Lecture Slides for

INTRODUCTION TO

Machine Learning

4e

ETHEM ALPAYDIN
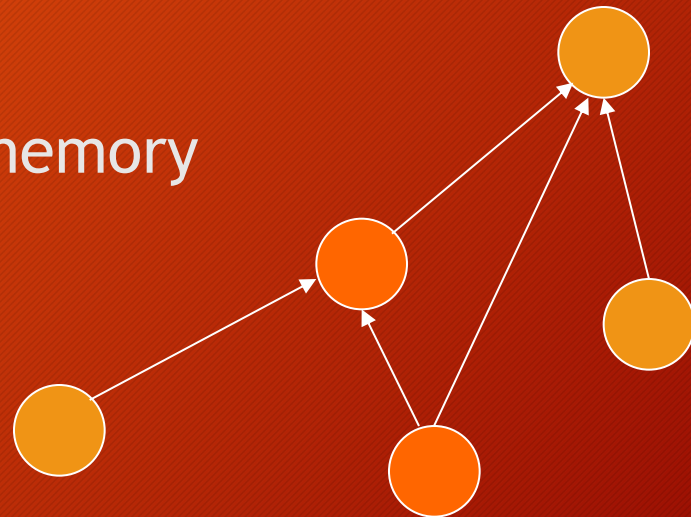The MIT Press, 2020

*ethem.alpaydin@gmail.com*

CHAPTER 11:

# Multilayer Perceptrons

# Neural Networks

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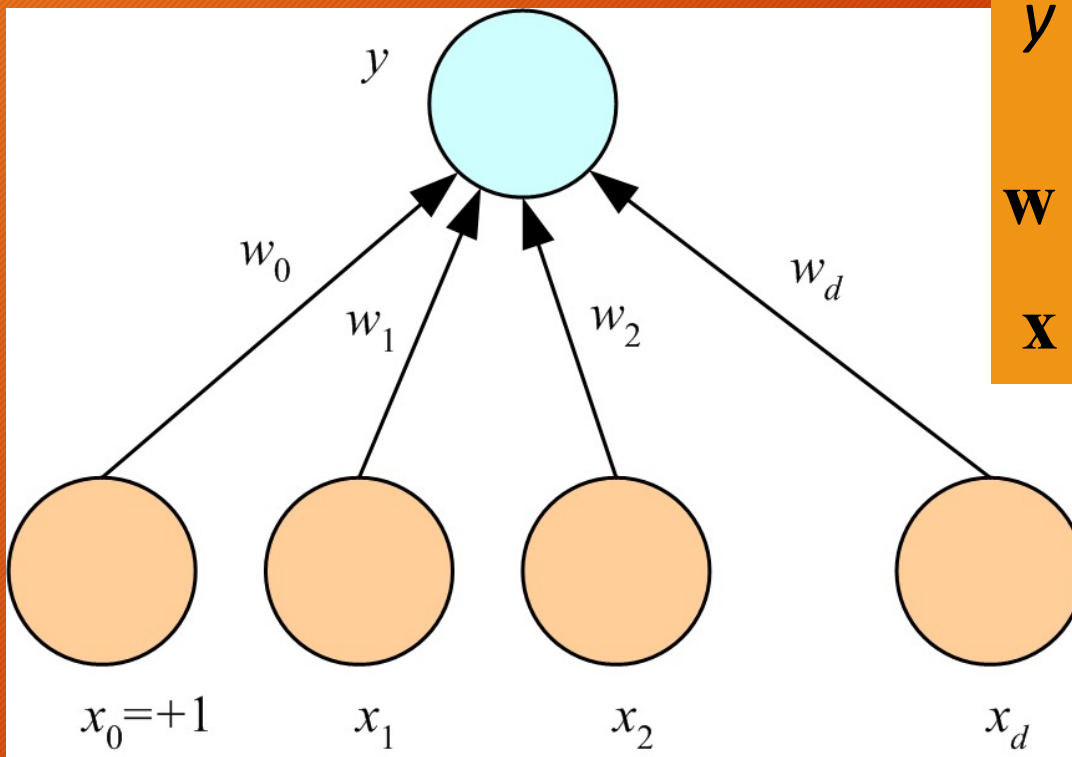- Distributed computation/memory
- Robust to noise, failures

# 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 vs MIMD

  Neural net: SIMD with modifiable local memory

  Learning: Update by training/experience

# Perceptron

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

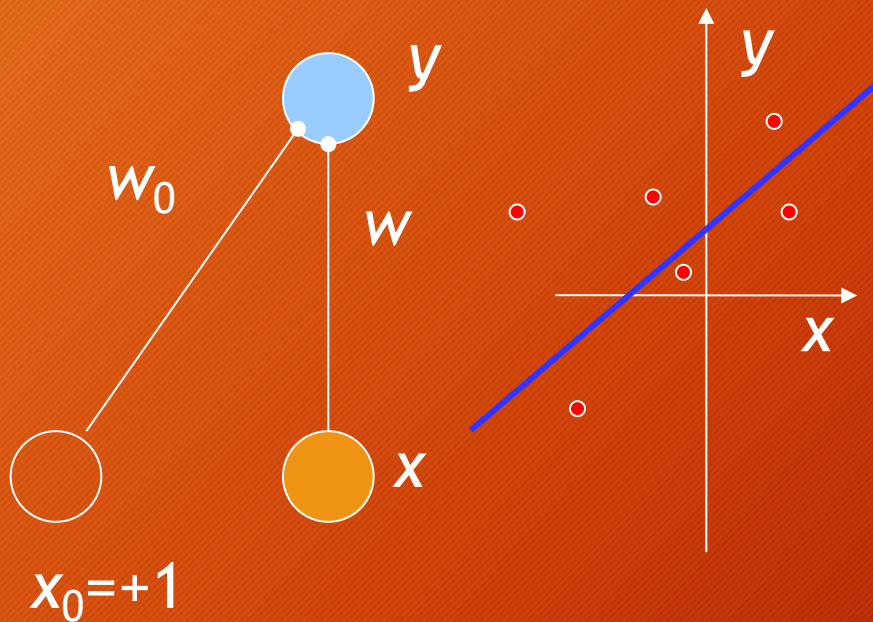$$\mathbf{w} = [w_0, w_1, ..., w_d]^T$$

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

(Rosenblatt, 1962)
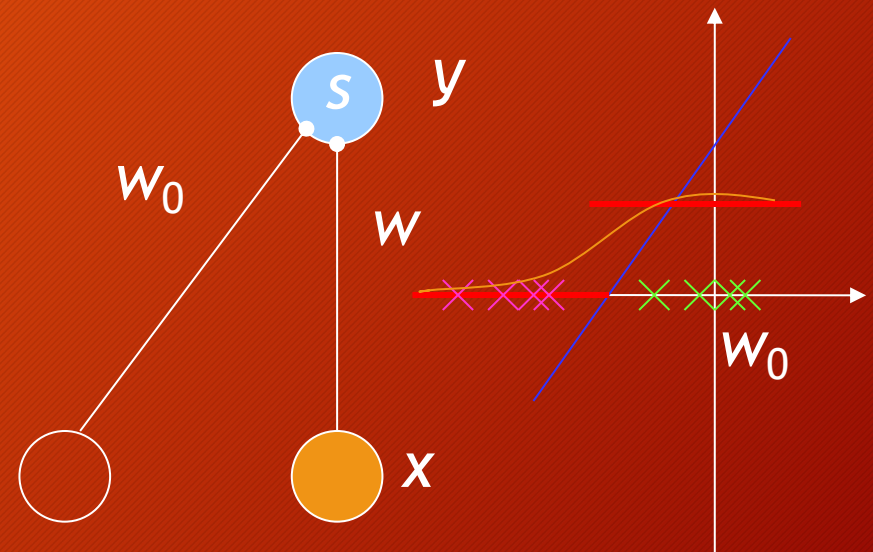
- Regression: $y = wx + w_0$
- Classification: $y = 1(wx + w_0 > 0)$

$w_0$

$w$

$y$

$x$

$x_0 = +1$

$s$

$w_0$

$w$

$y$

$x$

$y$

$x$

$y$

$w_0$

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

# *K* Outputs

Regression:

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

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

Classification:

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

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

choose $C_i$

if $y_i = \max_k y_k$

# 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):

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

$$Update = LearningFactor \cdot \left( DesiredOutput - ActualOutput \right) \cdot Input$$

# Training a Perceptron: Regression

- Regression (Linear output):

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

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

# Classification

- Single sigmoid output

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

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

$$\Delta w_j^t = \eta\left(r^t - y^t\right) x_j^t$$
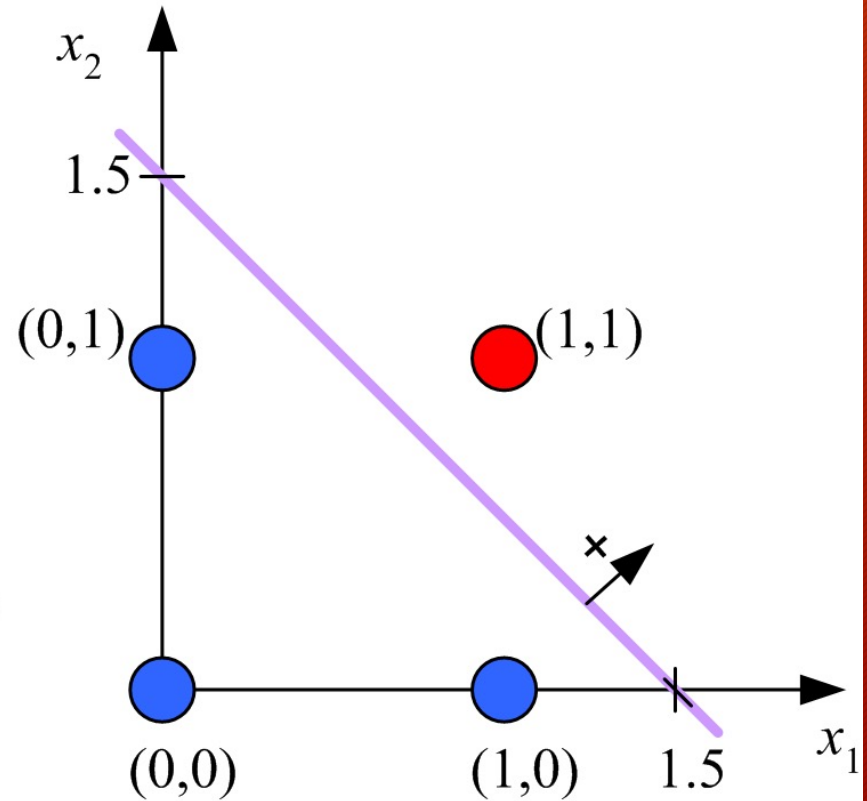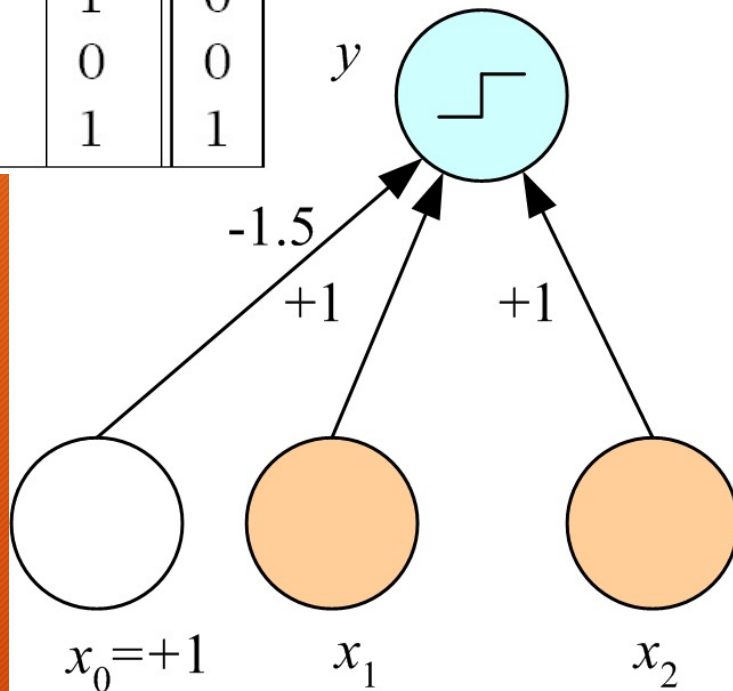
- $K>2$ softmax outputs

$$y^t = \frac{\exp \mathbf{w}_i^T \mathbf{x}^t}{\sum_k \exp \mathbf{w}_k^T \mathbf{x}^t} \qquad E^t\left(\{\mathbf{w}_i\}_i \mid \mathbf{x}^t, \mathbf{r}^t\right) = -\sum_i r_i^t \log y_i^t$$

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

| $x_1$ | $x_2$ | $r$ |
|-------|-------|-----|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

- No $w_0$, $w_1$, $w_2$ satisfy:

$$
\begin{aligned}
w_0 &\leq 0 \\
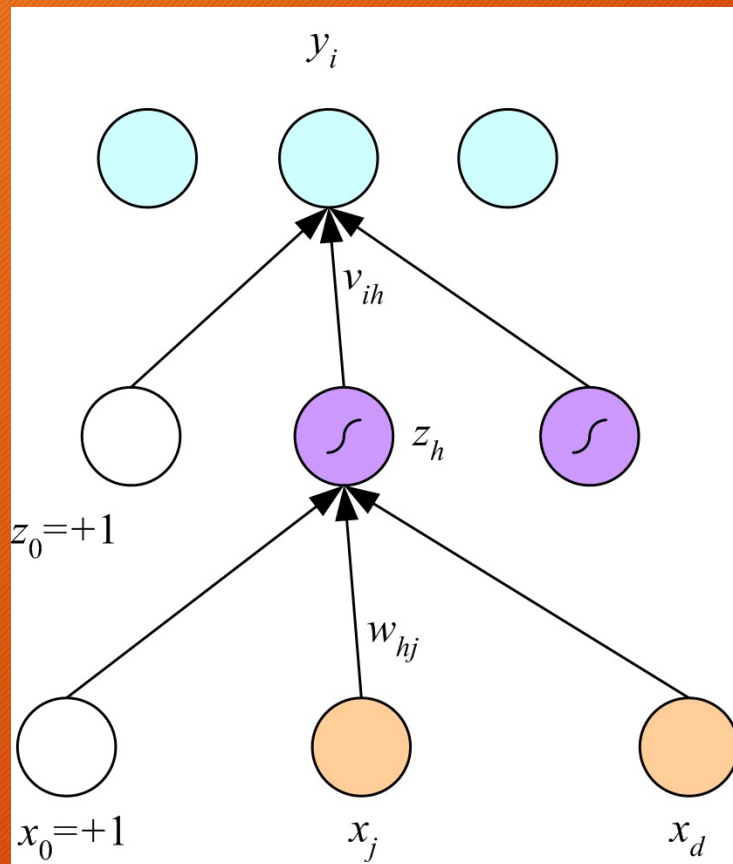w_2 + w_0 &> 0 \\
w_1 + w_0 &> 0 \\
w_1 + w_2 + w_0 &\leq 0
\end{aligned}
$$



(Minsky and Papert, 1969)

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

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

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

(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 |

$x_1$ XOR $x_2$ = ($x_1$ AND ~$x_2$) OR (~$x_1$ AND $x_2$)
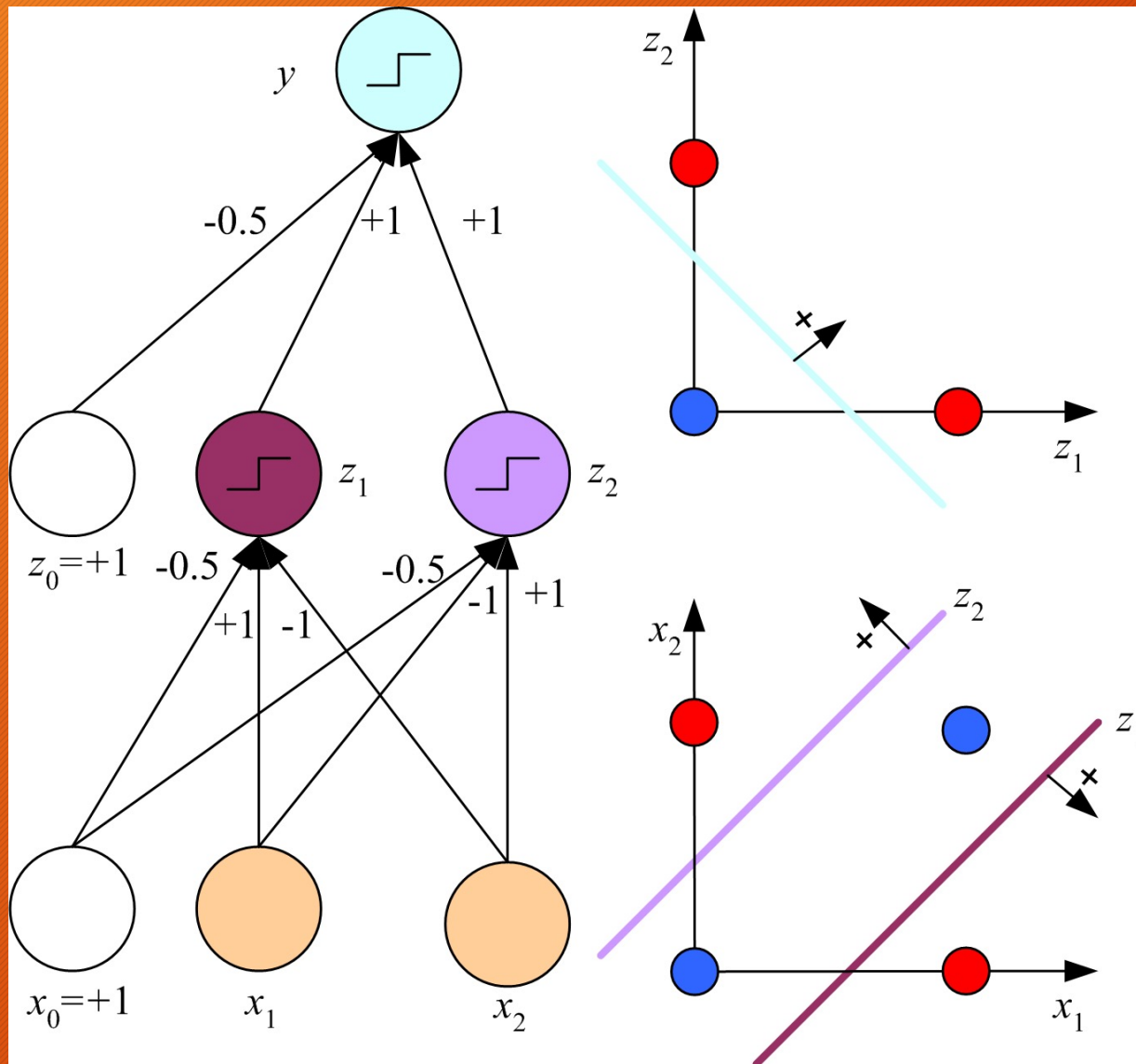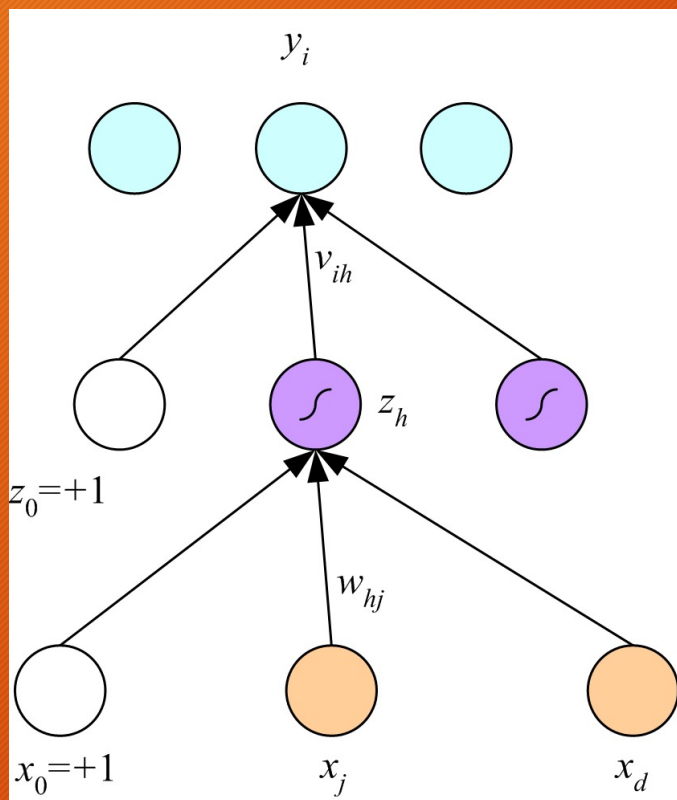
$$y_i = \mathbf{v}_i^T \mathbf{z} = \sum_{h=1}^{H} v_{ih} z_h + v_{i0}$$

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

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

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

# Regression

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

$$y^t = \sum_{h=1}^{H} v_h z_h^t + v_0$$

*Forward*

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

$\mathbf{x}$

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

*Backward*

$$\Delta w_{hj} = -\eta \frac{\partial E}{\partial w_{hj}}$$

$$= -\eta \sum_t \frac{\partial E}{\partial y^t} \frac{\partial y^t}{\partial z_h^t} \frac{\partial z_h^t}{\partial w_{hj}}$$

$$= -\eta \sum_t -\left( r^t - y^t \right) v_h z_h^t \left( 1 - z_h^t \right) x_j^t$$

$$= \eta \sum_t \left( r^t - y^t \right) v_h z_h^t \left( 1 - z_h^t \right) x_j^t$$

# Regression with Multiple Outputs
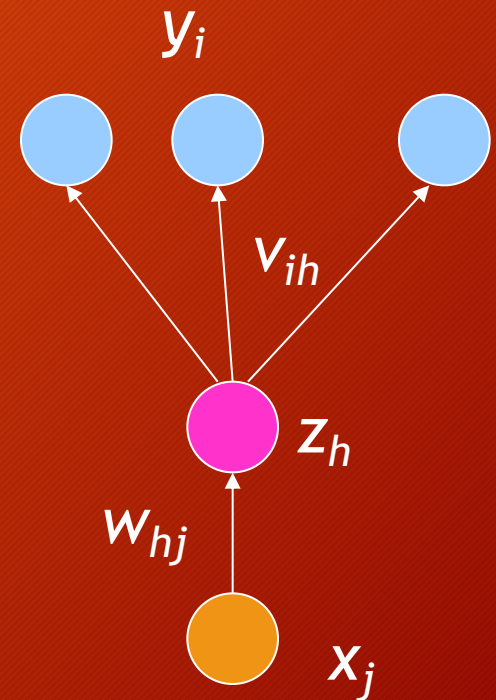
$$E(\mathbf{W}, \mathbf{V} \mid \mathcal{X}) = \frac{1}{2} \sum_t \sum_i \left( r_i^t - y_i^t \right)^2$$

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

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

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

$y_i$

$v_{ih}$

$z_h$

$w_{hj}$

$x_j$

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

Repeat

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

        For $h = 1, \ldots, H$

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

        For $i = 1, \ldots, K$

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

        For $i = 1, \ldots, K$

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

        For $h = 1, \ldots, H$

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

        For $i = 1, \ldots, K$

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

        For $h = 1, \ldots, H$

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

Until convergence

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

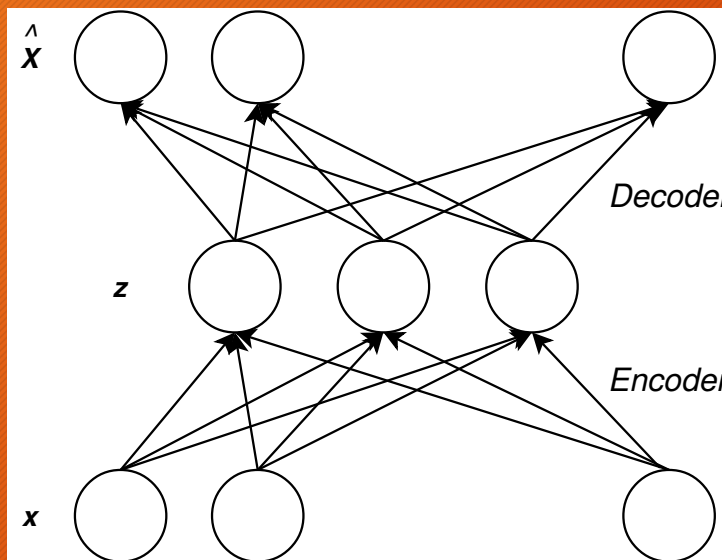$$y = \sum_{h=1}^{H} v_h \phi(\boldsymbol{x}|\boldsymbol{w}_h)$$

where

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

- 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

- Transfer learning: Use code in another task

- Semisupervised learning: Transfer from an unsupervised to supervised problem

$$\hat{\boldsymbol{x}}^t = \mathrm{Dec}(\boldsymbol{z}^t \,|\, \mathbf{V})$$

*Decoder*

*Encoder*

$$\boldsymbol{z}^t = \mathrm{Enc}(\boldsymbol{x}^t \,|\, \mathbf{W})$$

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

- Variants: Denoising, sparse autoencoders

# Word2vec

- Learn an embedding for words for NLP
- Skipgram: An autencoder with linear encoder and decoder where input is the center word and output is a context word
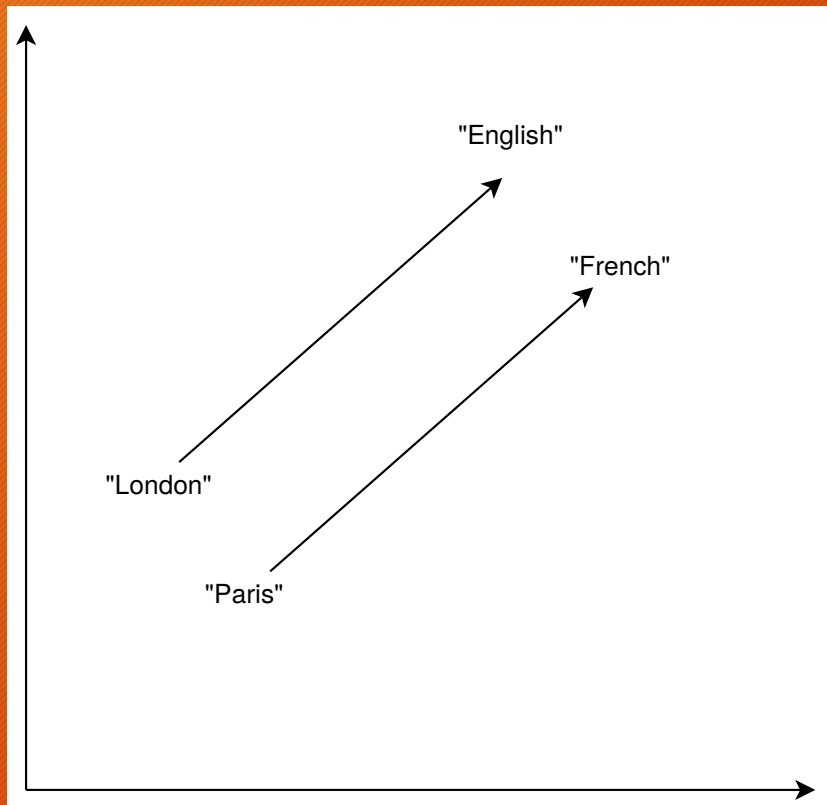
"Protesters in Paris clash with the French police."

Center

Context

- Similar words appear in similar contexts, so similar codes will be learned for them

# Vector Algebra

Because they will always appear in similar contexts in pairs in a large corpus, we expect

vec("English")-vec("London")

to be similar to

vec("French")-vec("Paris")

so,

vec("Paris") + vec("English")-vec("London")

will be similar to

vec("French")

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

$$y^t = \text{sigmoid}\left(\sum_{h=1}^{H} v_h z_h^t + v_0\right)$$

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

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

$$\Delta w_{hj} = \eta \sum_t \left(r^t - y^t\right) v_h z_h^t \left(1 - z_h^t\right) x_j^t$$