

Intro to ML

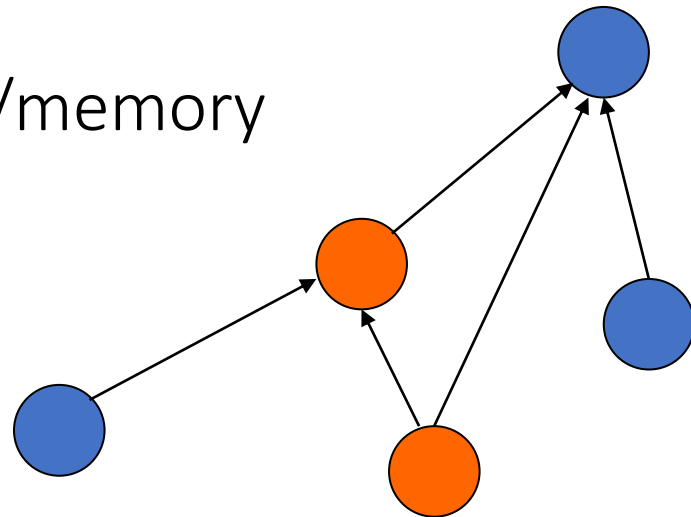
November 29th, 2021

CHAPTER 11:

Multilayer Perceptrons

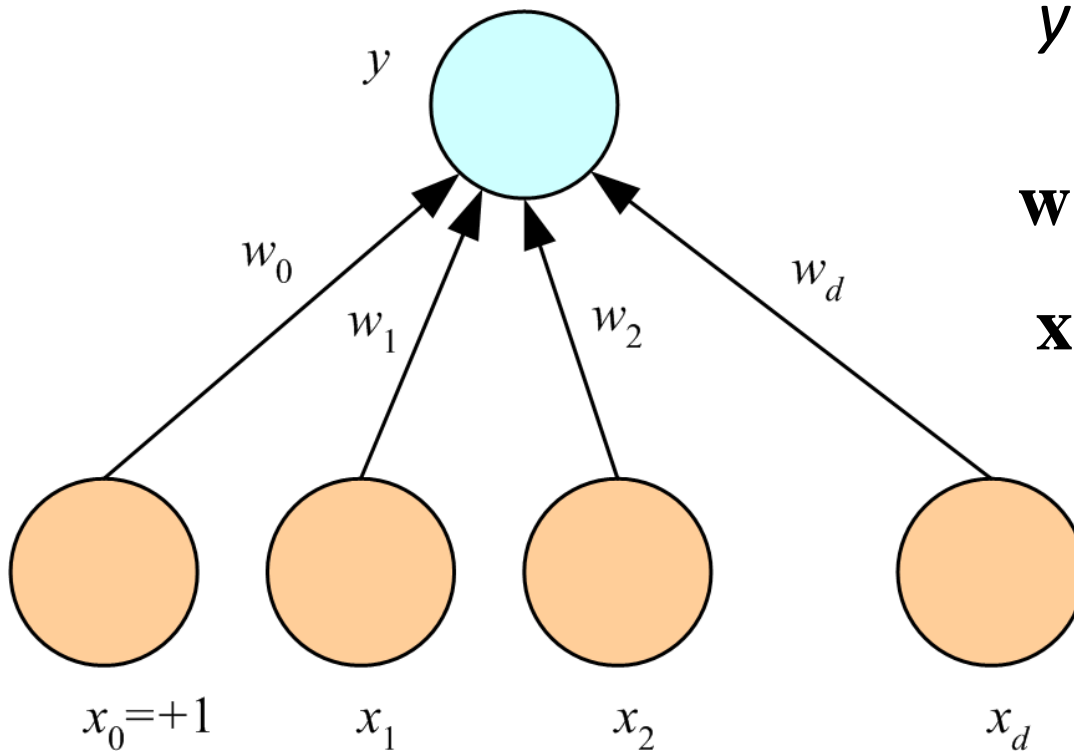
Neural Networks (human brain)

- 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



Perceptron

Basic processing unit



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

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

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

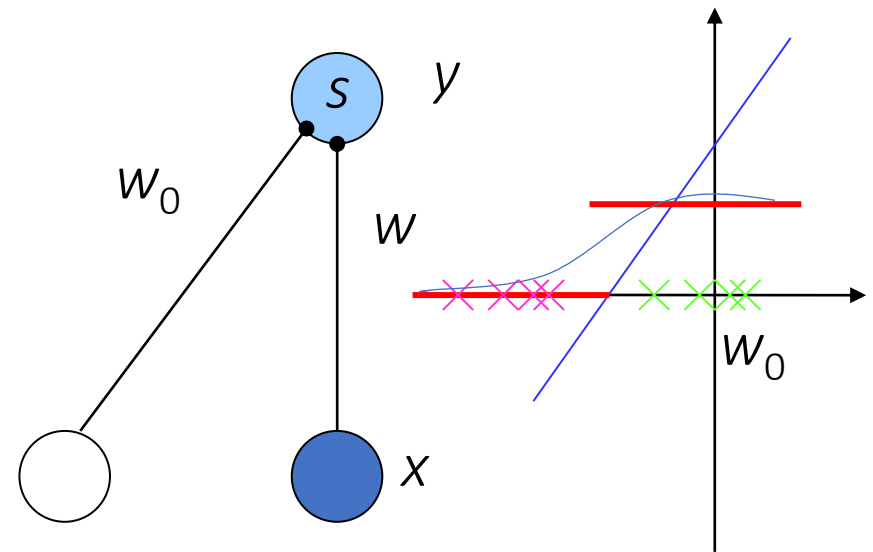
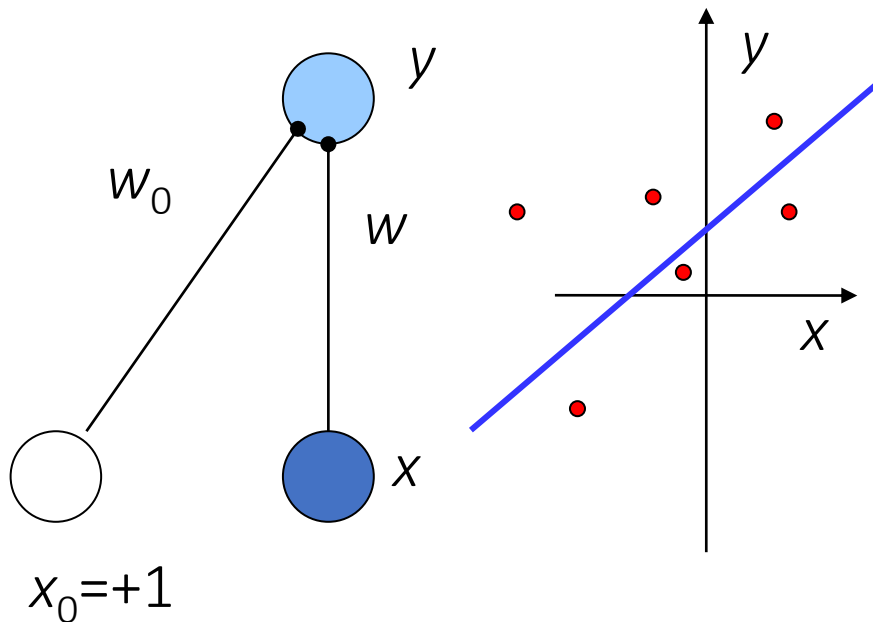
(Rosenblatt, 1962)

What a Perceptron Does

Simple
discriminant

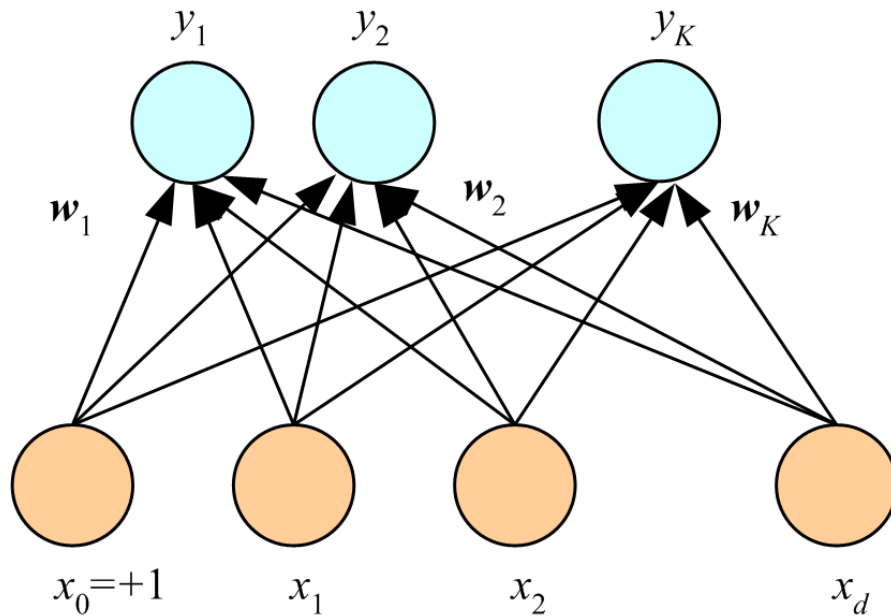
- Regression: $y = wx + w_0$

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



posterior $\rightarrow y = \text{sigmoid}(o) = \frac{1}{1 + \exp[-\mathbf{w}_5^T \mathbf{x}]}$

K Outputs



Regression:

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

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

Changing from d dimensions to K dimensions

Classification:

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

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

choose C_i

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

Imagine this as two steps:


1. Weighted sum
2. Softmax

Training

- Online (instances seen one by one) vs batch (whole sample) learning:
 - No need to store the whole sample
 - Problem may change in time
 - Wear and degradation in system components
- Stochastic gradient-descent: Update after a single pattern
- Generic update rule (LMS rule):

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

Update rule for
each training
sample



Update=LearningFactor·(DesiredOutput–ActualOutput)·Input

Training a Perceptron: Regression

- Regression (Linear output):

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

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



Update one by one
Stochastic gradient descent

Classification

- Single sigmoid output

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

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

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

The same as
previous
chapters
(logistic
discrimination)

- $K > 2$ softmax outputs

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

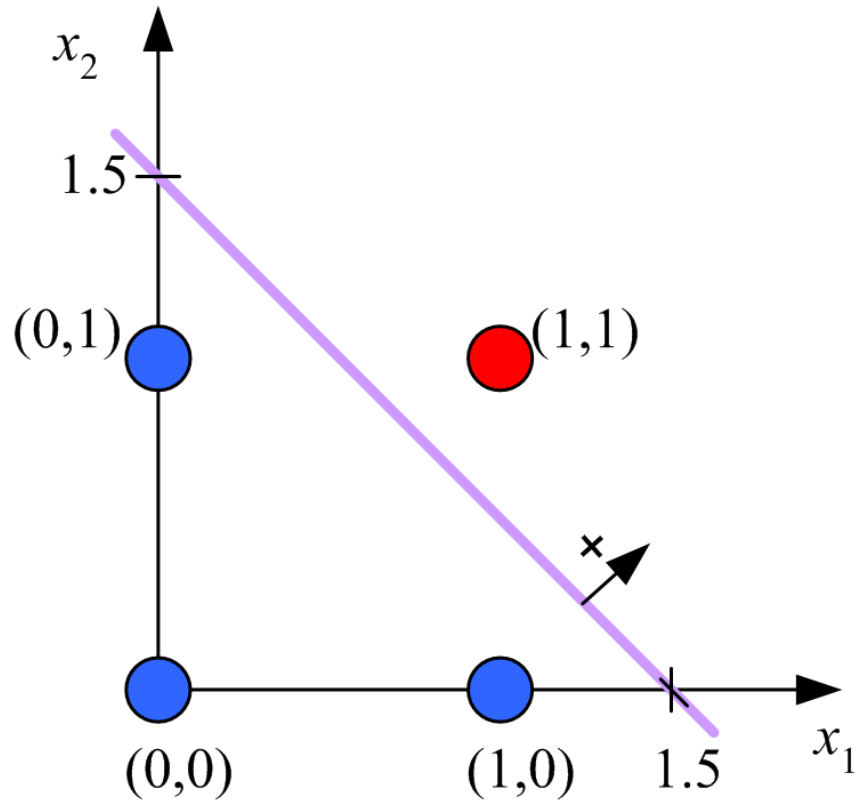
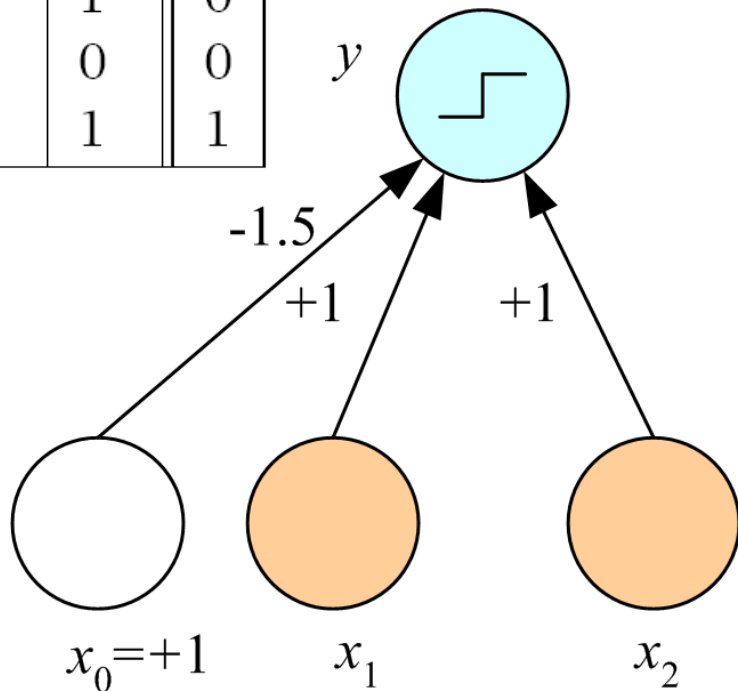
$$E^t(\{\mathbf{w}_i\}_i | \mathbf{x}^t, \mathbf{r}^t) = -\sum_i r_i^t \log y_i^t$$

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

Just for each
sample

Learning Boolean AND

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



XOR

Not linearly separable

Can't be solved by having a line as the perceptron

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

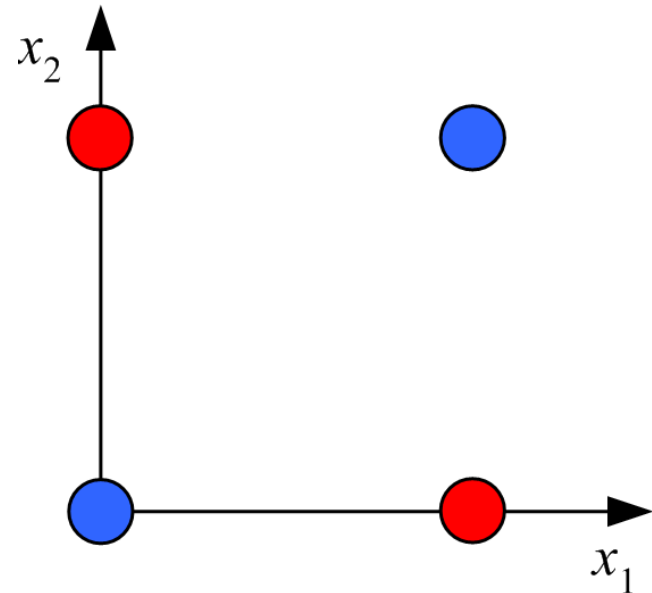
- No w_0, w_1, w_2 satisfy:

$$w_0 \leq 0$$

$$w_2 + w_0 > 0$$

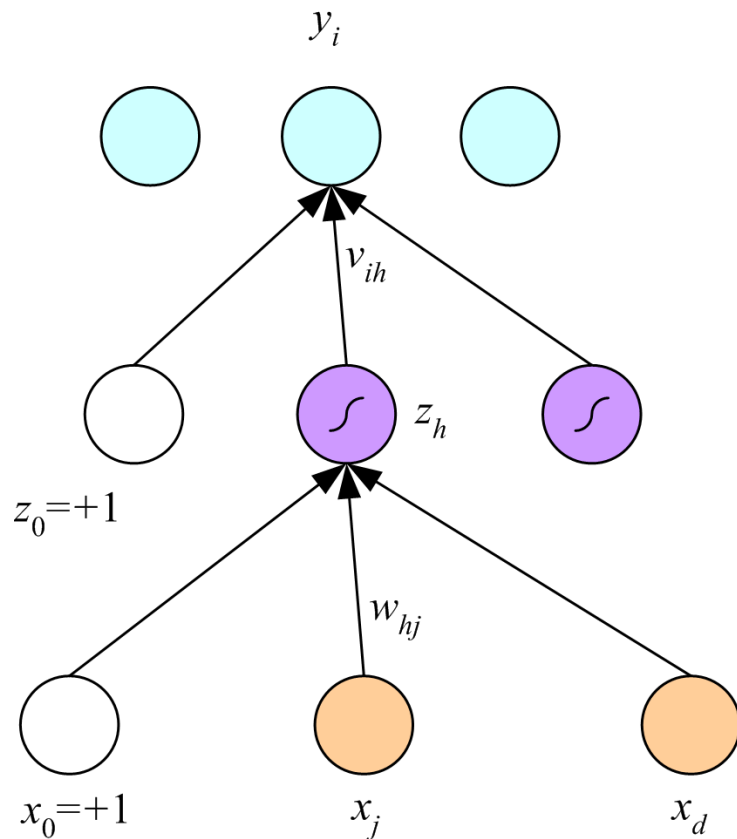
$$w_1 + w_0 > 0$$

$$w_1 + w_2 + w_0 \leq 0$$



(Minsky and Papert, 1969)

Multilayer Perceptrons

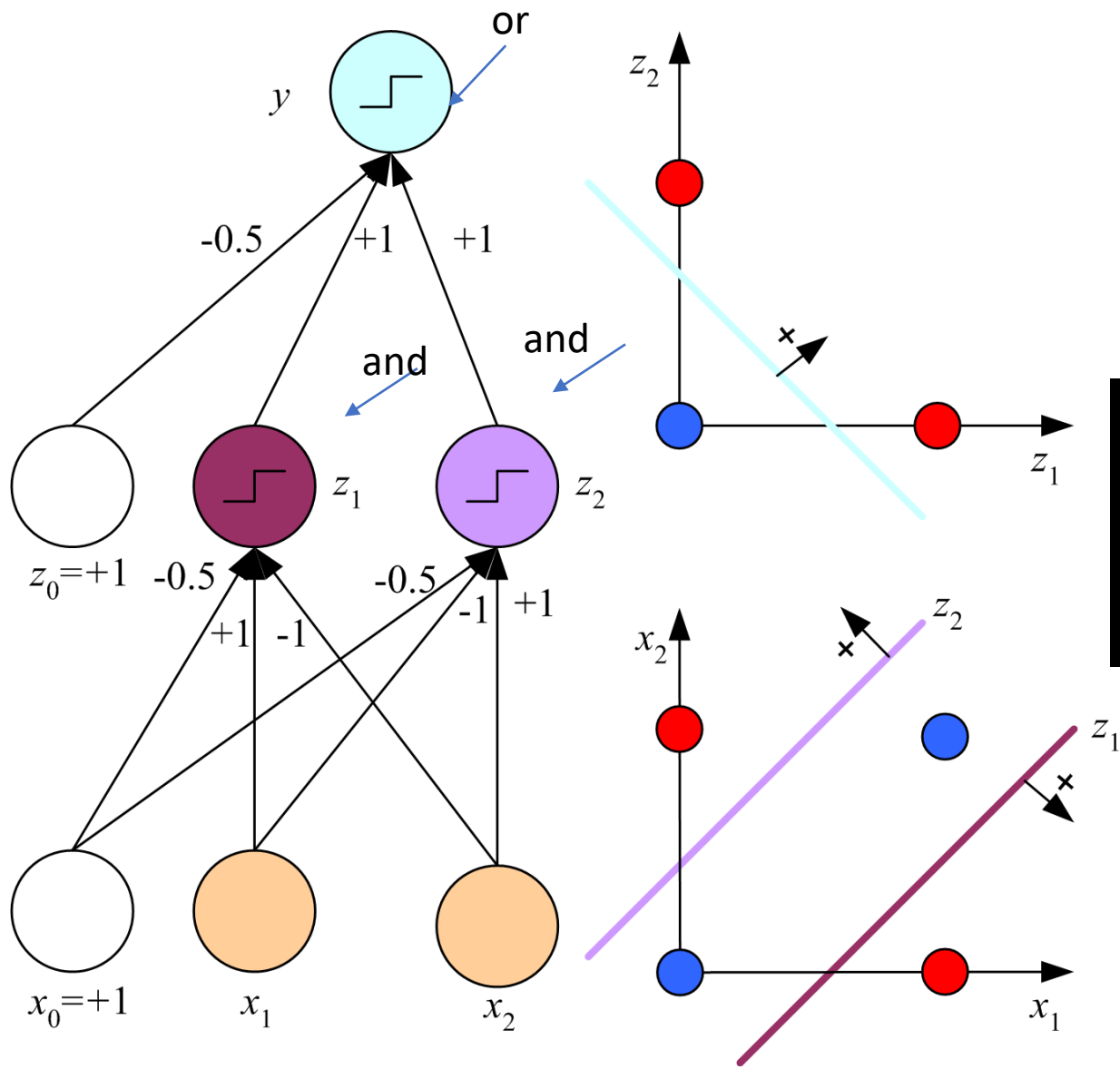


$$y_i = \mathbf{v}_i^T \mathbf{z} = \sum_{h=1}^H v_{ih} z_h + v_{i0}$$

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

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

(Rumelhart et al., 1986)

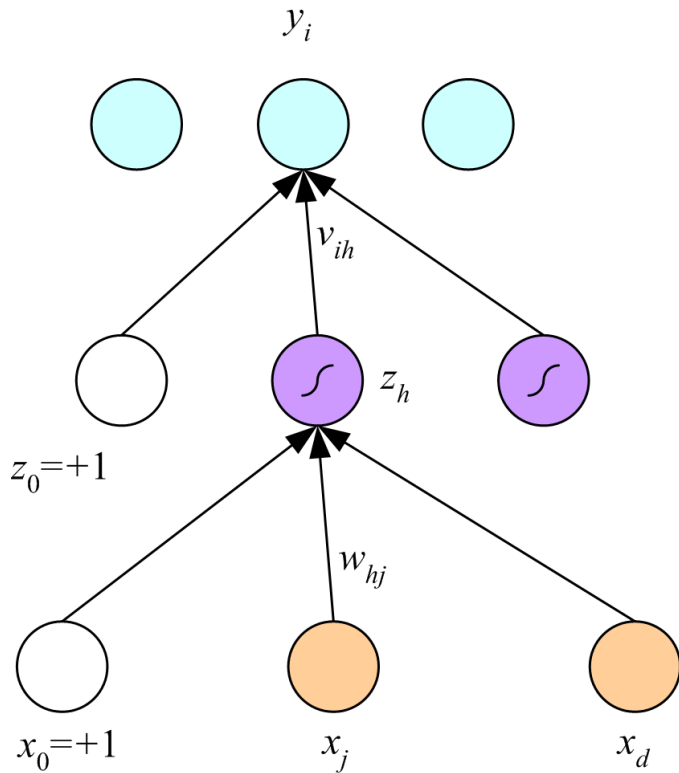


x_1	x_2	z_1	z_2	y
0	0	0	0	0
0	1	0	1	1
1	0	1	0	1
1	1	0	0	0

This implements xor

$$x_1 \text{ XOR } x_2 = (x_1 \text{ AND } \sim x_2) \text{ OR } (\sim x_1 \text{ AND } x_2)$$

Backpropagation



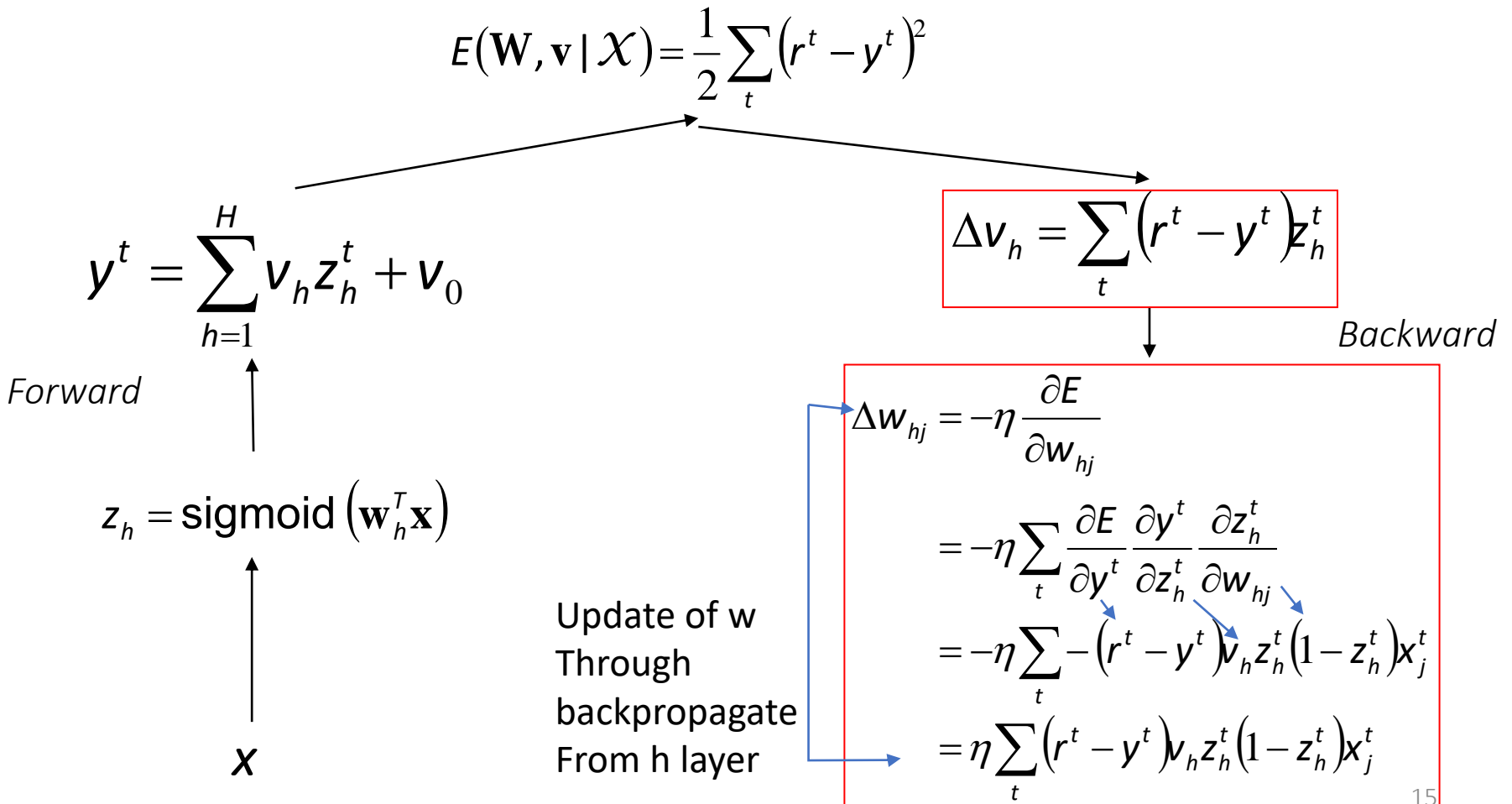
$$y_i = \mathbf{v}_i^T \mathbf{z} = \sum_{h=1}^H v_{ih} z_h + v_{i0}$$

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

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

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

Regression



Regression with Multiple Outputs

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

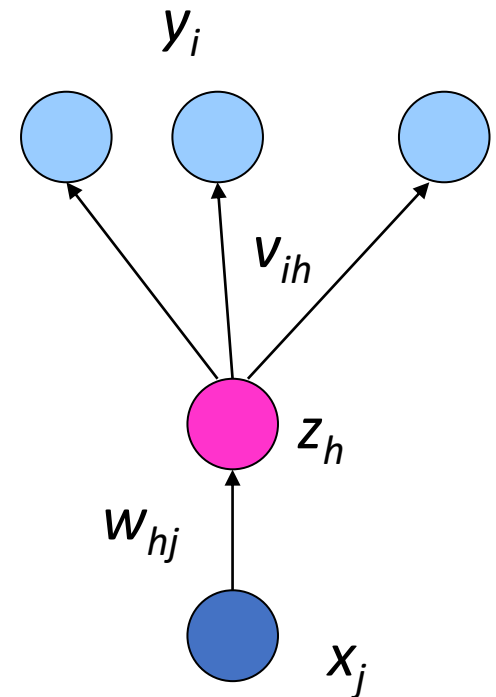
$$y_i^t = \sum_{h=1}^H v_{ih} z_h^t + v_{i0}$$

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

$$\Delta w_{hj} = \eta \sum_t \left[\sum_i (r_i^t - y_i^t) v_{ih} \right] z_h^t (1 - z_h^t) x_j^t$$

↓

Accumulated backpropagated error of
hidden unit h from all outputs



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 $h = 1, \dots, H$

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

For $i = 1, \dots, K$

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

For $i = 1, \dots, K$

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

For $h = 1, \dots, H$

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

For $i = 1, \dots, K$

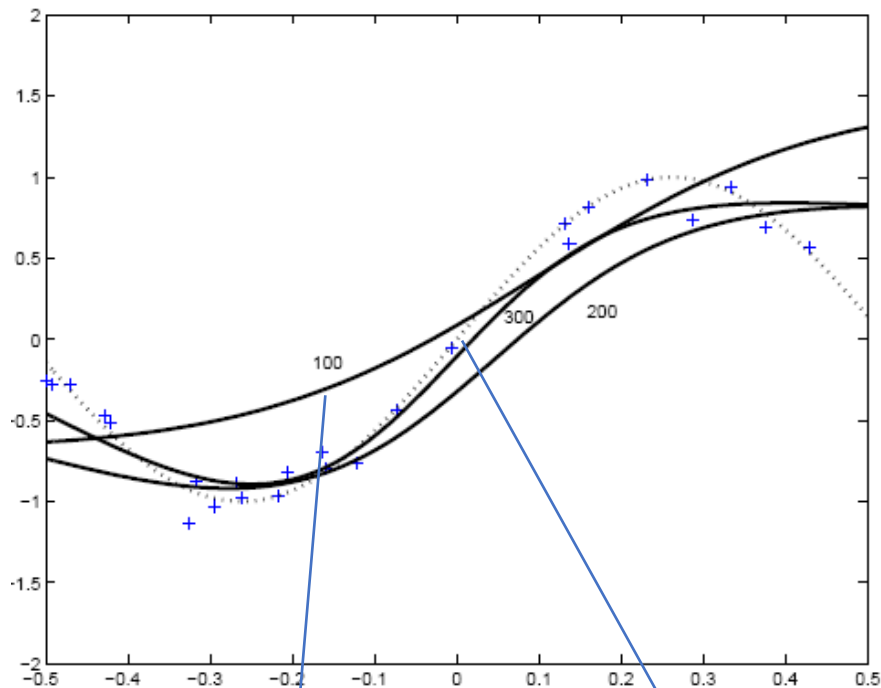
$$\mathbf{v}_i \leftarrow \mathbf{v}_i + \Delta \mathbf{v}_i$$

For $h = 1, \dots, H$

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

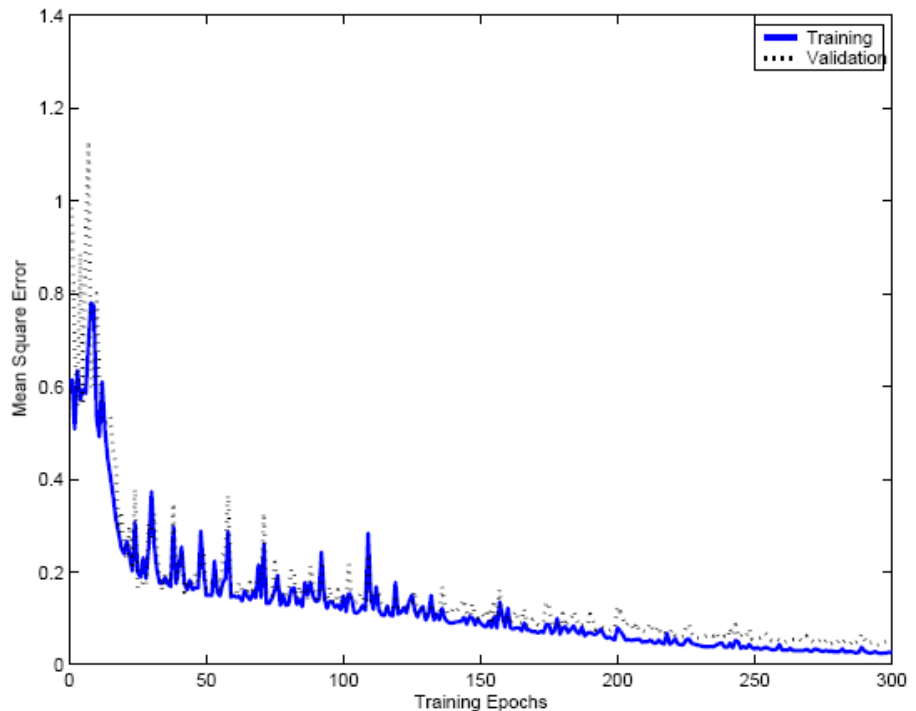
Until convergence

1d Regression: Convergence



2 hidden units
100 epochs..

$$y^t = f(x^t) + N(0, 0.1), f(x) = \sin(6x)$$



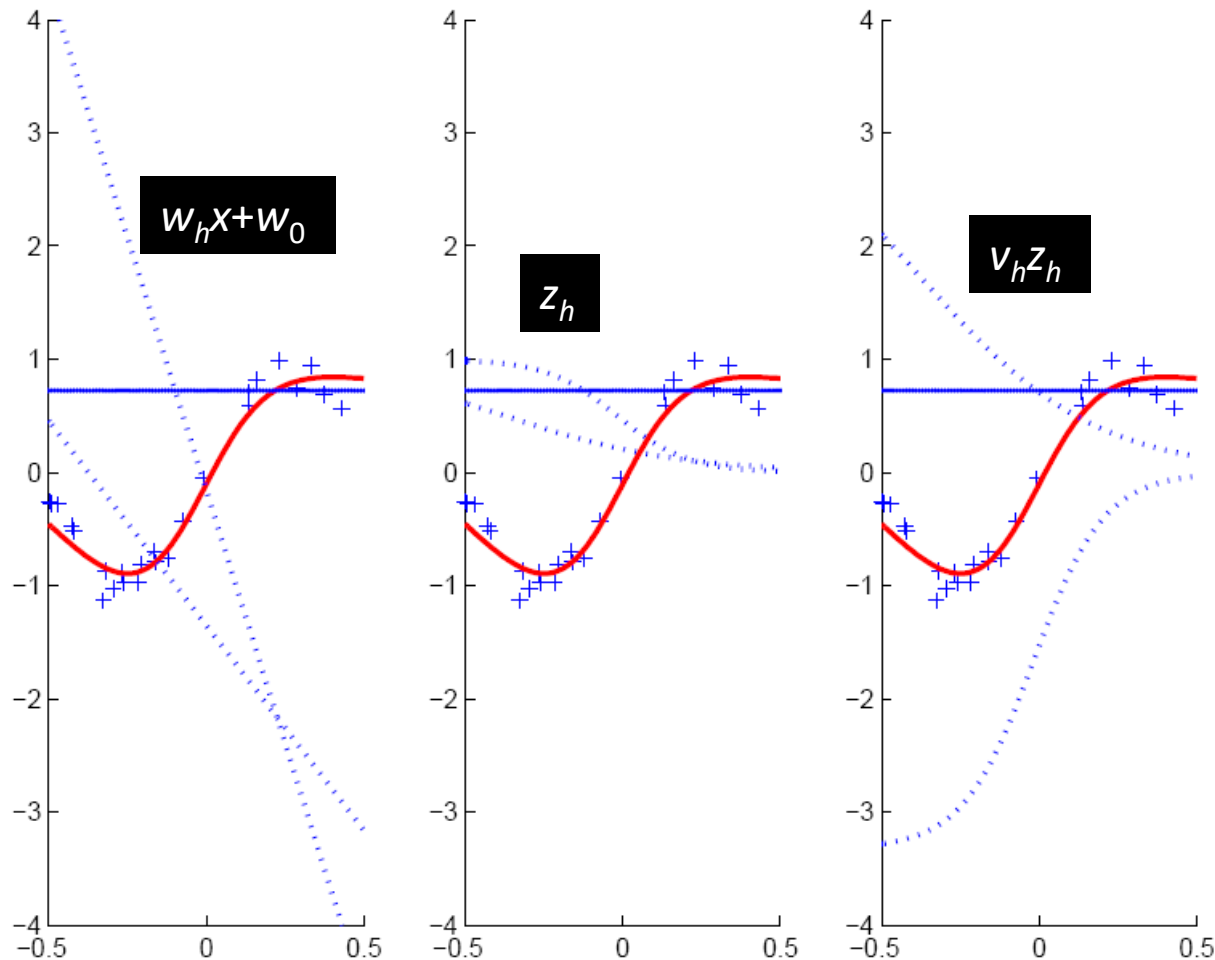
Epoch:

A complete pass over
all training samples

Minibatch:

Say 32 samples a time

1d Regression: What hidden units do



Two-Class Discrimination

- One sigmoid output y^t for $P(C_1|\mathbf{x}^t)$ and $P(C_2|\mathbf{x}^t) \equiv 1-y^t$

$$y^t = \text{sigmoid} \left(\sum_{h=1}^H \mathbf{v}_h z_h^t + \mathbf{v}_0 \right)$$

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

$$\Delta \mathbf{v}_h = \eta \sum_t (r^t - y^t) \mathbf{z}_h^t$$

$$\Delta \mathbf{w}_{hj} = \eta \sum_t (r^t - y^t) \mathbf{v}_h z_h^t (1 - z_h^t) \mathbf{x}_j^t$$

$K > 2$ Classes

softmax



$$o_i^t = \sum_{h=1}^H v_{ih} z_h^t + v_{i0} \quad y_i^t = \frac{\exp o_i^t}{\sum_k \exp o_k^t} \equiv P(C_i | \mathbf{x}^t)$$

$$E(\mathbf{W}, \mathbf{v} | \mathcal{X}) = - \sum_t \sum_i r_i^t \log y_i^t$$

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

$$\Delta w_{hj} = \eta \sum_t \left[\sum_i (r_i^t - y_i^t) v_{ih} \right] z_h^t (1 - z_h^t) x_j^t$$

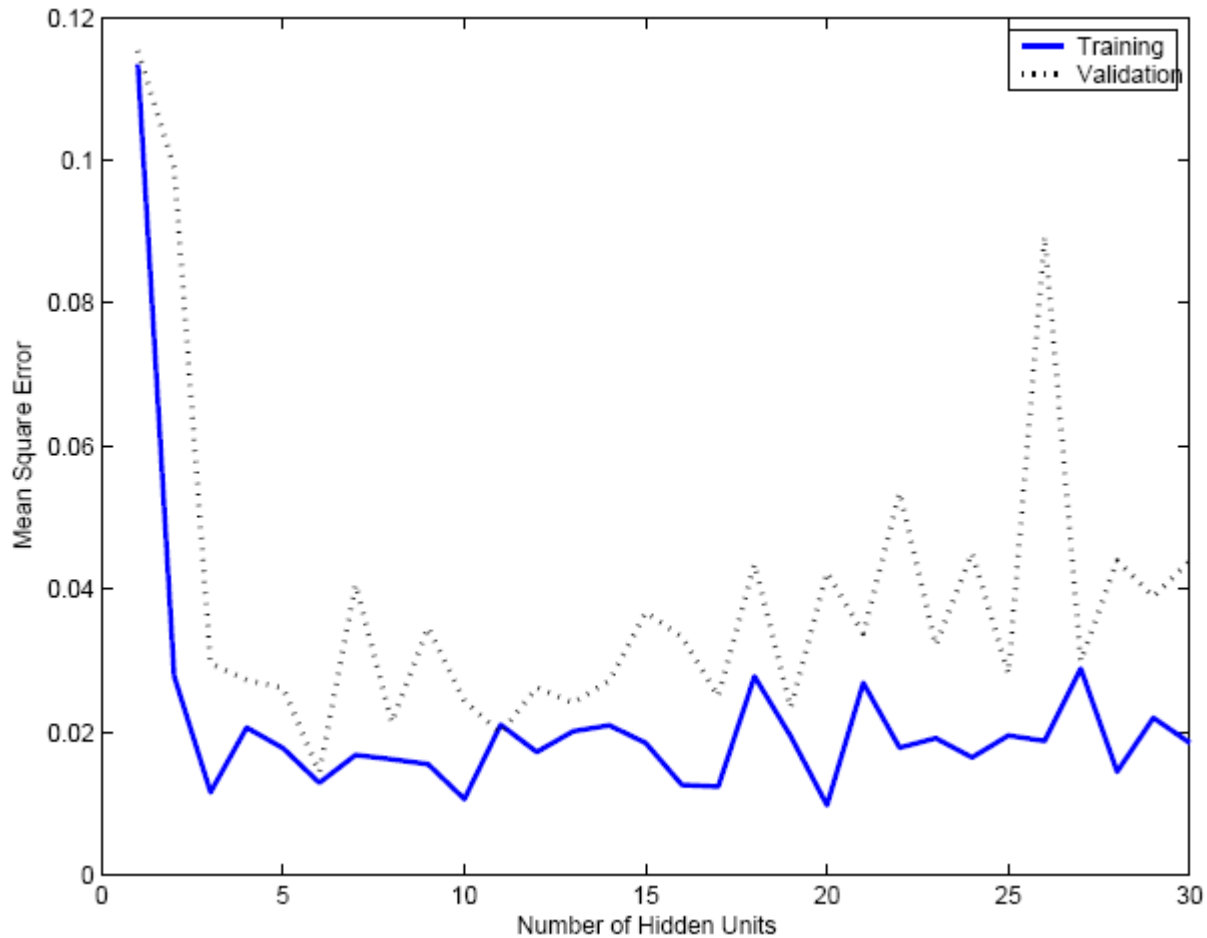
Overfitting

MLP:

d inputs

H hidden units

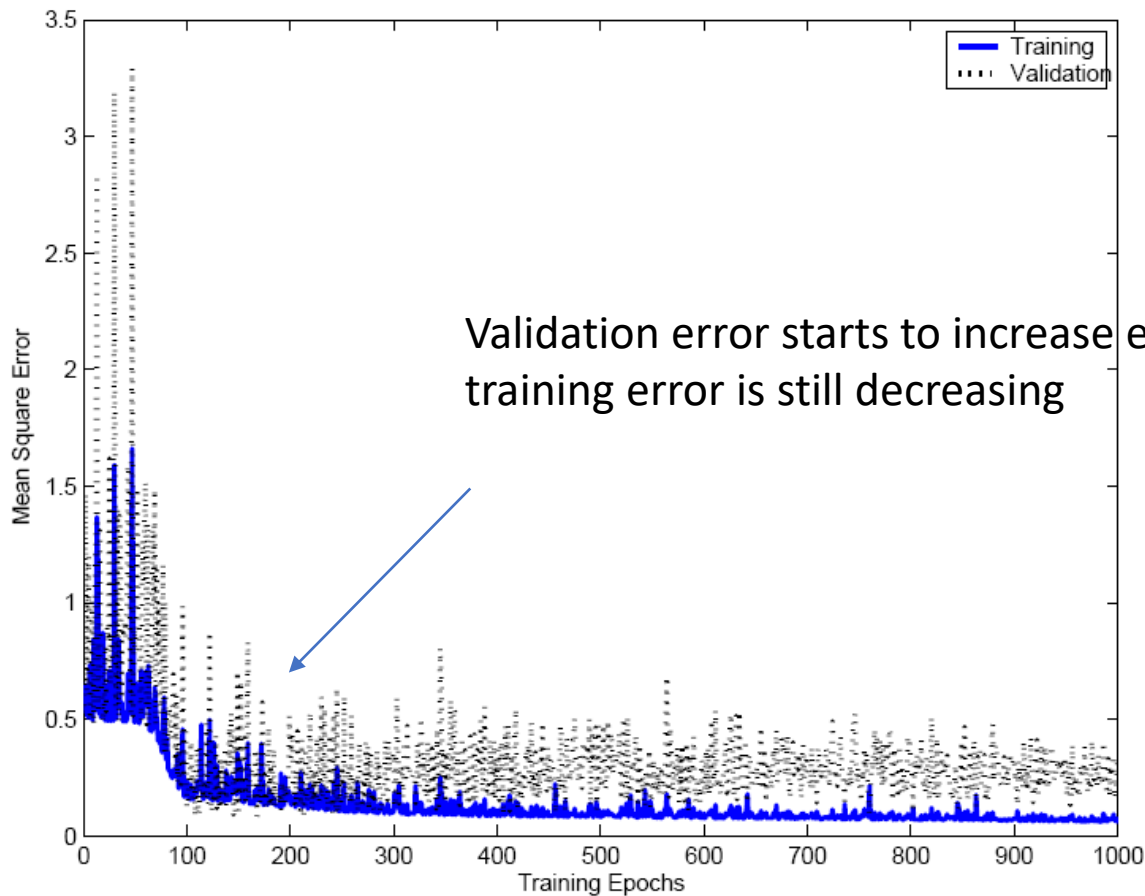
K outputs



Number of weights:

$$H(d+1)+(H+1)K$$

Overtraining



Most weights starts close to zero
Few jumps get updated
If keep training, those close to zeros will also be updated
(adding complexity to the model)

Learning Hidden Representations

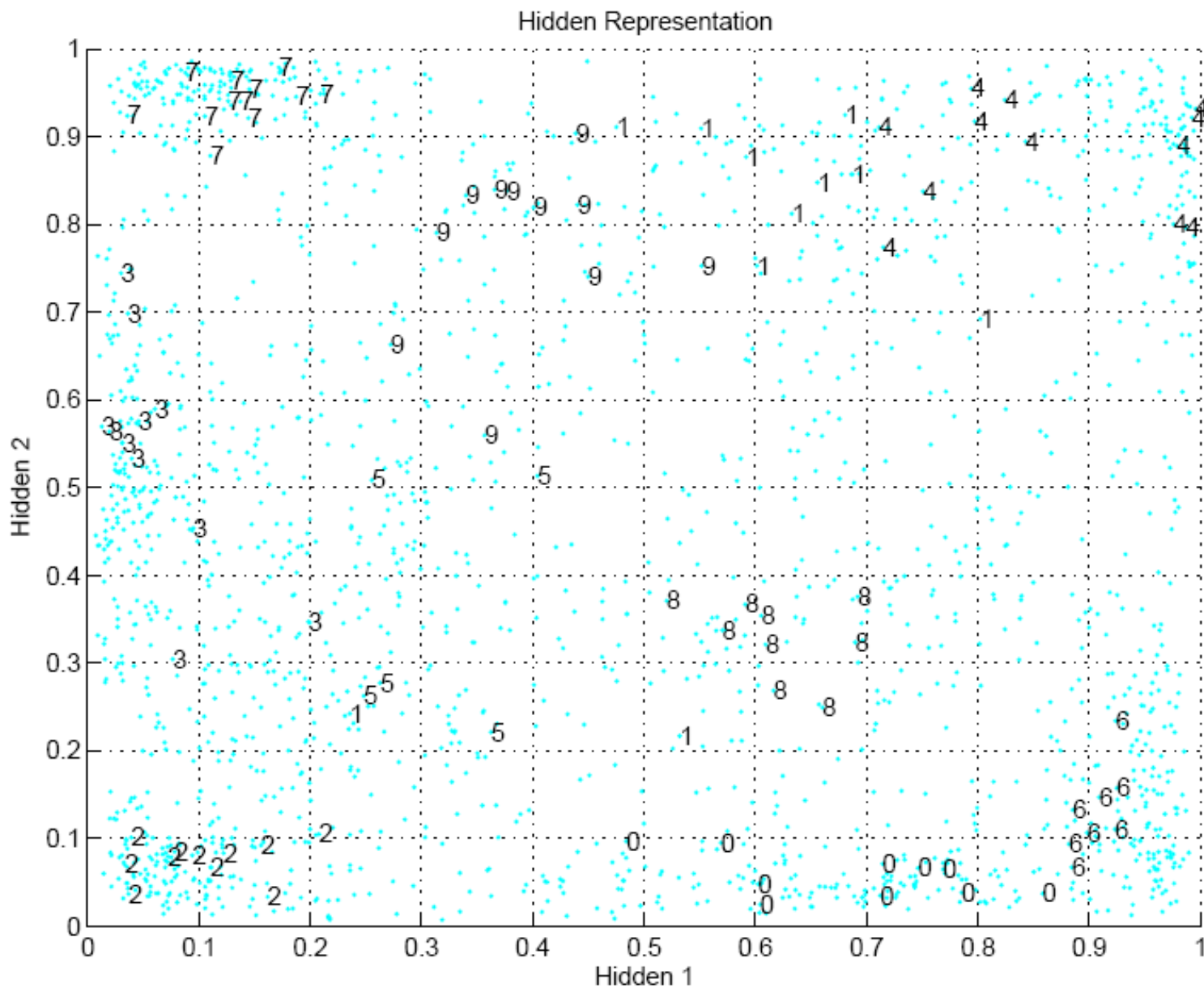
- MLP is a generalized linear model where hidden units are the nonlinear basis functions:

$$y = \sum_{h=1}^H v_h \phi(\mathbf{x} | \mathbf{w}_h)$$

where

$$\phi(\mathbf{x} | \mathbf{w}_h) \equiv \text{sigmoid}(\mathbf{w}_h^T \mathbf{x})$$

- The advantage is that the basis function parameters can also be learned from data.
- The hidden units, z_h , learn a *code/embedding*, a representation in the hidden space
 - If H space is < original d dimension: dimension reduction
- **Transfer learning**: Use code in another task
- Semisupervised learning: Transfer from an unsupervised to supervised problem

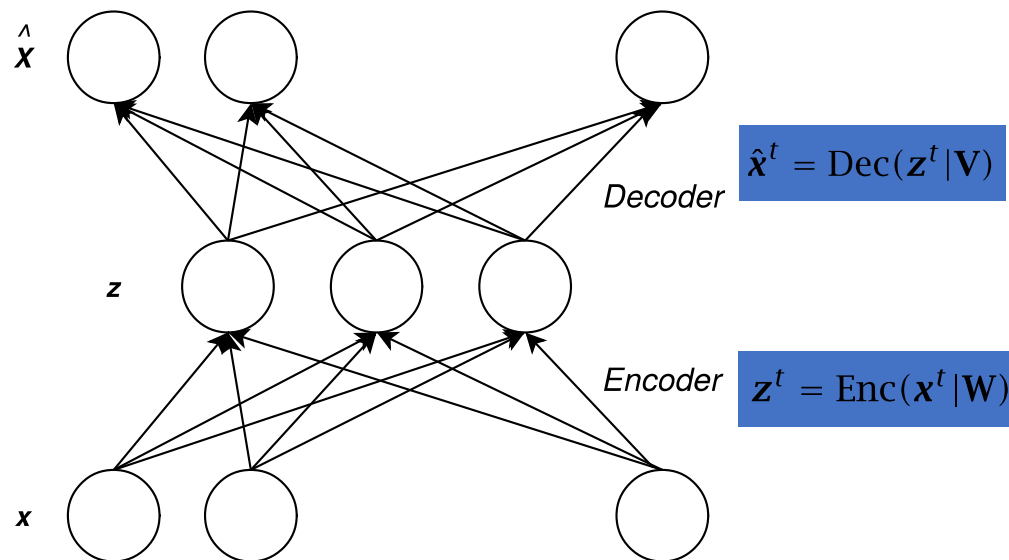


100 data
points

64 inputs
2 hidden
units

10 output

Autoencoders



An MLP structure

Equal number of input and output

Dimension reduction

$$E(\mathbf{W}, \mathbf{V} | \mathcal{X}) = \sum_t \|\mathbf{x}^t - \hat{\mathbf{x}}^t\|^2 = \sum_t \|\mathbf{x}^t - \text{Dec}(\text{Enc}(\mathbf{x}^t | \mathbf{W}) | \mathbf{V})\|^2$$

- Variants: Denoising, sparse autoencoders