*The question of whether computers can think is like the question of whether submarines can swim.*

*- E. W. Dijsktra*

# CSE455 & CSE552
# Machine Learning

Spring 2025

**Neural Networks**

© 2013-2025 Yakup Genc

1

---

## Supervised Machine Learning

$$\arg\max_{f \in F} \sum_{i=1}^{n} g(y_i, f(x_i))$$

2

---

## Supervised Machine Learning

$$\arg\max_{W} \sum_{i=1}^{n} g(y_i, f_{nn}(W, x_i))$$

3

---

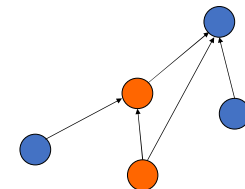## Neural Networks

- 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

4

## 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

5

## Supervised Machine Learning

$$\arg\max_{W} \sum_{i=1}^{n} g(y_i, f_{nn}(W, x_i))$$

6

## Perceptron



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

$$\mathbf{w} = [w_0, w_1, \ldots, w_d]^T$$
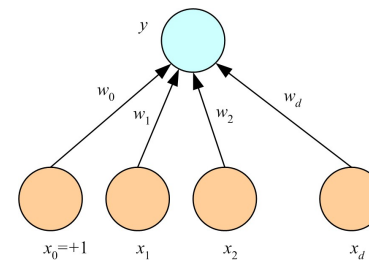
$$\mathbf{x} = [1, x_1, \ldots, x_d]^T$$

(Rosenblatt, 1962)

7

## Perceptron



$$y^{(i)} = \mathbf{w}^T \mathbf{x}^{(i)}$$

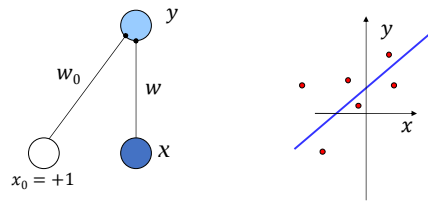$$\arg\max_{\mathbf{w}} \sum_{i=1}^{n} g(y^{(i)}, \mathbf{w}^T \mathbf{x}^{(i)})$$

8

## What a Perceptron Does

Regression: $y = wx + w_0$

9

## What a Perceptron Does

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

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

10

## Regression with K Outputs

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

$x_0 = +1 \quad x_1 \quad x_2 \quad x_d$
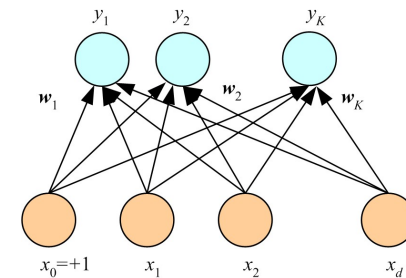
11

## Classification with K Outputs

Classification:
$$o_i = \boldsymbol{w}_i^T \boldsymbol{x}$$
$$y_i = \frac{e^{o_i}}{\sum_k e^{o_k}}$$

Choose $C_i$
   if $y_i = \max_k y_k$

$x_0 = +1 \quad x_1 \quad x_2 \quad x_d$

12

## Training

13

## 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):

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

- Update = Learning Factor * (Desired Output – Actual Output) * Input

14

## Training a Perceptron: Regression

- Regression (Linear output):

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

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

15

## Classification

- Single sigmoid output

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

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

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

16

## Classification

- $K>2$ softmax outputs

$$y^t = \frac{\exp w_i^T x^t}{\sum_k \exp w_k^T x^t} \quad E^t\left(\{w_i\}_i \mid x^t, r^t\right) = -\sum_i r_i^t \log y_i^t$$
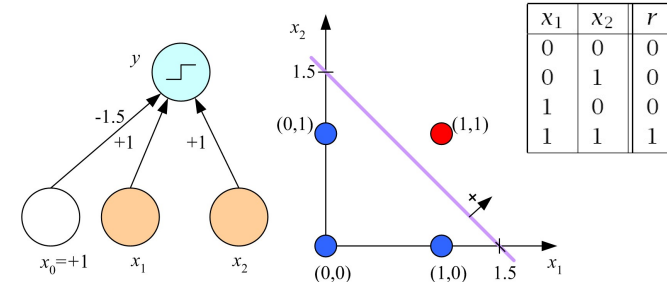
$$\Delta w_{ij}^t = \eta\left(r_i^t - y_i^t\right) x_j^t$$

17

## Learning Boolean AND



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

18

## Learning Boolean AND

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

19

## Learning Boolean AND

1*1-1.5*1+-1*0=-0.5>0 → 0
Deltaw1 = 1*(0-0)*1
Deltaw2 = 1*(0-0)*0



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

20

## Slide 21: XOR

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

21

## Slide 22: XOR

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

- No $w_0, w_1, w_2$ satisfy:

$$w_0 \le 0$$
$$w_2 + w_0 > 0$$
$$w_1 + w_0 > 0$$
$$w_1 + w_2 + w_0 \le 0$$

(Minsky and Papert, 1969)

22

## Slide 23: Multilayer Perceptrons

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

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

$$= \frac{1}{1 + \exp\left[-\left(\sum_{j=1}^{d} w_{hj} x_j + w_{h0}\right)\right]}$$
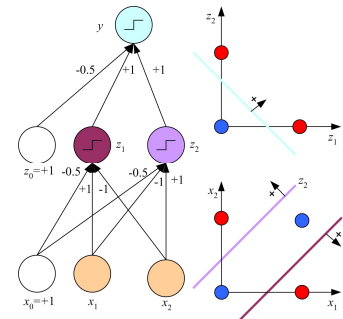
(Rumelhart et al., 1986)

23

## Slide 24: XOR with Multilayer Perceptrons

$x_1$ XOR $x_2$ = ($x_1$ AND $\sim x_2$) OR ($\sim x_1$ AND $x_2$)

24

## Backpropagation



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

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

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

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

25

## Regression

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

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

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

*Backward*

*Forward*

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

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

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

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

$$\mathbf{x}$$

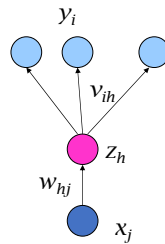$$= \eta \sum_t (r^t - y^t) v_h z_h^t (1 - z_h^t) x_j^t$$

26

## Regression with Multiple Outputs

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

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

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

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

27

## Algorithm
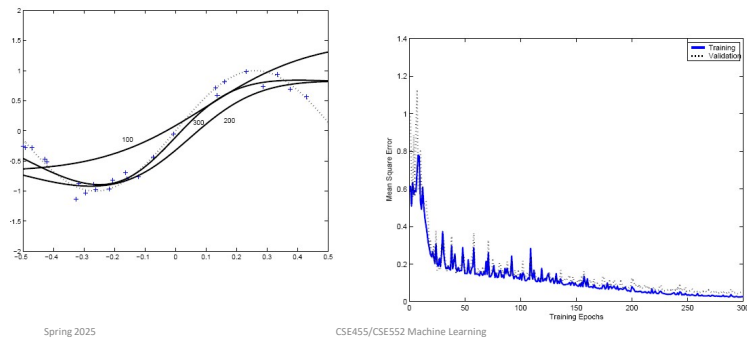
Initialize all $v_{ih}$ and $w_{hj}$ to rand$(-0.01, 0.01)$
Repeat
    For all $(\mathbf{x}^t, r^t) \in X$ in random order
        For $h = 1, \ldots, H$
            $z_h \leftarrow \text{sigmoid}(\mathbf{w}_h^T \mathbf{x}^t)$
        For $i = 1, \ldots, K$
            $y_i = \mathbf{v}_i^T \mathbf{z}$
        For $i = 1, \ldots, K$
            $\Delta \mathbf{v}_i = \eta (r_i^t - y_i^t) \mathbf{z}$
        For $h = 1, \ldots, H$
            $\Delta \mathbf{w}_h = \eta (\sum_i (r_i^t - y_i^t) v_{ih}) z_h (1 - z_h) \mathbf{x}^t$
        For $i = 1, \ldots, K$
            $\mathbf{v}_i \leftarrow \mathbf{v}_i + \Delta \mathbf{v}_i$
        For $h = 1, \ldots, H$
            $\mathbf{w}_h \leftarrow \mathbf{w}_h + \Delta \mathbf{w}_h$
Until convergence
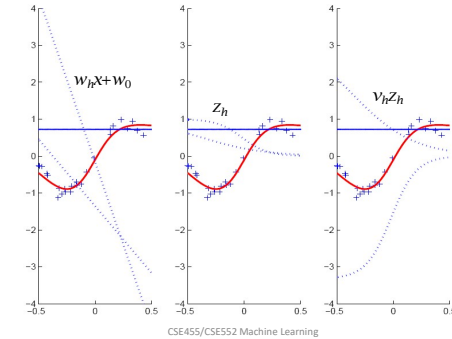
28

## Slide 29

### Results

29

## Slide 30

### Results



$w_h x + w_0$

$z_h$

$v_h z_h$

30

## Slide 31

### Two-Class Discrimination

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

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

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

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

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

31

## Slide 32

### K>2 Classes

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

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

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

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

32

## Multiple Hidden Layers

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

$$z_{1h} = \text{sigmoid}(\boldsymbol{w}_{1h}^T \boldsymbol{x}) = \text{sigmoid}\left(\sum_{j=1}^{d} w_{1hj} x_j + w_{1h0}\right), h = 1,\ldots,H_1$$
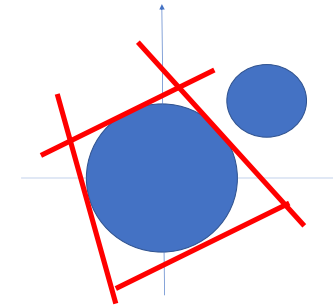
$$z_{2l} = \text{sigmoid}(\boldsymbol{w}_{2l}^T \boldsymbol{z}_1) = \text{sigmoid}\left(\sum_{h=1}^{H_1} w_{2lh} z_{1h} + w_{2l0}\right), l = 1,\ldots,H_2$$

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

33

34

## Improving Convergence

- Momentum

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

- Adaptive learning rate

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

35

## Overfitting/Overtraining

Number of weights: $H(d+1)+(H+1)K$

36

37

## Structured MLP



(Le Cun et al, 1989)

38

## Weight Sharing

39

## Hints

- Invariance to translation, rotation, size     (Abu-Mostafa, 1995)



- Virtual examples
- Augmented error: $E'=E+\lambda_h E_h$

If $\boldsymbol{x'}$ and $\boldsymbol{x}$ are the "same": $E_h=[g(x|\theta)- g(x'|\theta)]^2$

Approximation hint: 
$$E_h = \begin{cases} 0 & \text{if } g(x|\theta) \in [a_x, b_x] \\ (g(x|\theta)-a_x)^2 & \text{if } g(x|\theta) < a_x \\ (g(x|\theta)-b_x)^2 & \text{if } g(x|\theta) > b_x \end{cases}$$

40

## Slide 41

# Tuning the Network Size

- Destructive
- Weight decay:

- Constructive
- Growing networks

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i} - \lambda w_i$$

$$E' = E + \frac{\lambda}{2}\sum_i w_i^2$$



*Dynamic Node Creation*      *Cascade Correlation*

(Ash, 1989)    (Fahlman and Lebiere, 1989)

41

## Slide 42

# Bayesian Learning

- Consider weights $w_i$ as random vars, prior $p(w_i)$

$$p(w \mid X) = \frac{p(X \mid w)p(w)}{p(X)} \quad \hat{w}_{MAP} = \arg\max_w \log p(w \mid X)$$

$$\log p(w \mid X) = \log p(X \mid w) + \log p(w) + C$$

$$p(w) = \prod_i p(w_i) \text{ where } p(w_i) = c \cdot \exp\left[-\frac{w_i^2}{2(1/2\lambda)}\right]$$
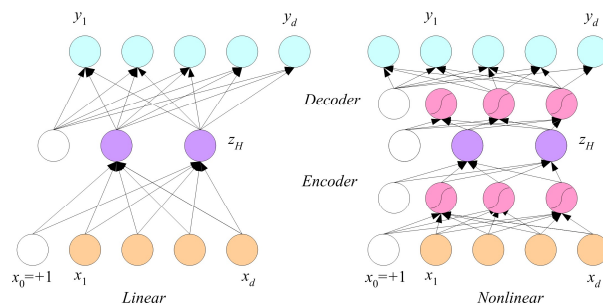
$$E' = E + \lambda \|w\|^2$$

- Weight decay, ridge regression, regularization
  cost=data-misfit + λ complexity

42

## Slide 43

# Dimensionality Reduction



*Linear*      *Nonlinear*

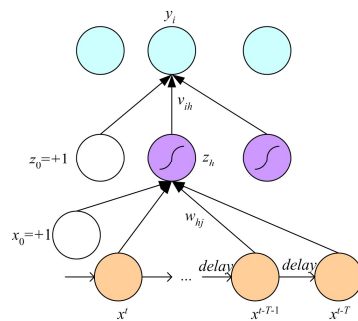43

## Slide 44

44

45

# Learning Time

- Applications:
  - Sequence recognition: Speech recognition
  - Sequence reproduction: Time-series prediction
  - Sequence association
- Network architectures
  - Time-delay networks (Waibel et al., 1989)
  - Recurrent networks (Rumelhart et al., 1986)

46

# Time-Delay Neural Networks

47

# Recurrent Networks

48

## Unfolding in Time

49

Thanks for listening!

50