_i) (y_i - r_i^{out}) r_j^{in})$$

with

$$\delta_i^{out} = g^{out \prime}(h_i^h) (r_i^{out} - y_i) \quad \text{delta rule for output weights}$$

## Hidden layer weights

$$\frac{\partial E}{\partial w_{ij}^h} = \frac{1}{2} \frac{\partial}{\partial w_{ij}^h} \sum_i (r_i^{out} - y_i)^2$$

$$\begin{aligned}\frac{\partial E}{\partial w_{ij}^h} &= \frac{1}{2} \frac{\partial}{\partial w_{ij}^h} \sum_i (g^{out}(\sum_j w_{ij}^{out} g^h(\sum_k w_{jk}^h r_k^{in})) - y_i)^2 \\ &= \delta_i^h r_j^{in}\end{aligned}$$

with

$$\delta_i^h = g^{h \prime}(h_i^{in}) \sum_k w_{ik}^{out} \delta_k^{out} \quad \begin{array}{l} \text{delta rule for hidden wts} \\ \text{(depends on delta rule for output wts)} \end{array}$$

# Training through the layers...

## Generalized delta rule (output wts)

$$\begin{aligned}\frac{\partial E}{\partial w_{ij}^{out}} &= \frac{1}{2} \frac{\partial}{\partial w_{ij}^{out}} \sum_i (r_i^{out} - y_i)^2 \\ &= \delta_i^{out} r_j^h\end{aligned}$$

$$E = \frac{1}{2} \sum_i (r_i^{out} - y_i)^2$$

$$\Delta w_{ij} = \varepsilon (g'(h_i) (y_i - r_i^{out}) r_j^{in})$$

with

$$\delta_i^{out} = g^{out \prime}(h_i^h) (r_i^{out} - y_i) \quad \text{delta rule for output weights}$$

## Hidden layer weights

$$\frac{\partial E}{\partial w_{ij}^h} = \frac{1}{2} \frac{\partial}{\partial w_{ij}^h} \sum_i (r_i^{out} - y_i)^2$$

$$\begin{aligned}\frac{\partial E}{\partial w_{ij}^h} &= \frac{1}{2} \frac{\partial}{\partial w_{ij}^h} \sum_i (g^{out}(\sum_j w_{ij}^{out} g^h(\sum_k w_{jk}^h r_k^{in})) - y_i)^2 \\ &= \delta_i^h r_j^{in}\end{aligned}$$

with

$$\delta_i^h = g^{h \prime}(h_i^{in}) \sum_k w_{ik}^{out} \delta_k^{out} \quad \text{delta rule for hidden wts}$$

(depends on delta rule for output wts)

# Training through the layers...

**Generalized delta rule (output wts)**

$$\begin{aligned}\frac{\partial E}{\partial w_{ij}^{out}} &= \frac{1}{2} \frac{\partial}{\partial w_{ij}^{out}} \sum_i (r_i^{out} - y_i)^2 \\ &= \delta_i^{out} r_j^h\end{aligned}$$

$$E = \frac{1}{2} \sum_i (r_i^{out} - y_i)^2$$

$$\Delta w_{ij} = \varepsilon (g'(h_i) (y_i - r_i^{out}) r_j^{in})$$

with

$$\delta_i^{out} = (r_i^{out} - y_i) g'(h_i)$$

**Error propagates BACKWARDS through the layers!**

**Hidden layer weights**

$$\frac{\partial E}{\partial w_{ij}^h} = \frac{1}{2} \frac{\partial}{\partial w_{ij}^h} \sum_i (r_i^{out} - y_i)^2$$

$$\begin{aligned}\frac{\partial E}{\partial w_{ij}^h} &= \frac{1}{2} \frac{\partial}{\partial w_{ij}^h} \sum_i (g^{out}(\sum_j w_{ij}^{out} g^h(\sum_k w_{jk}^h r_k^{in})) - y_i)^2 \\ &= \delta_i^h r_j^{in}\end{aligned}$$

with

$$\delta_i^h = g^{h'}(h_i^{in}) \sum_k w_{ik}^{out} \delta_k^{out}$$

**delta rule for hidden wts**

**(depends on delta rule for output wts)**

# Training protocol (gradient descent)



**Training set**

**Test set**

# Training protocol (gradient descent)



**Training set**

**Test set**

1. Train until gradient of error function reaches minimum.
  - A. Batch:** use full training set with each iteration (smooth convergence, but more prone to local minima)
  - B. Online:** use different sample of training set with each iteration (more memory efficient, but messy convergence)

# Training protocol (gradient descent)



**Training set**

**Test set**

1. Train until gradient of error function reaches minimum.
  - A. Batch:** use full training set with each iteration (smooth convergence, but more prone to local minima)
  - B. Online:** use different sample of training set with each iteration (more memory efficient, but messy convergence)

# Training protocol (gradient descent)



**Training set**

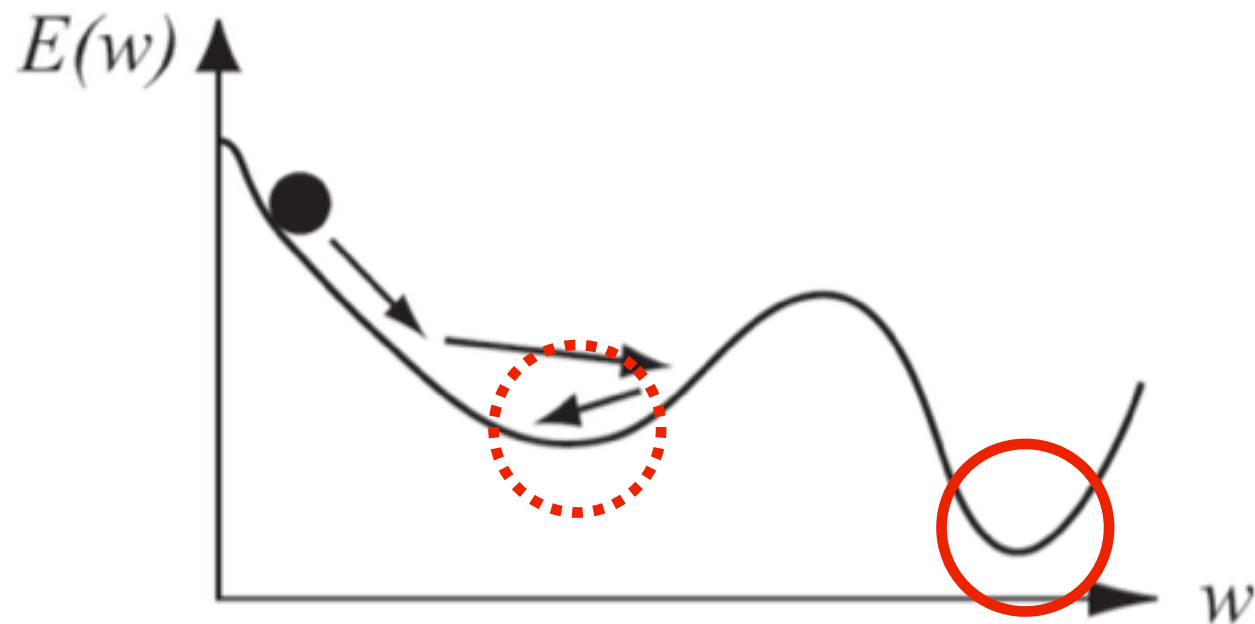
**Test set**

1. Train until gradient of error function reaches minimum.
  - A. Batch:** use full training set with each iteration (smooth convergence, but more prone to local minima)
  - B. Online:** use different sample of training set with each iteration (more memory efficient, but messy convergence)
2. Test generalization of network using previously unseen test data set.

# Caveats



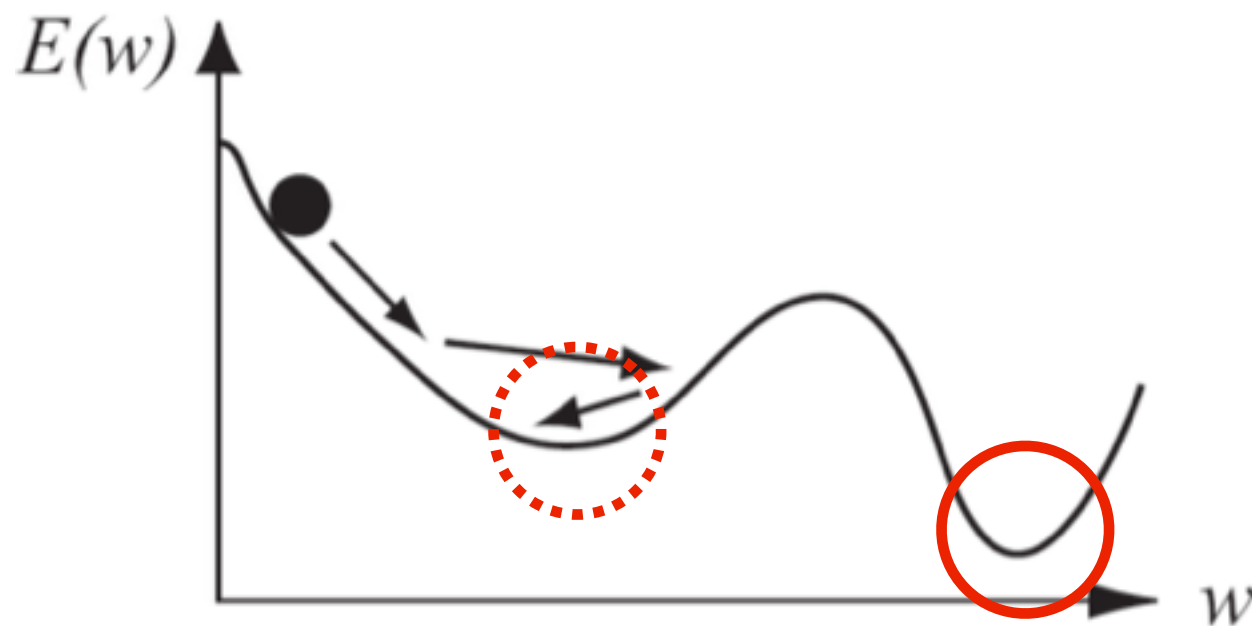
# Caveats



## Local minima

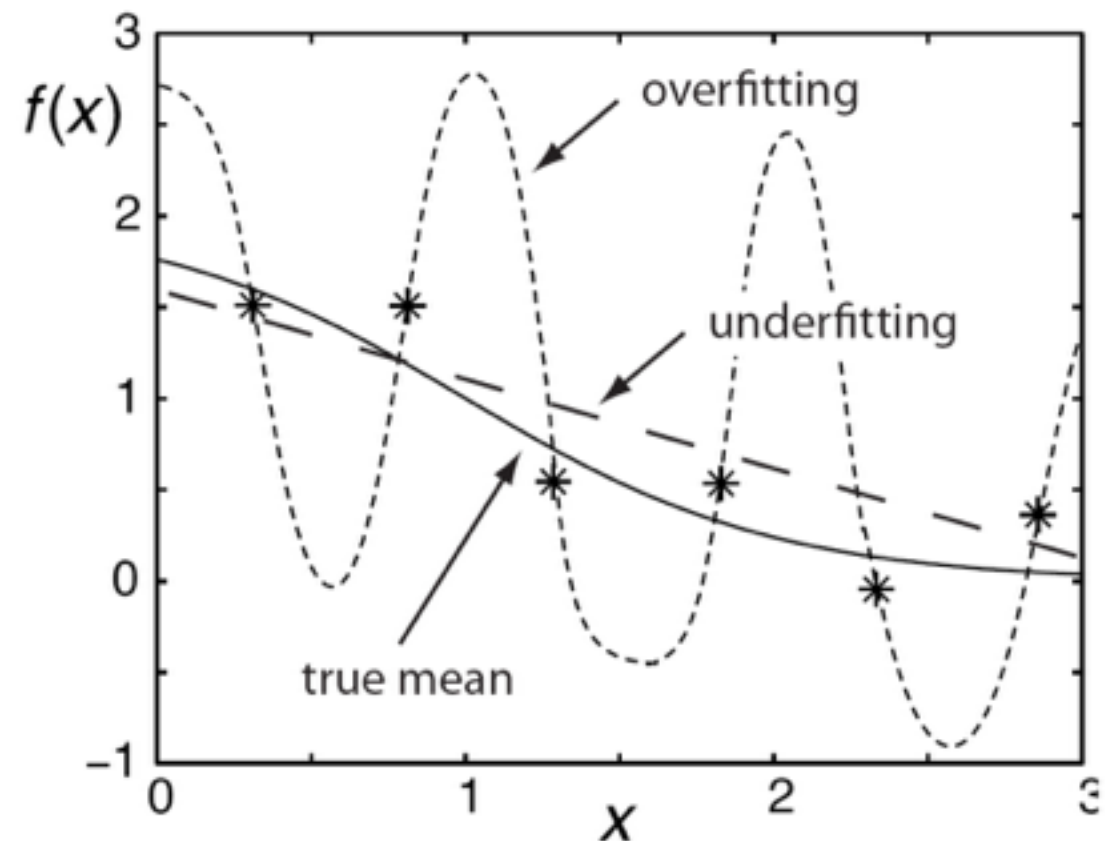
Can use momentum term in weight update to incorporate history of weight changes.

# Caveats



## Local minima

Can use momentum term in weight update to incorporate history of weight changes.



Trappenberg 2010

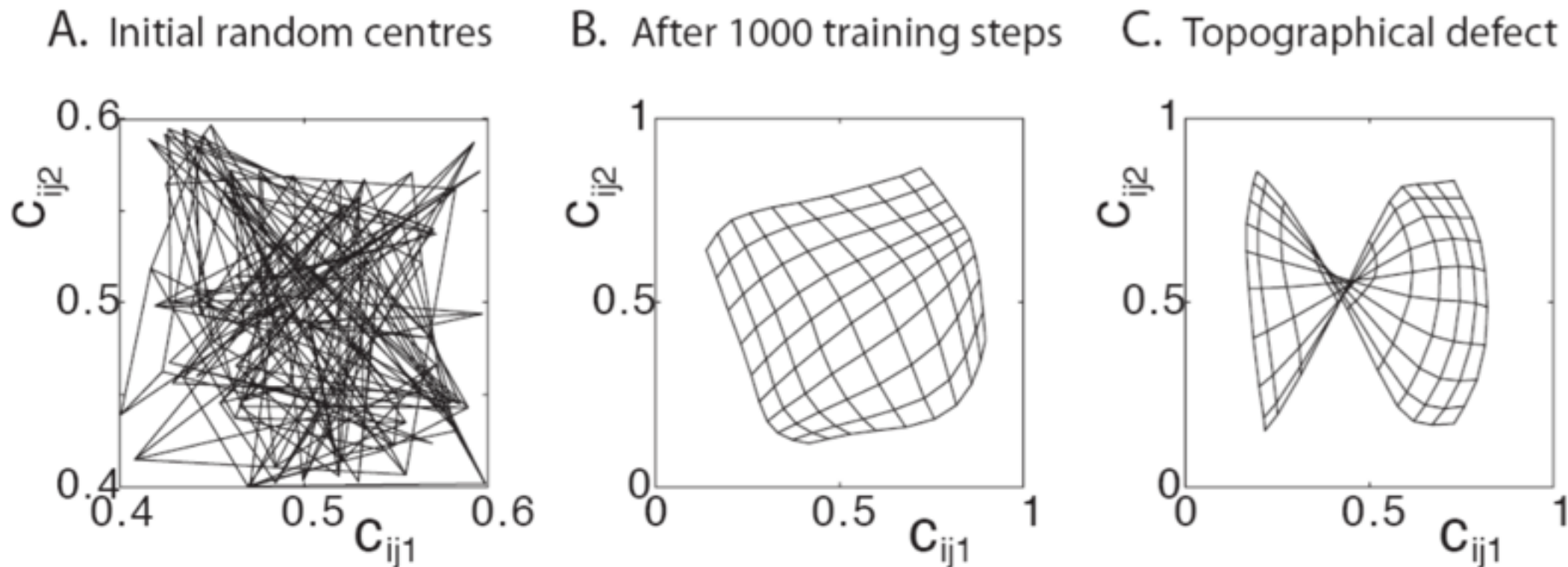
## Overfitting

Use heuristics to determine appropriate number of nodes for solving particular problem. Can usually use  $2 \times$  number of nodes for training set.

# Recurrent networks are a whole new game!

# Recurrent networks are a whole new game!

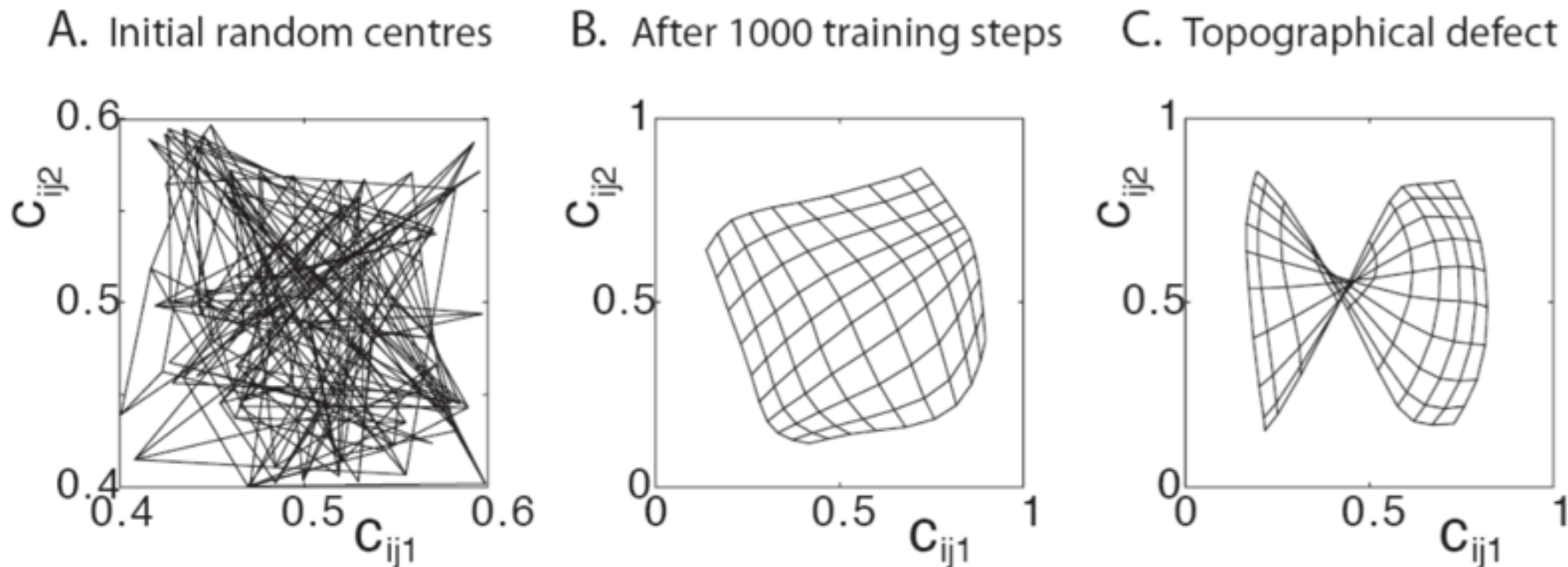
*but I'll spare you*



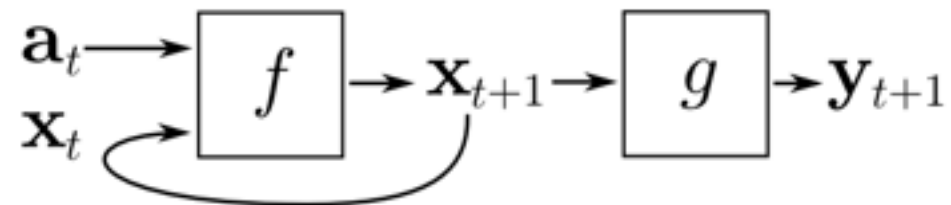
Trappenberg 2010

# Recurrent networks are a whole new game!

*but I'll spare you*



**Backprop. through time  
(BPTT)**



Trappenberg 2010

↓ unfold through time ↓

