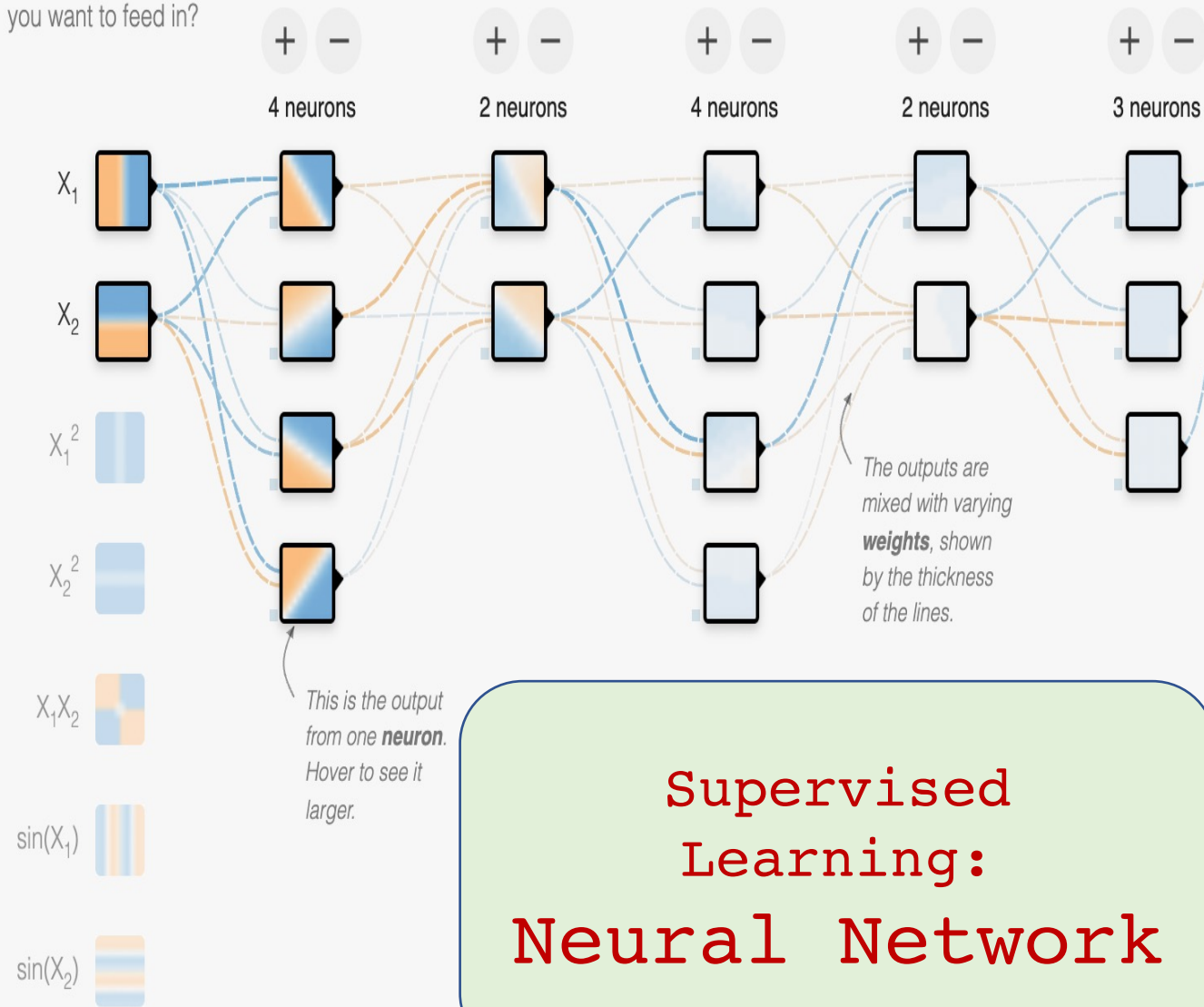


FEATURES

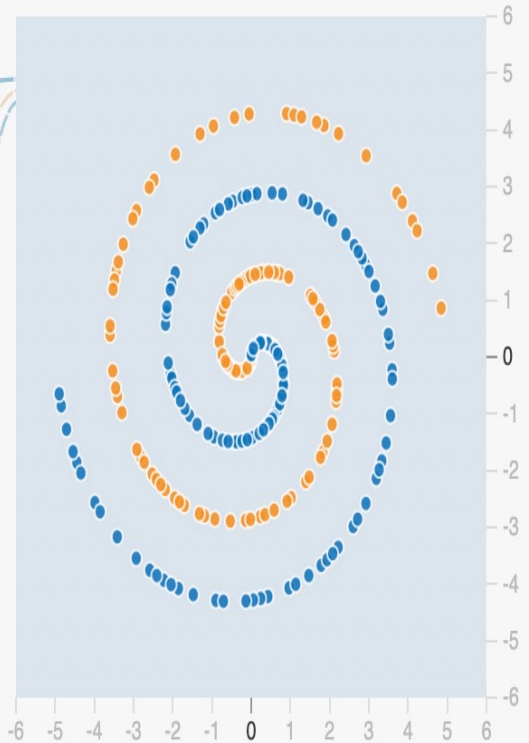
Which properties do you want to feed in?



+ - 5 HIDDEN LAYERS

OUTPUT

Test loss 0.516  
Training loss 0.501



Supervised Learning:  
Neural Network



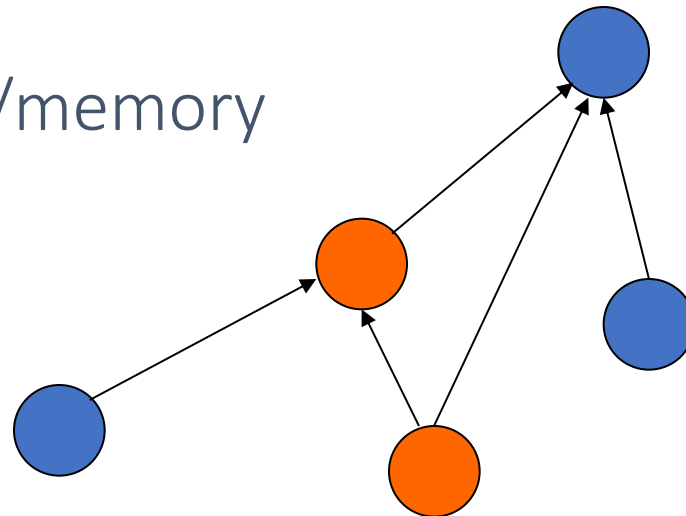
☐ Show test data ☐ Discretize output

# Roadmap

- Overview of neural networks
- Single-layer NN: Perceptron
  - Training, delta rule, etc.
  - Limitation
- Multi-layered NN: Multilayered Perceptron (MLP)
  - Chain rule
  - Back-propagation algorithm
  - Other issues
- Take-home messages

# Multilayered Perceptron (MLP): A Very First Neural Networks

- Networks of processing units (neurons) with connections (synapses) between them
- Large number of neurons:  $10^{10}$
- Large connectivity:  $10^5$
- Parallel processing
- Distributed computation/memory
- Robust to noise, failures



# Understanding the Brain

- Levels of analysis (Marr, 1982)
  1. Computational theory
  2. Representation and algorithm
  3. Hardware implementation
- Reverse engineering: From hardware to theory
- Parallel processing:
  - SIMD (Single instruction, multiple data) VS
  - MIMD (Multiple instruction, multiple data)
- Neural net: SIMD with modifiable local memory
- Learning: Update by training/experience ( $E$ )

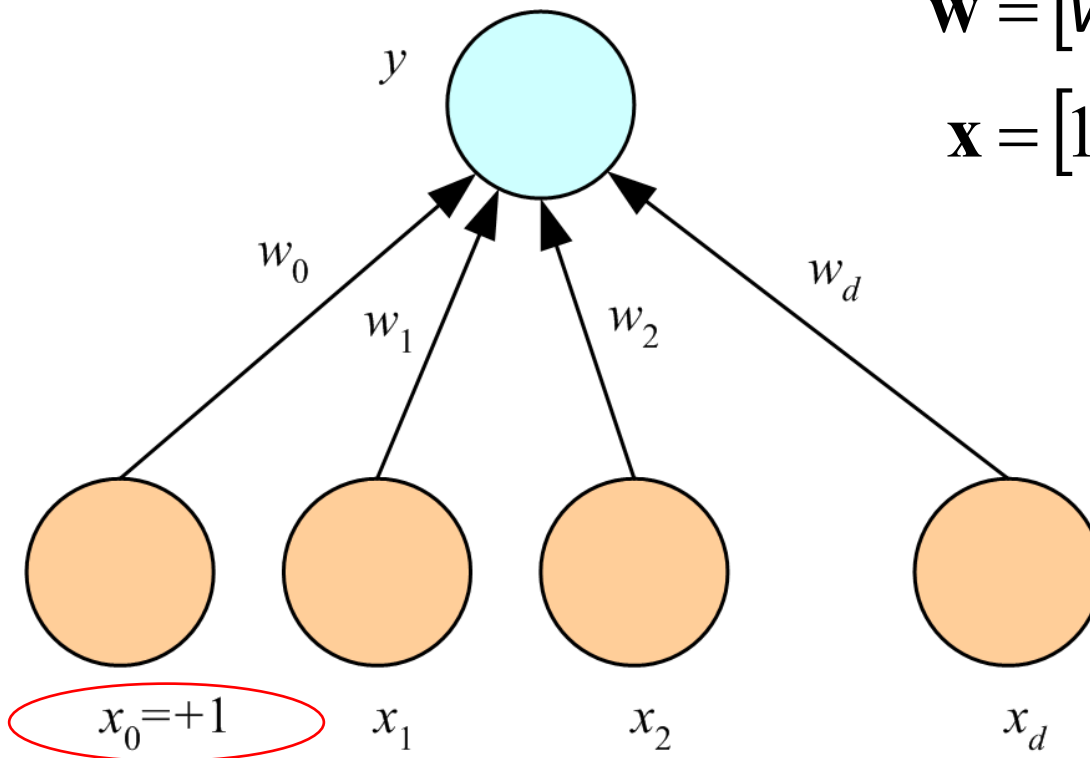
# A Perceptron

$$y = \sum_{j=1}^d w_j x_j + w_0 = \mathbf{w}^T \mathbf{x}$$

$$\mathbf{w} = [w_0, w_1, \dots, w_d]^T$$

$$\mathbf{x} = [1, x_1, \dots, x_d]^T$$

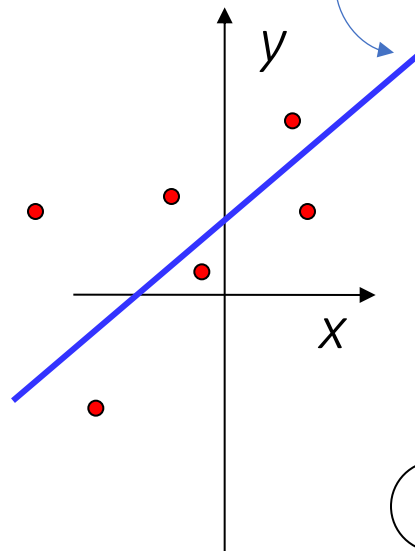
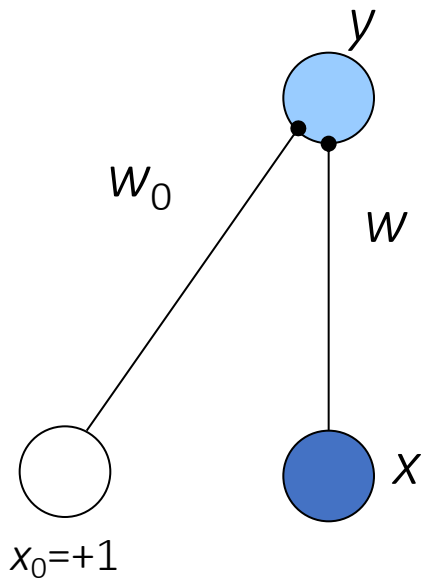
(Rosenblatt, 1962)



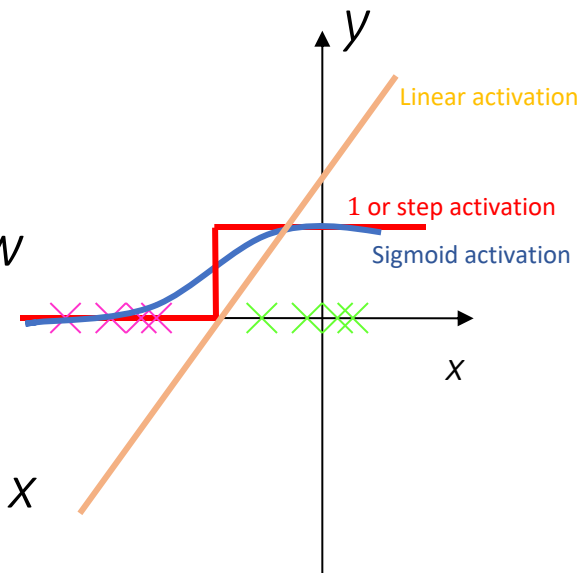
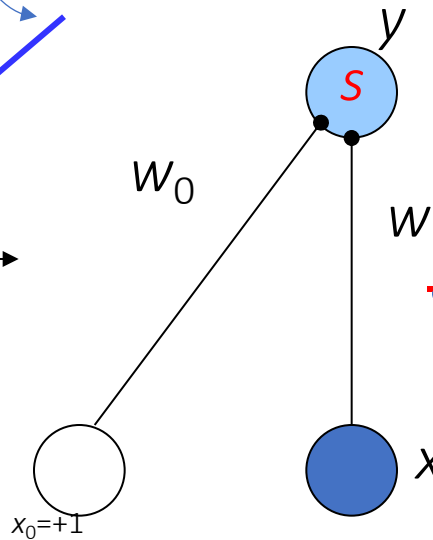
# What a Perceptron Does?

- Regression:  $y = wx + w_0$

Line fitting



- Classification:  $y = 1(w x + w_0 > 0)$



Activation Functions:

1 or step activation:  $y = 1(o) = \begin{cases} 1 & \text{if } o > 0 \\ 0 & \text{otherwise} \end{cases}$

Sigmoid activation:  $y = \text{sigmoid}(o) = \frac{1}{1 + e^{-w^T x}}$

Linear activation:  $y = o = w^T x$

Neural Networks

# K Perceptrons for K Outputs

Classification:

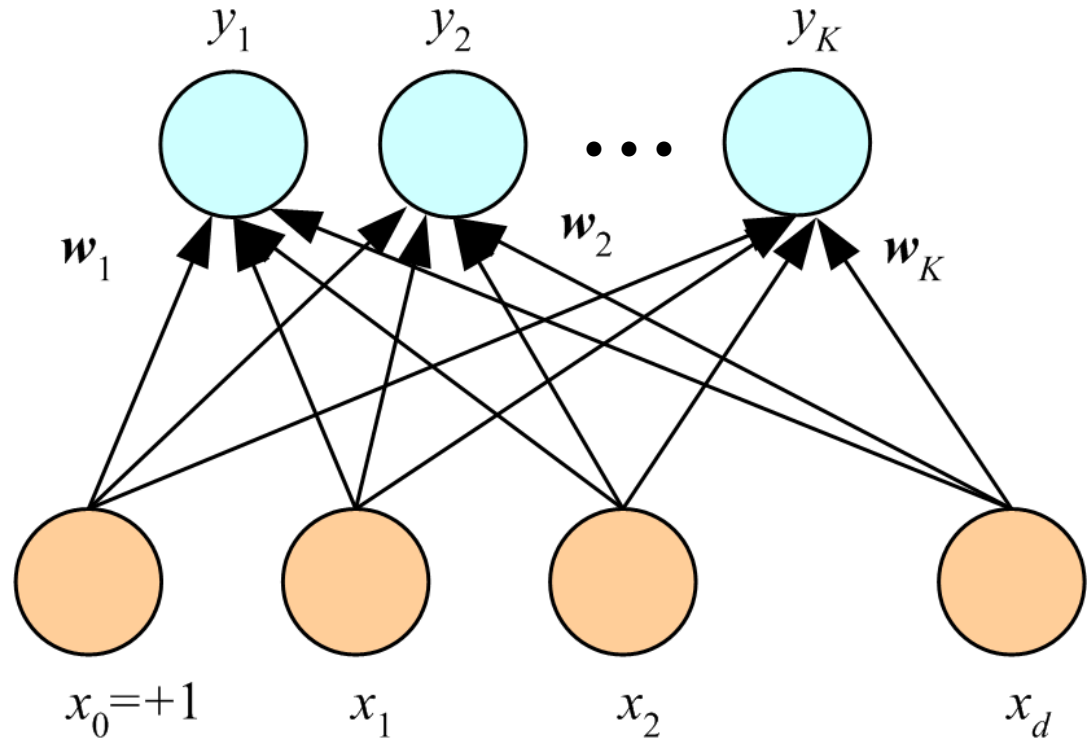
$$o_i = \mathbf{w}_i^T \mathbf{x}$$

$$y_i = \frac{\exp o_i}{\sum_k \exp o_k}$$

choose  $C_i$

$$\text{if } y_i = \max_k y_k$$

Softmax  
function



Regression: 
$$y_i = \sum_{j=1}^d w_{ij} x_j + w_{i0} = \mathbf{w}_i^T \mathbf{x}$$

$$\mathbf{y} = \mathbf{W}\mathbf{x}$$

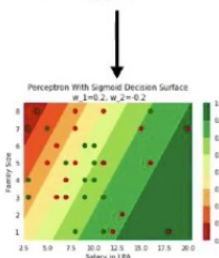
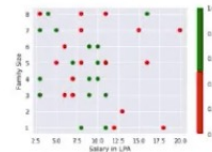
# Supervised Learning with NN



Real inputs  
 $\in \mathbb{R}$



Classification  
/Regression



Objective function



$$Loss = \sum_{i=1}^n (y - \hat{y})^2$$



$$w = w + \Delta w$$

$$b = b + \Delta b$$

OR

$$w_{t+1} = w_t - \eta \Delta w_t$$

$$b_{t+1} = b_t - \eta \Delta b_t$$



$$RMSE$$

$$Accuracy = \frac{\text{Number of correct predictions}}{\text{Total number of predictions}}$$



# Learning/Training

- Online (instances seen one by one) learning
  - No need to store the whole sample as in batch learning
  - Problem may change in time
- Stochastic gradient-descent (**SGD**): Update after a single instance
- Generic update rule (delta rule):

$$\Delta w_{ij}^t = \eta (r_i^t - y_i^t) x_j^t$$

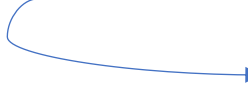
$t$  denotes current training instance

Weight\_update = learning\_rate x (desired\_output - predicted\_output) x input

# Training a Perceptron: Regression Task

- Regression (Linear output): Linear activation:  $y = o = \mathbf{w}^T \mathbf{x}$

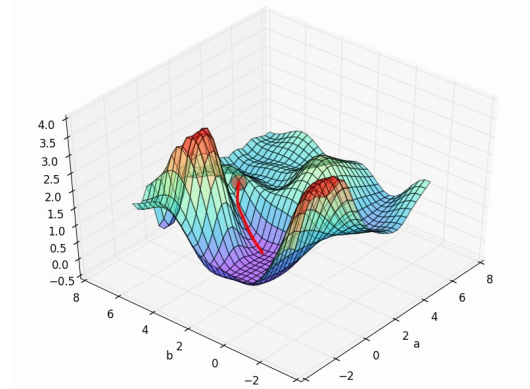
Half squared error/loss


$$E^t(\mathbf{w} | \mathbf{x}^t, r^t) = \frac{1}{2} (r^t - y^t)^2 = \frac{1}{2} [r^t - (\mathbf{w}^T \mathbf{x}^t)]^2$$

$$\Delta w_j^t = \eta (r^t - y^t) x_j^t$$


$$\Delta w_j^t = -\eta \frac{\partial E^t}{\partial w_j^t} \quad \leftarrow \dots \text{Gradient descent}$$

# Gradient Descent



- $E(\mathbf{w} | X)$  is error with parameters  $\mathbf{w}$  on sample  $X$   
 $\mathbf{w}^* = \arg \min_{\mathbf{w}} E(\mathbf{w} | X)$

- Gradient

Partial derivatives

Gradient in n-D space, i.e., vector of Partial derivatives; del or nabla operator

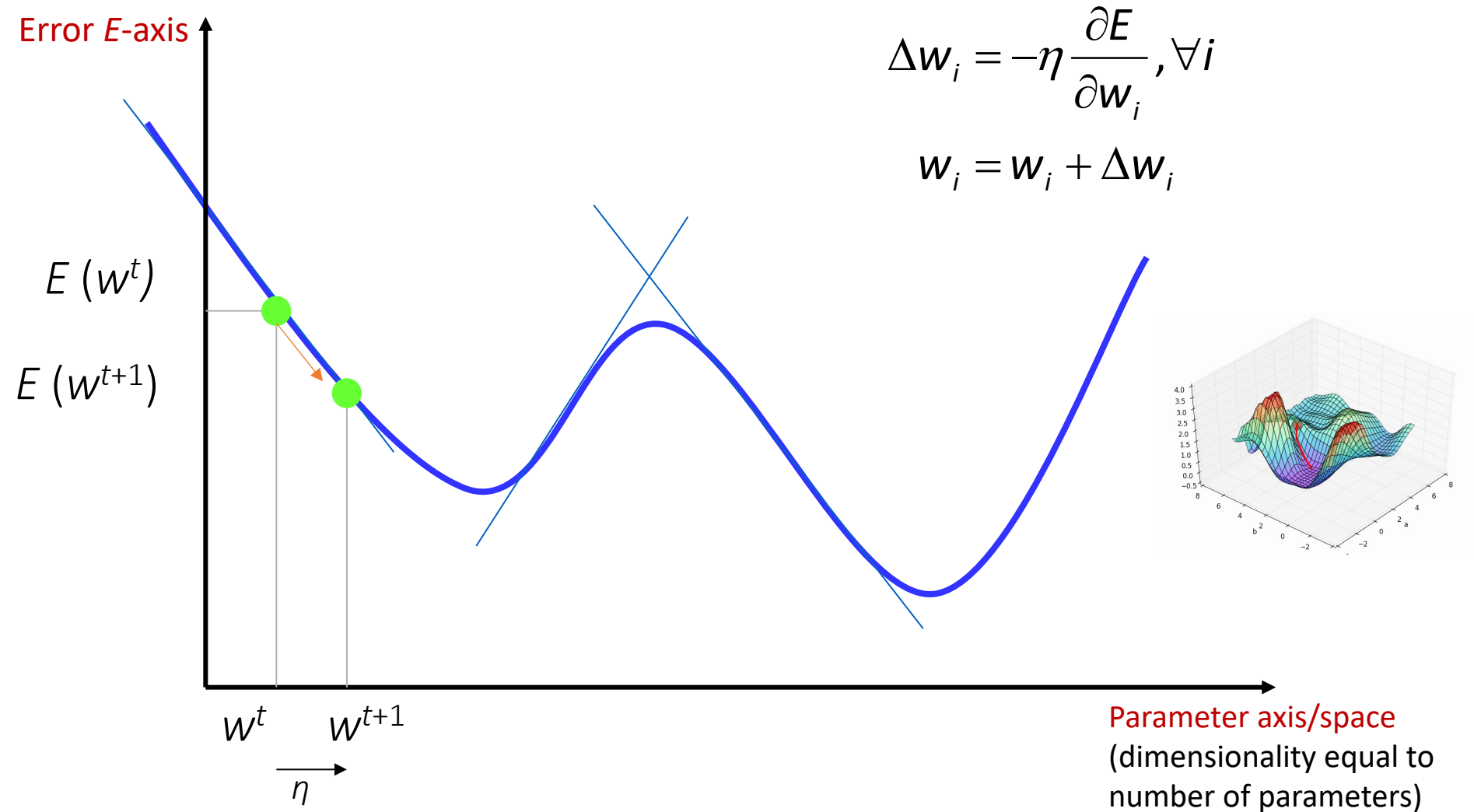
$$\nabla_{\mathbf{w}} E = \left[ \frac{\partial E}{\partial w_1}, \frac{\partial E}{\partial w_2}, \dots, \frac{\partial E}{\partial w_d} \right]^T$$

Parameter space

- Gradient-descent:

Starts from random  $\mathbf{w}$  and updates  $\mathbf{w}$  iteratively in the **negative direction of gradient**. (Note that gradient is a vector pointing at the greatest increase of a function, negative gradient is a vector pointing at the greatest decrease of a function!)

# Gradient Descent



# Training a Perceptron: Classification Task

- $K=1$  Single sigmoid output

Sigmoid activation:  $y = \text{sigmoid}(o) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$

$$y^t = \text{sigmoid}(\mathbf{w}^T \mathbf{x}^t)$$

$$E^t(\mathbf{w} | \mathbf{x}^t, \mathbf{r}^t) = -r^t \log y^t - (1 - r^t) \log (1 - y^t)$$

Logistic loss

$$\Delta \mathbf{w}_j^t = \eta (r^t - y^t) \mathbf{x}_j^t$$

Weight\_update = learning\_rate x (desired\_output - predicted\_output) x input

- $K \geq 2$  Softmax outputs

$$y^t = \frac{\exp \mathbf{w}_i^T \mathbf{x}^t}{\sum_k \exp \mathbf{w}_k^T \mathbf{x}^t}$$

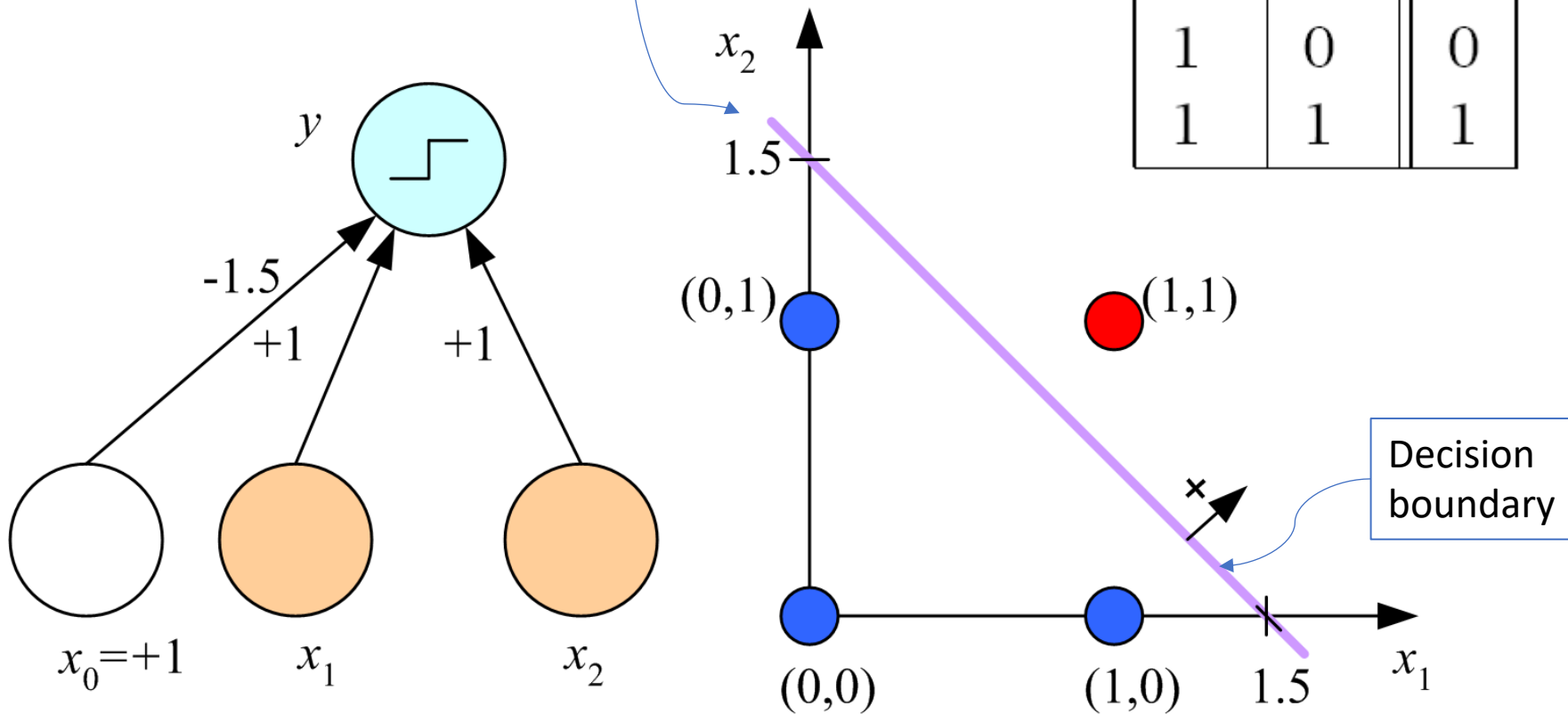
$$E^t(\{\mathbf{w}_i\}_i | \mathbf{x}^t, \mathbf{r}^t) = -\sum_i r_i^t \log y_i^t$$

$$\Delta \mathbf{w}_{ij}^t = \eta (r_i^t - y_i^t) \mathbf{x}_j^t$$

Weight\_update = learning\_rate x (desired\_output - predicted\_output) x input

# Learning Boolean AND

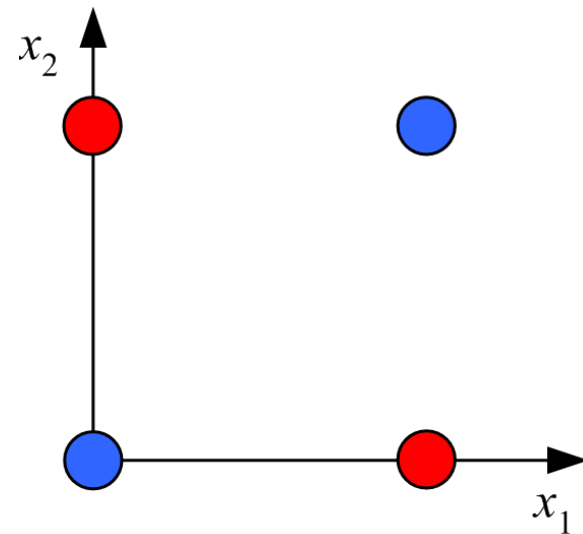
Linear separable!



# How about learning XOR?

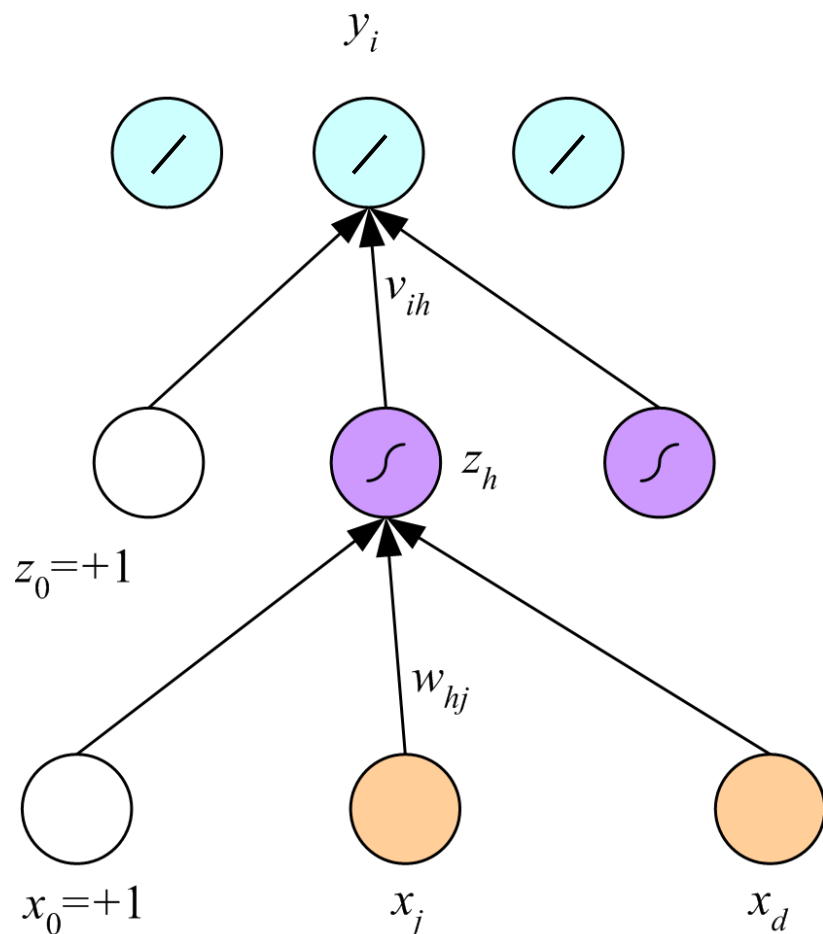
Linear non-separable!

$x_1$	$x_2$	$r$
0	0	0
0	1	1
1	0	1
1	1	0



Any straight line (one line only) to perfectly separate the four pts?

# Empowering perceptron with more layers: Multilayer Perceptrons (MLP)



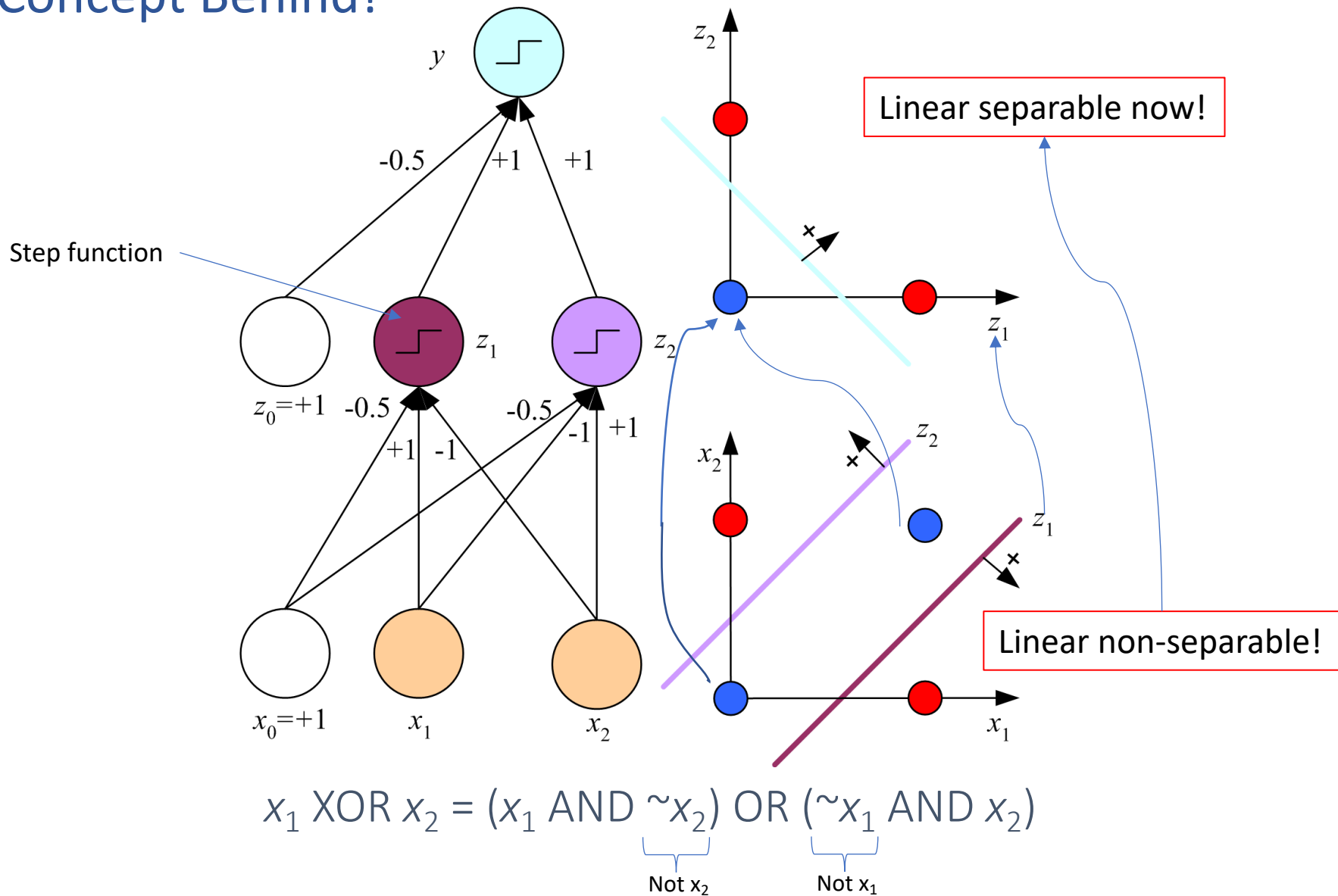
$$y_i = \mathbf{v}_i^T \mathbf{z} = \sum_{h=1}^H v_{ih} z_h + v_{i0}$$

$$z_h = \text{sigmoid}(\mathbf{w}_h^T \mathbf{x})$$
$$= \frac{1}{1 + \exp\left[-\left(\sum_{j=1}^d w_{hj} x_j + w_{h0}\right)\right]}$$

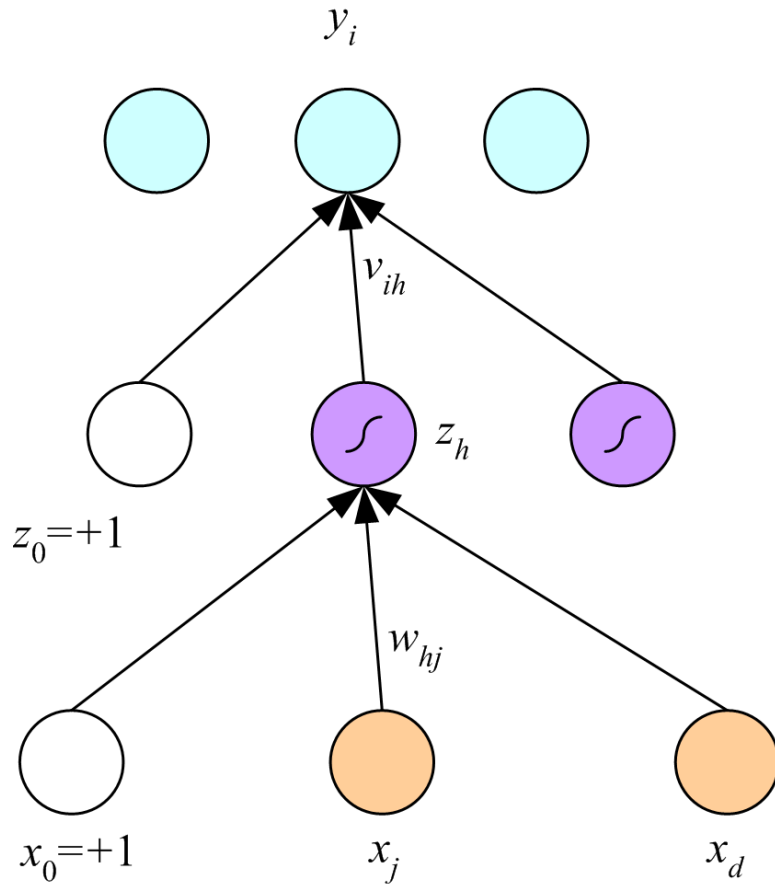
(Rumelhart et al., 1986)



# Concept Behind!



# Backpropagation (BP) Algorithm for Training MLP



$$y_i = \mathbf{v}_i^T \mathbf{z} = \sum_{h=1}^H v_{ih} z_h + v_{i0}$$

$$z_h = \text{sigmoid}(\mathbf{w}_h^T \mathbf{x})$$

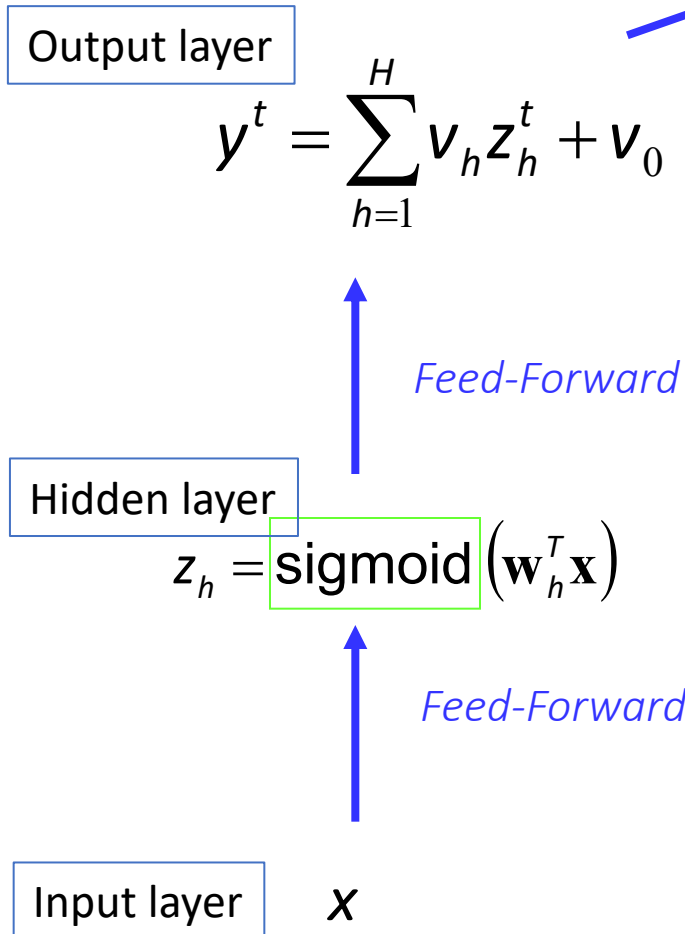
$$= \frac{1}{1 + \exp\left[-\left(\sum_{j=1}^d w_{hj} x_j + w_{h0}\right)\right]}$$

$$\frac{\partial E}{\partial w_{hj}} = \frac{\partial E}{\partial y_i} \frac{\partial y_i}{\partial z_h} \frac{\partial z_h}{\partial w_{hj}}$$

Very famous “[Chain rule](#)” in calculus!

# Regression

with (K=)1 Output



*Error (to backpropagate):*

$$E(\mathbf{W}, \mathbf{v} | \mathcal{X}) = \frac{1}{2} \sum_t (r^t - y^t)^2$$

$$\Delta v_h = \eta \sum_t (r^t - y^t) z_h^t$$

*Backward*

$$\begin{aligned} \Delta \mathbf{w}_{hj} &= -\eta \frac{\partial E}{\partial \mathbf{w}_{hj}} \\ &= -\eta \sum_t \frac{\partial E}{\partial y^t} \frac{\partial y^t}{\partial z_h^t} \frac{\partial z_h^t}{\partial \mathbf{w}_{hj}} \\ &= -\eta \sum_t -(r^t - y^t) v_h z_h^t (1 - z_h^t) x_j^t \\ &= \eta \sum_t (r^t - y^t) v_h z_h^t (1 - z_h^t) x_j^t \end{aligned}$$

# Regression with Multiple ( $K \geq 2$ ) Outputs

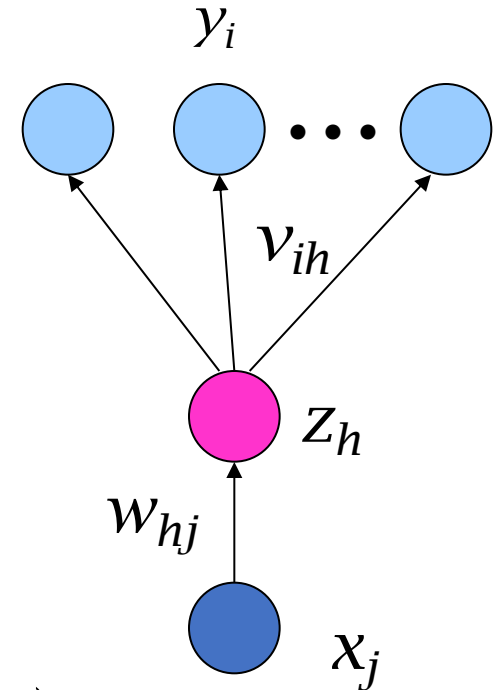
$$E(\mathbf{W}, \mathbf{V} | \mathcal{X}) = \frac{1}{2} \sum_t \sum_i (r_i^t - y_i^t)^2$$

$$y_i^t = \sum_{h=1}^H v_{ih} z_h^t + v_{i0}$$

$$\Delta v_{ih} = \eta \sum_t (r_i^t - y_i^t) z_h^t$$

$$\Delta w_{hj} = \eta \sum_t \left[ \underbrace{\sum_i (r_i^t - y_i^t) v_{ih}}_{\text{Aggregated errors to backpropagate!}} \right] z_h^t (1 - z_h^t) x_j^t$$

Aggregated errors to backpropagate!



# Back Propagation Algorithm

Initialize all  $v_{ih}$  and  $w_{hj}$  to  $\text{rand}(-0.01, 0.01)$

Repeat

For all  $(\mathbf{x}^t, r^t) \in \mathcal{X}$  in random order

For  $h = 1, \dots, H$  /\* for all hidden neurons \*/

$$z_h \leftarrow \text{sigmoid}(\mathbf{w}_h^T \mathbf{x}^t)$$

For  $i = 1, \dots, K$  /\* for all output neurons \*/

$$y_i = \mathbf{v}_i^T \mathbf{z}$$

For  $i = 1, \dots, K$  /\* Compute updates for all output neurons \*/

$$\Delta \mathbf{v}_i = \eta(r_i^t - y_i^t) \mathbf{z}$$

For  $h = 1, \dots, H$  /\* Compute updates for all hidden neurons \*/

$$\Delta \mathbf{w}_h = \eta \left( \sum_i (r_i^t - y_i^t) v_{ih} \right) z_h (1 - z_h) \mathbf{x}^t$$

For  $i = 1, \dots, K$  /\* Update parameters for all output neurons \*/

$$\mathbf{v}_i \leftarrow \mathbf{v}_i + \Delta \mathbf{v}_i$$

For  $h = 1, \dots, H$  /\* Update parameters for all hidden neurons \*/

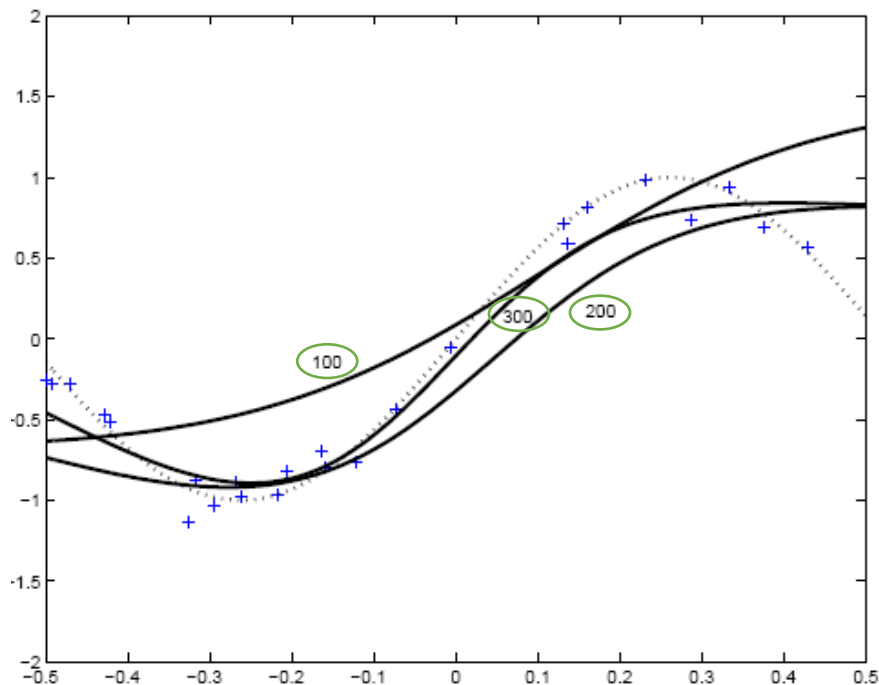
$$\mathbf{w}_h \leftarrow \mathbf{w}_h + \Delta \mathbf{w}_h$$

Until convergence

Feedforward

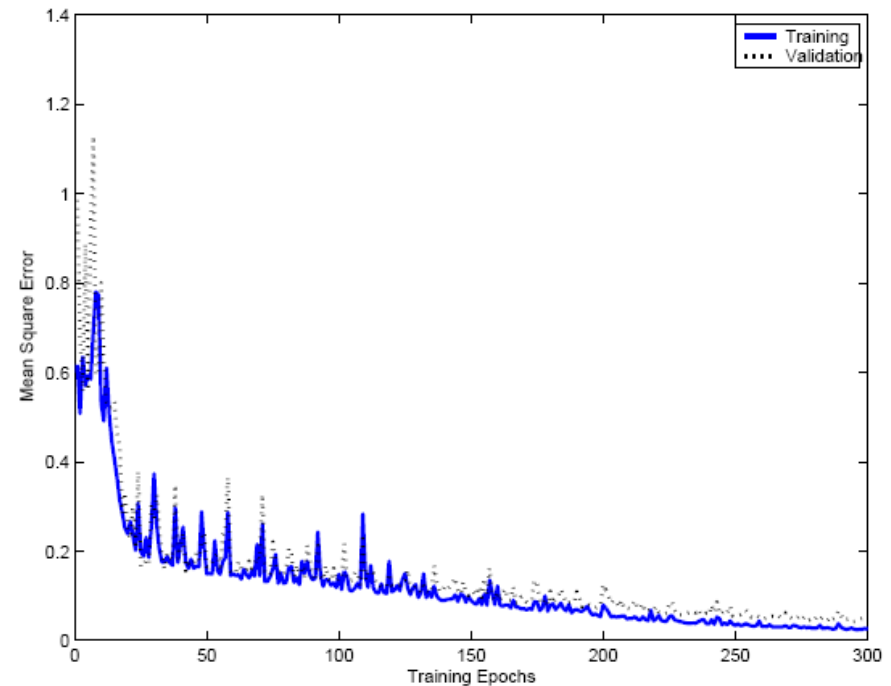
Feedback  
(Back-propagation)

# Learning Curve and Performance



**Figure 11.8** Sample training data shown as '+', where  $x^t \sim U(-0.5, 0.5)$ , and  $y^t = f(x^t) + \mathcal{N}(0, 0.1)$ .  $f(x) = \sin(6x)$  is shown by a dashed line. The evolution of the fit of an MLP with two hidden units after 100, 200, and 300 epochs is drawn.

Learning curve...error decreasing



**Figure 11.9** The mean square error on training and validation sets as a function of training epochs.

Special test sets

# Qualitative Description of Backpropagation Learning

- Iteratively process a set of training tuples & compare the network's prediction with the actual known target value
- For each training tuple, the weights are modified to **minimize the mean squared error** between the network's prediction and the actual target value
- Modifications are made in the “**backwards**” direction: from the output layer, through each hidden layer down to the first hidden layer, hence “**backpropagation**”
- Steps
  - Initialize weights to small random numbers, associated with biases
  - Propagate the inputs forward (by applying activation function)
  - Backpropagate the error (by updating weights and biases)
  - Terminating condition (when error is very small, etc.)

# Multiple Hidden Layers

- MLP with one hidden layer is a **universal approximator** (Hornik et al., 1989), but using more layers may lead to simpler networks

$$z_{1h} = \text{sigmoid}(\mathbf{w}_{1h}^T \mathbf{x}) = \text{sigmoid}\left(\sum_{j=1}^d w_{1hj} x_j + w_{1h0}\right), h = 1, \dots, H_1$$

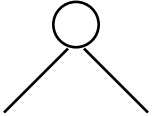
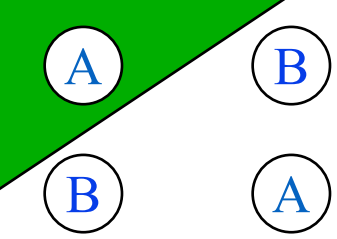
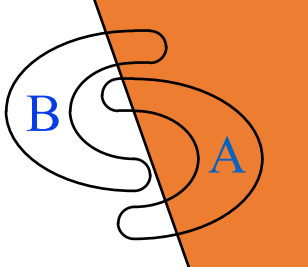

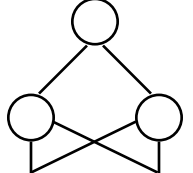
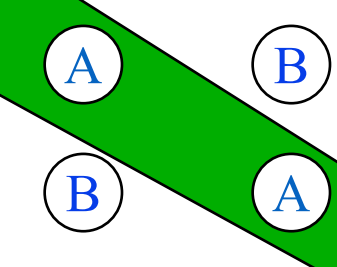
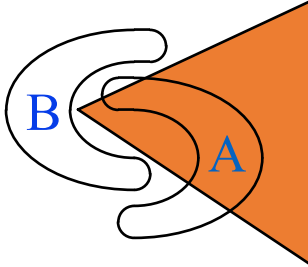
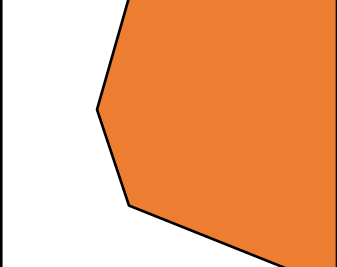
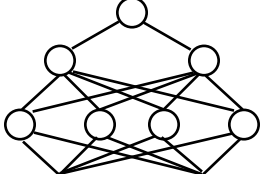
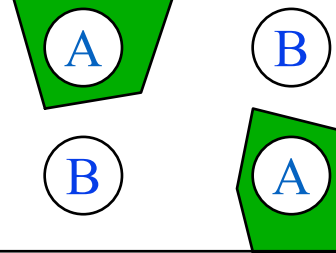
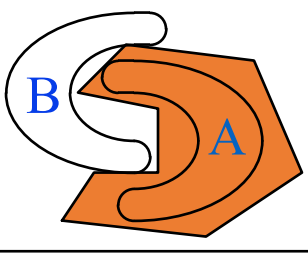
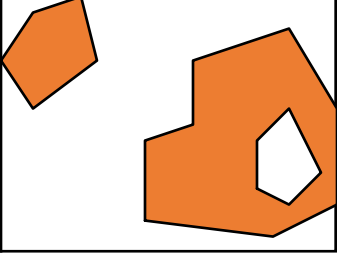
$$z_{2l} = \text{sigmoid}(\mathbf{w}_{2l}^T \mathbf{z}_1) = \text{sigmoid}\left(\sum_{h=1}^{H_1} w_{2lh} z_{1h} + w_{2l0}\right), l = 1, \dots, H_2$$

$$y = \mathbf{v}^T \mathbf{z}_2 = \sum_{l=1}^{H_2} v_l z_{2l} + v_0$$



# Different non-linearly separable problems

Neural Networks – An Introduction Dr. Andrew Hunter

<i>Structure</i>	<i>Types of Decision Regions</i>	<i>Exclusive-OR Problem</i>	<i>Classes with Meshed regions</i>	<i>Most General Region Shapes</i>
<i>Single-Layer</i> 	<i>Half Plane Bounded By Hyperplane</i>			
<i>Two-Layer</i> 	<i>Convex Open Or Closed Regions</i>			
<i>Three-Layer</i> 	<i>Arbitrary (Complexity Limited by No. of Nodes)</i>			

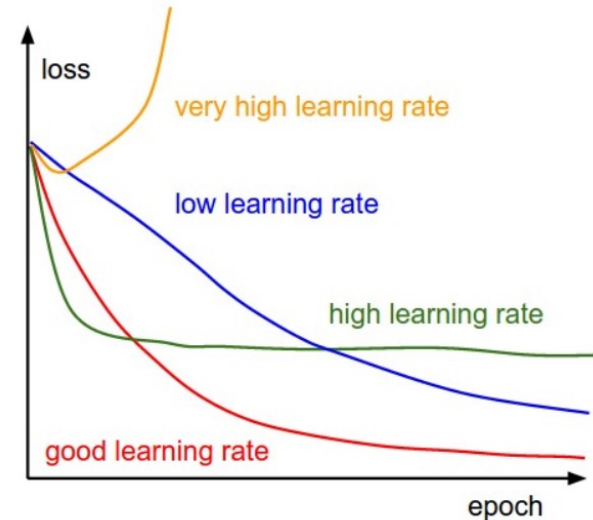
# Improving Convergence

- Momentum

$$\Delta w_i^t = -\eta \frac{\partial E^t}{\partial w_i} + \alpha \Delta w_i^{t-1}$$

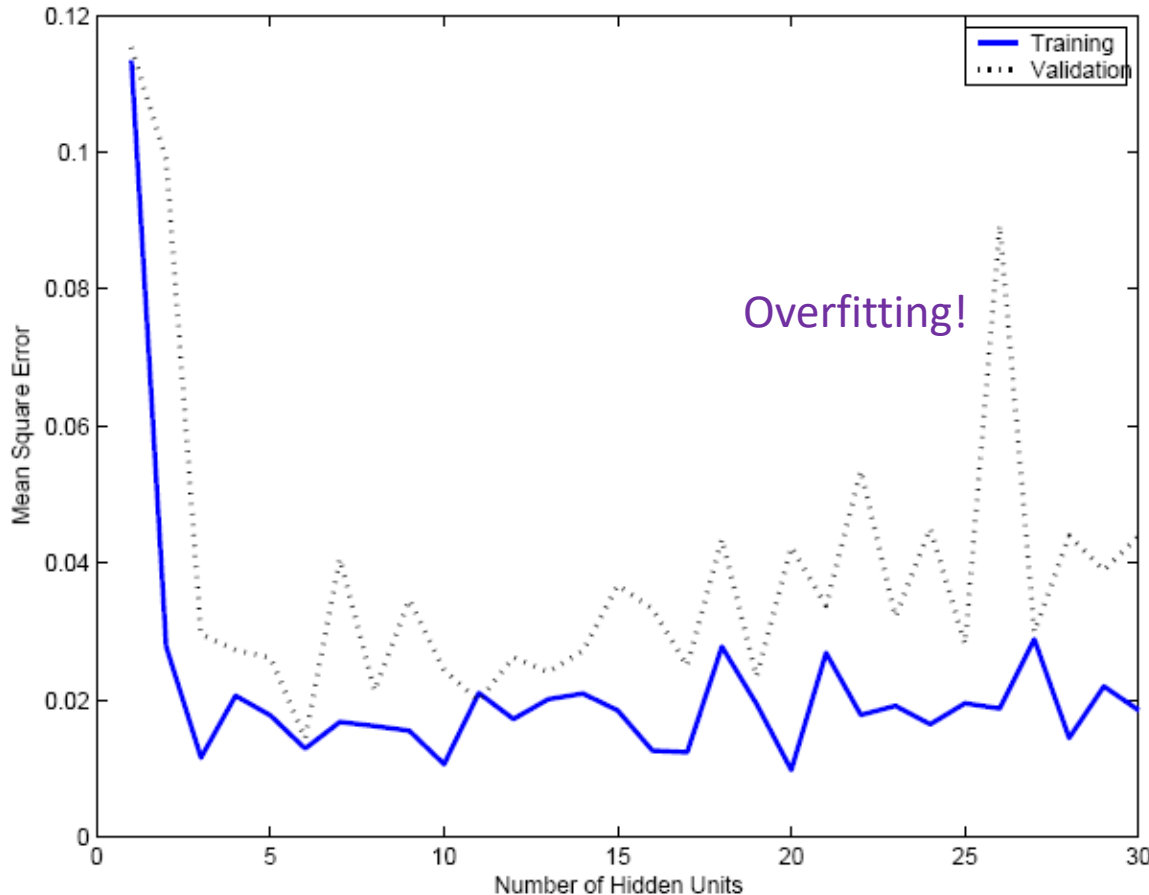
- Adaptive learning rate

$$\Delta \eta = \begin{cases} +a & \text{if } E^{t+\tau} < E^t \\ -b\eta & \text{otherwise} \end{cases}$$



Gradient descent with different learning rates. [Source](#)

# Overfitting (with more complex NN)



Number of weights:

$H(d+1)+(H+1)K$  for MLP  
with  $d$  inputs,  $H$  hidden  
units (1 hidden layer  
only), and  $K$  outputs

Space complexity:

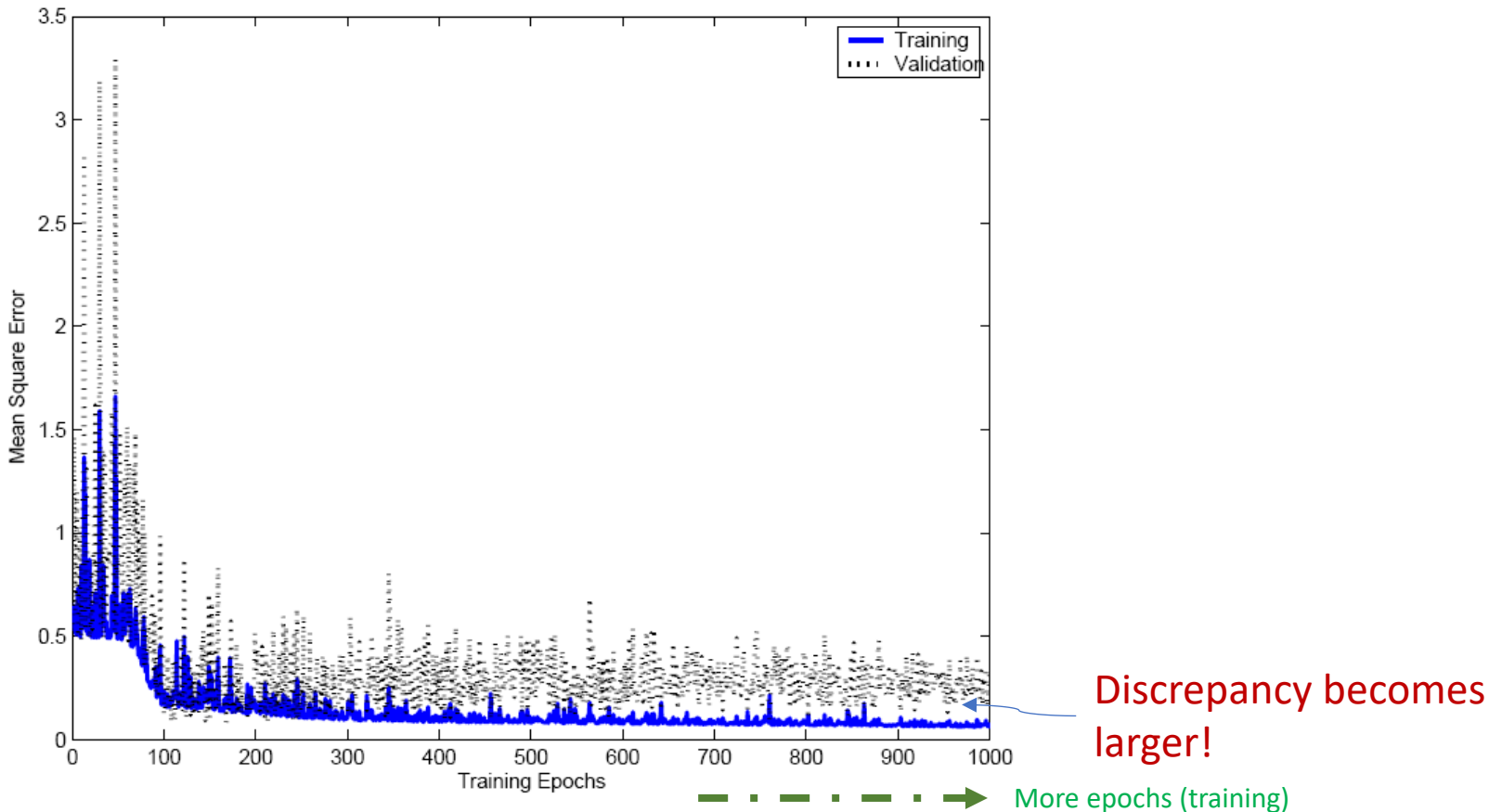
$O(H \cdot (d+K))$  for MLP

Time complexity:

$O(e \cdot H \cdot (d+K))$  where  $e$  is  
the number of training  
epochs

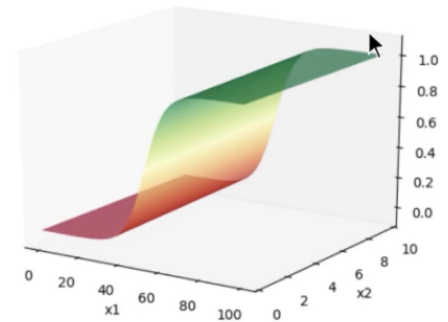
**Figure 11.12** As complexity increases, training error is fixed but the validation error starts to increase and the network starts to overfit.

# Overfitting (with Overtraining)



**Figure 11.13** As training continues, the validation error starts to increase and the network starts to overfit.

# Loss function of sigmoid modified linear discriminant



- General loss function

$$loss = \sum_i (y_i - \hat{y}_i)^2$$

i.e. the sum of the squared difference between the true output  $y_i$  and the predicted output  $\hat{y}_i$

- An example:

$x_1$	$x_2$	$y$	$\hat{y}$
1	1	0.5	0.6
2	1	0.8	0.7
1	2	0.2	0.2
2	2	0.9	0.5

$$\hat{y} = \frac{1}{1 + e^{-(wx+b)}}$$

$$Loss = \sum_{i=1}^4 (y - \hat{y})^2 = 0.18$$

↗  
y could be binary, i.e., 0 or 1.

# Gradient Descent Learning

## Gradient Descent Rule

- The direction  $u$  that we intend to move in should be at  $180^\circ$  w.r.t. the gradient.
- In other words, move in a direction opposite to the gradient.



$$w = w + \Delta w$$

$$b = b + \Delta b$$

OR

$$w_{t+1} = w_t - \eta \Delta w_t$$

$$b_{t+1} = b_t - \eta \Delta b_t$$

## Parameter Update Rule

$$w_{t+1} = w_t - \eta \Delta w_t$$

$$b_{t+1} = b_t - \eta \Delta b_t$$

$$\text{where } \Delta w_t = \frac{\partial \mathcal{L}(w,b)}{\partial w} \text{ at } w=w_t, b=b_t, \Delta b_t = \frac{\partial \mathcal{L}(w,b)}{\partial b} \text{ at } w=w_t, b=b_t$$

## Learning Algorithm

**Initialise**  $w, b$

**Iterate over data:**

*compute*  $\hat{y}$

*compute*  $\mathcal{L}(w, b)$

$$w_{t+1} = w_t - \eta \Delta w_t$$

$$b_{t+1} = b_t - \eta \Delta b_t$$

**till satisfied**

# Gradient Descent Optimization

## Derivative of Sigmoid Function

$$\nabla w = \frac{\partial}{\partial w} \left[ \frac{1}{2} * (f(x) - y)^2 \right]$$

$$= \frac{1}{2} * [2 * (f(x) - y) * \frac{\partial}{\partial w} (f(x) - y)]$$

$$= (f(x) - y) * \frac{\partial}{\partial w} (f(x))$$

$$= (f(x) - y) * \frac{\partial}{\partial w} \left( \frac{1}{1 + e^{-(wx+b)}} \right)$$

$$= (f(x) - y) * f(x) * (1 - f(x)) * x$$

$$\frac{\partial}{\partial w} \left( \frac{1}{1 + e^{-(wx+b)}} \right)$$

$$= \frac{-1}{(1 + e^{-(wx+b)})^2} \frac{\partial}{\partial w} (e^{-(wx+b)})$$

$$= \frac{-1}{(1 + e^{-(wx+b)})^2} * (e^{-(wx+b)}) \frac{\partial}{\partial w} (-(wx + b))$$

$$= \frac{-1}{(1 + e^{-(wx+b)})} * \frac{e^{-(wx+b)}}{(1 + e^{-(wx+b)})} * (-x)$$

$$= \frac{1}{(1 + e^{-(wx+b)})} * \frac{e^{-(wx+b)}}{(1 + e^{-(wx+b)})} * (x)$$

$$= f(x) * (1 - f(x)) * x$$

# An Example

Emergency Room Visits	Narcotics	Pain	Total Visits	Medical Claims	PoorCare
0	2	6	11	53	1
1	1	4	25	40	0
0	0	5	10	28	0
1	3	5	7	20	1

**Initialise**  $w_1, w_2, \dots, w_5, b$

**Iterate over data:**

$$w_1 = w_1 - \eta \Delta w_1$$

$$w_2 = w_2 - \eta \Delta w_2$$

$\vdots$

$$w_5 = w_5 - \eta \Delta w_5$$

$$b = b - \eta \Delta b$$

**till satisfied**

$$z = w_1 * ER\_visits + w_2 * Narcotics + w_3 * Pain + w_4 * TotalVisits + w_5 * MedicalClaims + b$$

$$z = w_1 * x_{i1} + w_2 * x_{i2} + w_3 * x_{i3} + w_4 * x_{i4} + w_5 * x_{i5} + b$$

$$\hat{y} = \frac{1}{1+e^{-z}}$$

$$\hat{y} = \frac{1}{1+e^{-(w_1*x_{i1}+w_2*x_{i2}+w_3*x_{i3}+w_4*x_{i4}+w_5*x_{i5}+b)}}$$

$$\Delta w = \sum_{i=1}^m (f(x) - y) * f(x) * (1 - f(x)) * x$$

$$\Delta w_1 = \sum_{i=1}^m (\hat{y} - y) * \hat{y} * (1 - \hat{y}) * x_{i1}$$

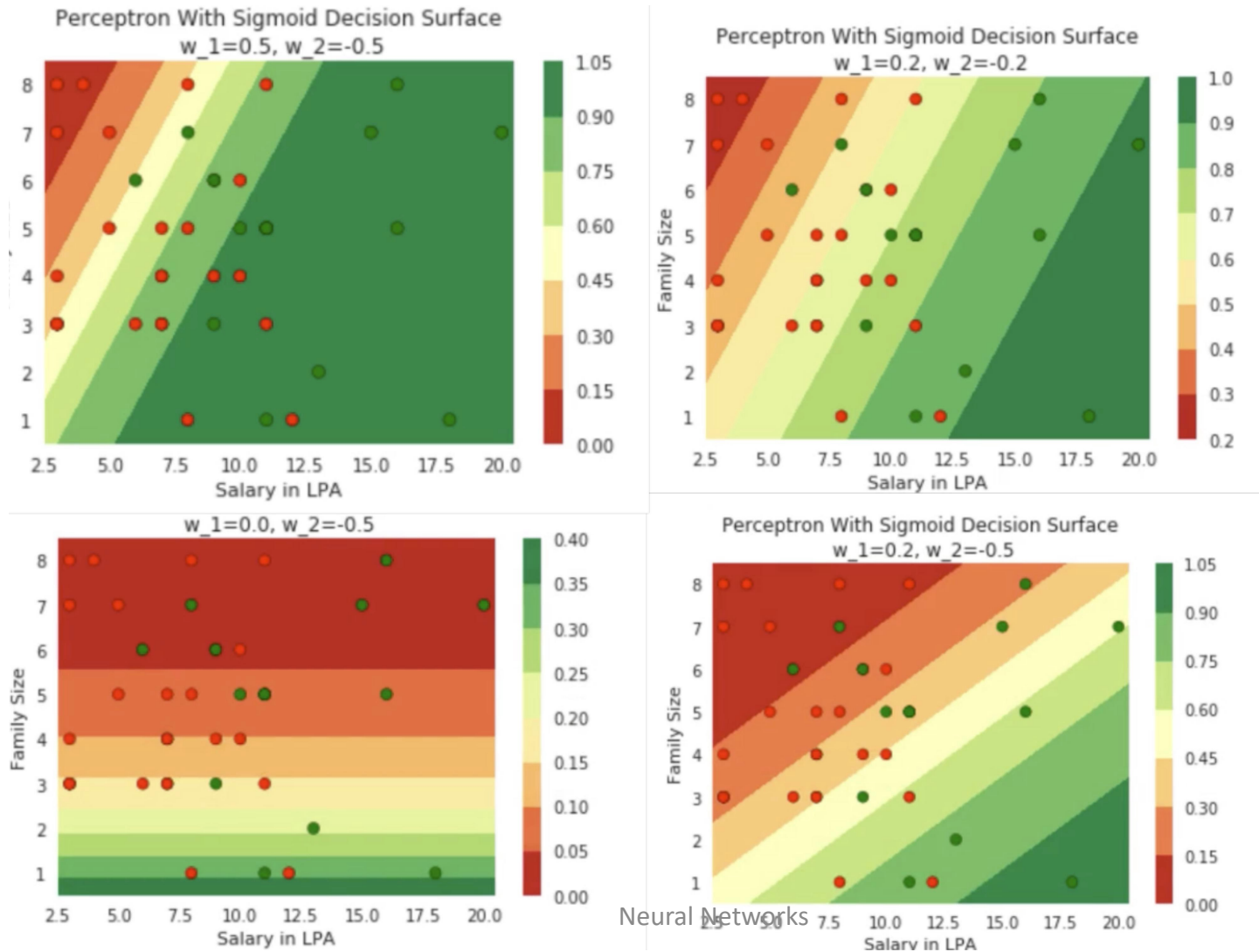
$$\Delta w_2 = \sum_{i=1}^m (\hat{y} - y) * \hat{y} * (1 - \hat{y}) * x_{i2}$$

$$\Delta w_j = \sum_{i=1}^m (\hat{y} - y) * \hat{y} * (1 - \hat{y}) * x_{ij}$$

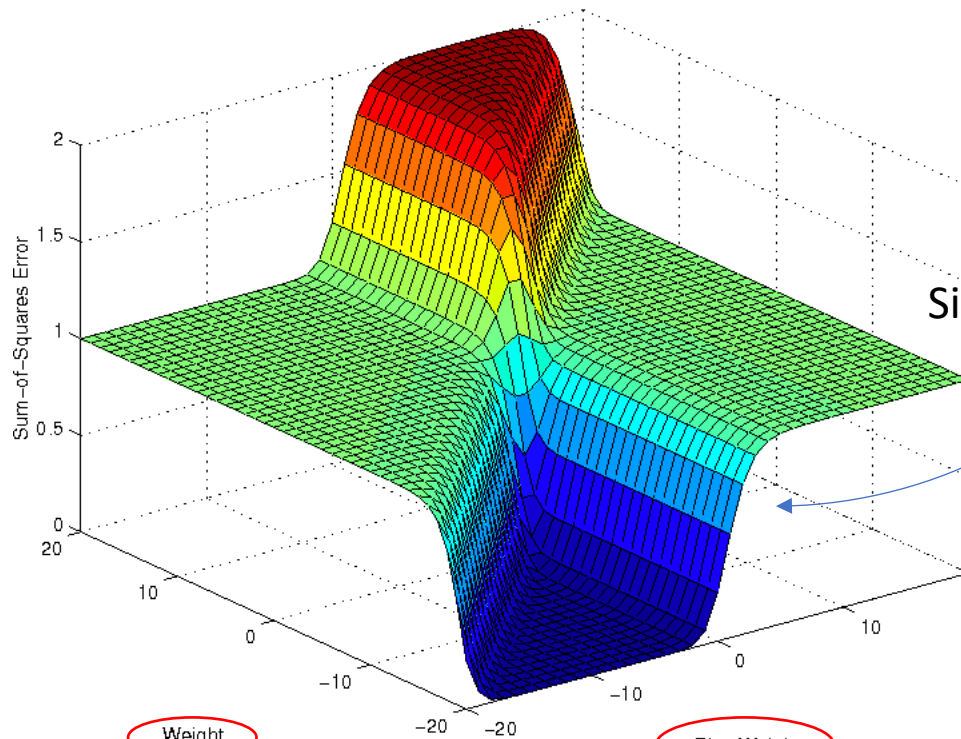
(c) One Fourth Labs



# Changes in sigmoid decision surface



Decision surface -> optimization space (yet another space you need to understand!)

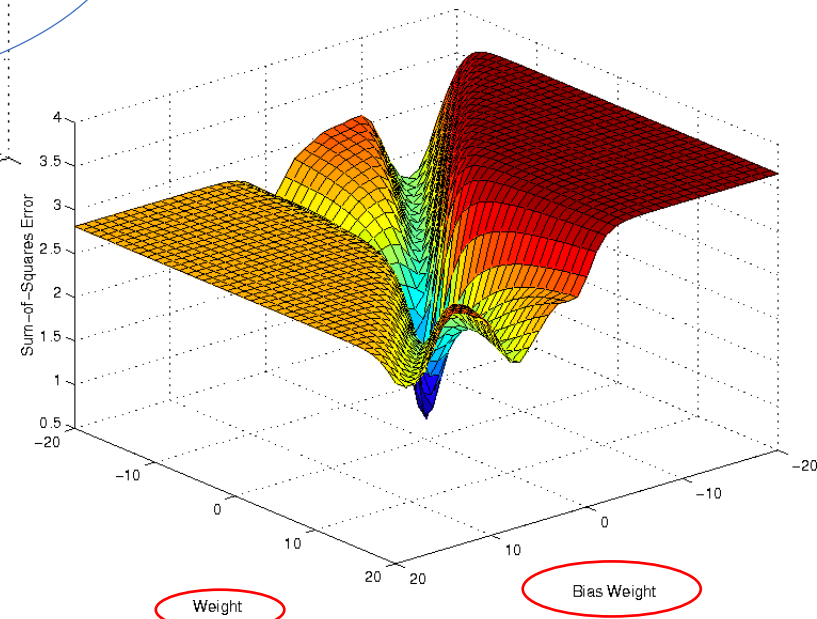


Sigmoid S shape

Weight

Bias Weight

Here, only 2 parameters!



Weight

Bias Weight

# How to iterate over data?

- **Batch Gradient Descent**

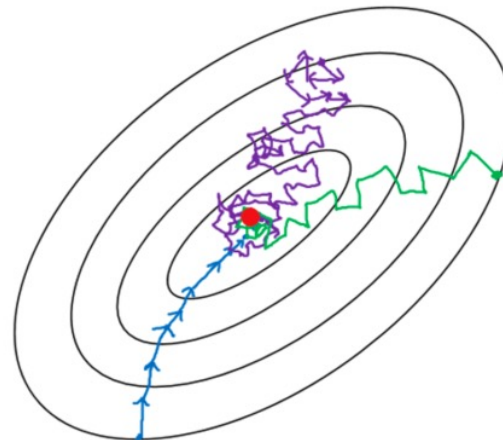
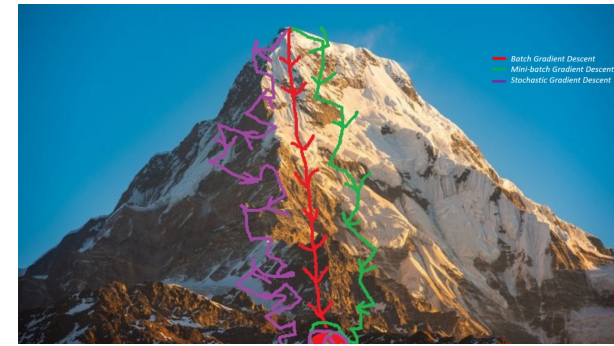
- uses the whole batch of training data at every step
- calculates the error for each record and takes an average to determine the gradient

- **Stochastic Gradient Descent (SGD)**

- just picks one instance from training set at every step
- update gradient only based on that single record

- **Mini-Batch Gradient Descent**

- Hybrid of Batch type and SGD
- Pick  $1 < n < N$  instance to update



- Batch gradient descent
- Mini-batch gradient Descent
- Stochastic gradient descent

# till satisfied...how to measure?

- Simple ones

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y - \hat{y})^2}$$

$$Accuracy = \frac{\text{Number of correct predictions}}{\text{Total number of predictions}}$$

# Take-home messages

- The **layered structure** of NN makes it powerful!
- NN, originated from neuroscience, has very nice mathematical structure.
- Machine learning through optimization of error/loss function -> mathematical structure is important!
- Gradient descent is simple and well-fit for ML!

# Acknowledgement

- Slides/Materials of

- [1] E. Alpaydin, Introduction to Machine Learning. 2<sup>nd</sup> Ed. MIT Press, 2010.

- [2] <https://medium.com/datadriveninvestor/simplified-sigmoid-neuron-a-building-block-of-deep-neural-network-5bfa75c8d8a9>

- Photos from Internet