# On the History and Present of Neural Networks

Thomas Stibor

GSI
Helmholtzzentrum für Schwerionenforschung GmbH

t.stibor@gsi.de

29$^{th}$ October 2024
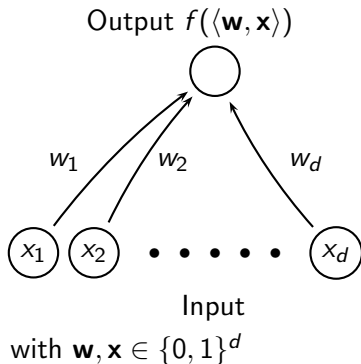
# McCulloch-Pitts Neuron

A LOGICAL CALCULUS OF THE
IDEAS IMMANENT IN NERVOUS ACTIVITY

WARREN S. MCCULLOCH AND WALTER PITTS

FROM THE UNIVERSITY OF ILLINOIS, COLLEGE OF MEDICINE,
DEPARTMENT OF PSYCHIATRY AT THE ILLINOIS NEUROPSYCHIATRIC INSTITUTE,
AND THE UNIVERSITY OF CHICAGO

Because of the "all-or-none" character of nervous activity, neural
events and the relations among them can be treated by means of propo-
sitional logic. It is found that the behavior of every net can be described
in these terms, with the addition of more complicated logical means for
nets containing circles; and that for any logical expression satisfying
certain conditions, one can find a net behaving in the fashion it describes.
It is shown that many particular choices among possible neurophysiologi-
cal assumptions are equivalent, in the sense that for every net behav-
ing under one assumption, there exists another net which behaves un-
der the other and gives the same results, although perhaps not in the
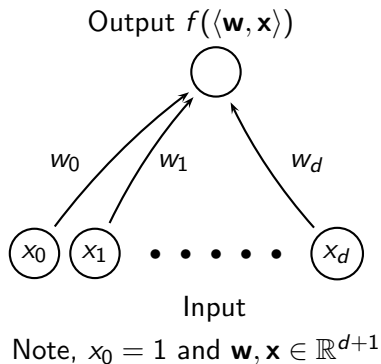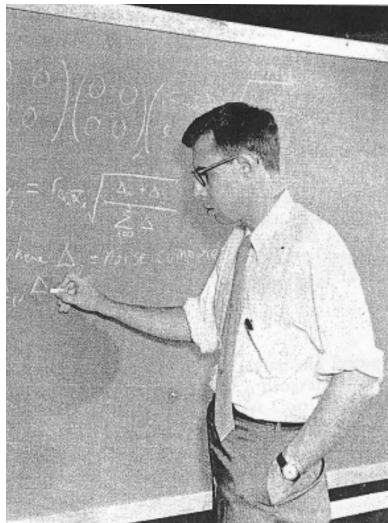same time. Various applications of the calculus are discussed.

Output $f(\langle \mathbf{w}, \mathbf{x} \rangle)$

$w_1$  $w_2$  $w_d$

$x_1$  $x_2$  $\bullet \bullet \bullet \bullet \bullet$  $x_d$

Input
with $\mathbf{w}, \mathbf{x} \in \{0, 1\}^d$

In algebraic notation: $w_1 x_1 + w_2 x_2 + \ldots + w_d x_d = \mathbf{w}^T \cdot \mathbf{x} \stackrel{def}{=} \langle \mathbf{w}, \mathbf{x} \rangle$, and
threshold $\Theta$ step function,
$$f(a) = \begin{cases} 1 & \text{if } \langle \mathbf{w}, \mathbf{x} \rangle \geq \Theta \\ 0 & \text{otherwise} \end{cases}$$

The McCulloch-Pitts neurons represent basic logical functions like *AND*,
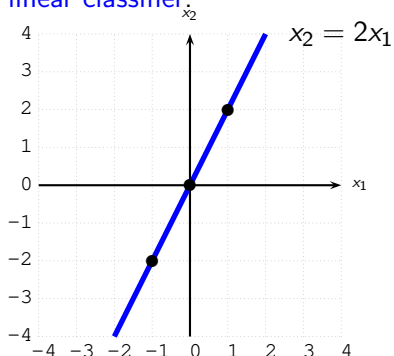*OR*, and *NOT*, but doesn't have a mechanism to learn the weights $\mathbf{w}$.

# Rosenblatt's Perceptron



Output $f(\langle \mathbf{w}, \mathbf{x} \rangle)$

$w_0$    $w_1$    $w_d$

$x_0$   $x_1$   • • • • •   $x_d$

Input

Note, $x_0 = 1$ and $\mathbf{w}, \mathbf{x} \in \mathbb{R}^{d+1}$

Proposed a *learning rule* to infer the weights values from training data.

# Linear Classifier

Rosenblatt's Perceptron (also called single layer neural networks) is a linear classifier.
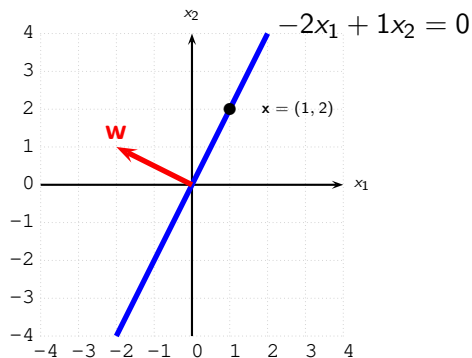


Observe, that $x_2 = 2x_1$ can also be expressed as

$$w_1 x_1 + w_2 x_2 = 0 \Leftrightarrow x_2 = -\frac{w_1}{w_2} x_1,$$

where for instance

$$w_1 = -2, \; w_2 = 1.$$

Furthermore, observe that all points lying on the line $x_2 = 2x_1$ satisfy $w_1 x_1 + w_2 x_2 = -2x_1 + 1x_2 = 0$.
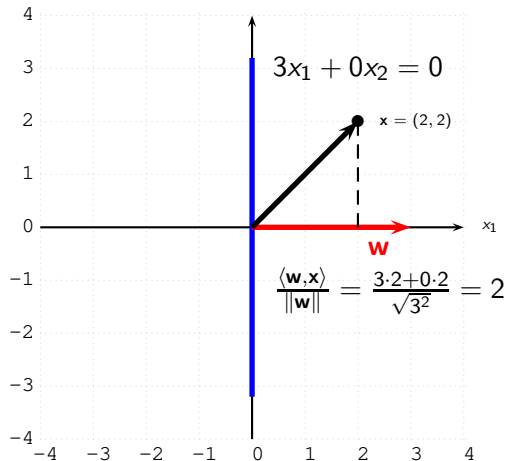
# Linear Classifier & Dot Product



$$-2x_1 + 1x_2 = 0$$

- What about the vector $\mathbf{w} = (w_1, w_2) = (-2, 1)$?
- Vector $\mathbf{w}$ is perpendicular to the line $-2x_1 + 1x_2 = 0$.
- Let us calculate the dot product of $\mathbf{w}$ and $\mathbf{x}$.
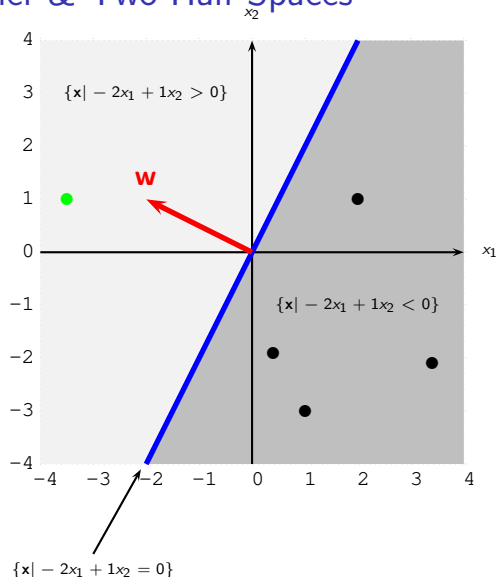
In our example $d = 2$ and we obtain $-2 \cdot 1 + 1 \cdot 2 = 0$.

## Linear Classifier & Dot Product (cont.)

Let us consider the *weight* vector $\mathbf{w} = (3, 0)$ and vector $\mathbf{x} = (2, 2)$.



$$3x_1 + 0x_2 = 0$$

$\mathbf{x} = (2, 2)$

$\mathbf{w}$

$x_1$

$$\frac{\langle \mathbf{w}, \mathbf{x} \rangle}{\|\mathbf{w}\|} = \frac{3 \cdot 2 + 0 \cdot 2}{\sqrt{3^2}} = 2$$

Geometric interpretation of the dot product: Length of the projection of $\mathbf{x}$ onto the unit vector $\mathbf{w}/\|\mathbf{w}\|$.
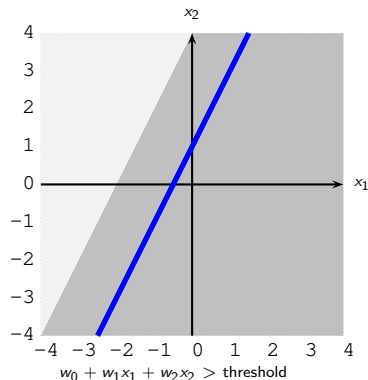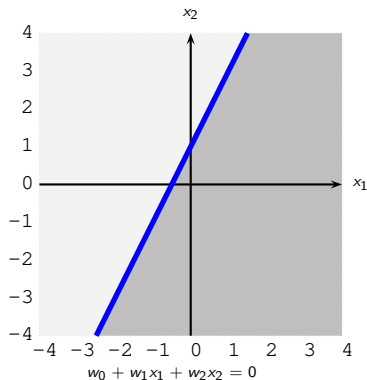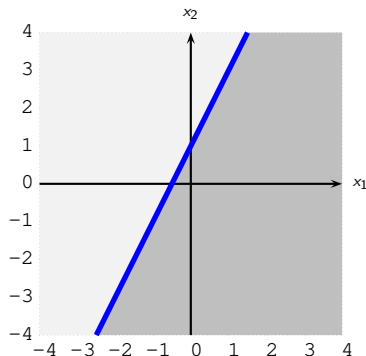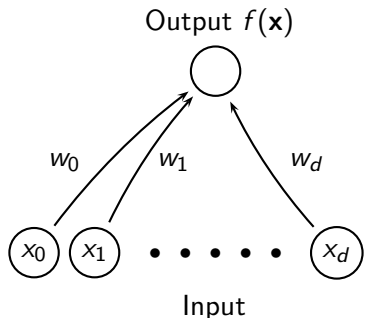
# Linear Classifier & Two Half-Spaces



The **x**-space is separated in two half-spaces.

# Linear Classifier & Dot Product (cont.)

- Observe, that $w_1 x_1 + w_2 x_2 = 0$ implies, that the separating line always goes through the origin.

- By adding an offset (bias), that is
  $w_0 + w_1 x_1 + w_2 x_2 = 0 \Leftrightarrow x_2 = -\frac{w_1}{w_2} x_1 - \frac{w_0}{w_2} \equiv y = mx + b$, one can shift the line arbitrary.

# Linear Classifier & Single Layer NN

Output $f(\mathbf{x})$



$\Leftrightarrow$

Note that $x_0 = 1$, $f(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle$.

Given data which we want to separate, that is, a sample
$\mathcal{X} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \ldots, (\mathbf{x}_N, y_N)\} \in \mathbb{R}^{d+1} \times \{-1, +1\}$.

How to determine the proper values of $\mathbf{w}$ such that the "minus" and "plus" points are separated by $f(\mathbf{x})$? Infer the values of $\mathbf{w}$ from the data by some learning algorithm.

# Perceptron Learning Algorithm

---

**input** : $(\mathbf{x}_1, y_1), \ldots, (\mathbf{x}_N, y_N) \in \mathbb{R}^{d+1} \times \{-1, +1\}, \eta \in$
  $\mathbb{R}_+, \text{max.epoch} \in \mathbb{N}$
**output**: $\mathbf{w}$
**begin**

    Randomly initialize $\mathbf{w}$ ;
    epoch $\leftarrow 0$ ;
    **repeat**
        **for** $i \leftarrow 1$ **to** $N$ **do**
            **if** $y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \leq 0$ **then**
               $\mathbf{w} \leftarrow \mathbf{w} + \eta \mathbf{x}_i \, y_i$
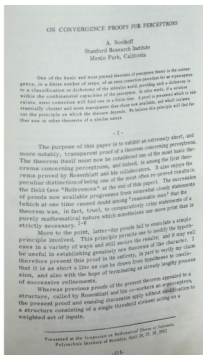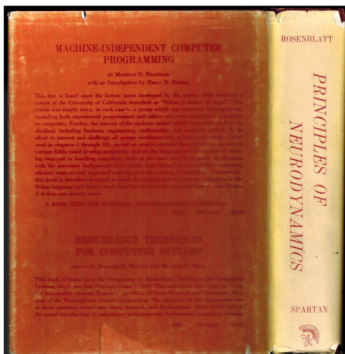
        epoch $\leftarrow$ epoch $+ 1$
    **until** (*epoch* = *max.epoch*) *or* (*no change in* $\mathbf{w}$);
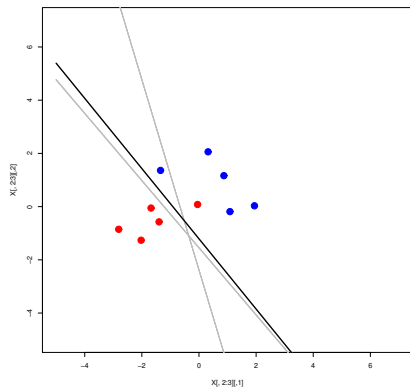    **return** $\mathbf{w}$

---

# Perceptron Convergence Theorem

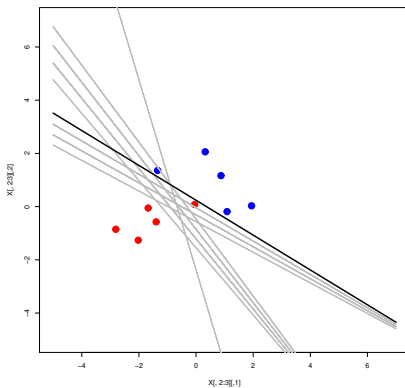How often one has to cycle through the patterns in the training set?

- If the training data is linearly separable, the perceptron learning algorithm will converge after a finite number of iterations, meaning it will find a set of weights that perfectly classify the data.
- If the data is not linearly separable, the perceptron will not converge and will continue updating its weights indefinitely.

# Perceptron Algorithm Visualization



One epoch                    terminate if no change in **w**

# From Perceptron Loss$_\Theta$ to Gradient Descent

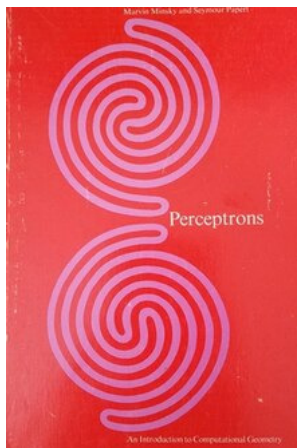The parameters to learn are: $(w_0, w_1, w_2) = \mathbf{w}$.

- What is our loss function Loss$_\Theta$ we would like to minimize?
- Where is term $\mathbf{w}_{\text{new}} = \mathbf{w} + \eta \mathbf{x}\, y$ coming from?

$$\text{Loss}_\Theta \widehat{=} E(\mathbf{w}) = -\sum_{m \in \mathcal{M}} \langle \mathbf{w}, \mathbf{x}_m \rangle y_m$$

where $\mathcal{M}$ denotes the set of all missclassified patterns. Moreover, Loss$_\Theta$ is *continuous* and *piecewise linear* and fits in the spirit iterative *gradient descent* method

$$\mathbf{w}_{\text{new}} = \mathbf{w} + \eta \nabla E(\mathbf{w}) = \mathbf{w} + \eta \mathbf{x}\, y$$

# The Neural Network Winter



Perceptrons: An Introduction to Computational Geometry. Marvin Minsky and Seymour Papert, 1969.

- Analyzed the capabilities and limitations of the single-layer perceptron.
- Proved that single-layer perceptrons are fundamentally limited in their ability to solve non-linearly separable problems, such as the XOR problem.
- AI shifted their focus to other methods, particularly symbolic AI and rule-based systems.
- Funding agencies and academic institutions also deprioritized neural network research (dead-end field).

# Hopfield Network

## Neural networks and physical systems with emergent collective computational abilities
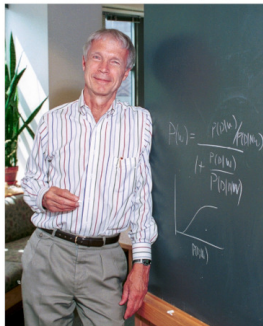
(associative memory/parallel processing/categorization/content-addressable memory/fail-soft devices)

J. J. HOPFIELD

Division of Chemistry and Biology, California Institute of Technology, Pasadena, California 91125; and Bell Laboratories, Murray Hill, New Jersey 07974

**ABSTRACT**     Computational properties of use to biological organisms or to the construction of computers can emerge as collective properties of systems having a large number of simple equivalent components (or neurons). The physical meaning of content-addressable memory is described by an appropriate phase space flow of the state of a system. A model of such a system is given, based on aspects of neurobiology but readily adapted to integrated circuits. The collective properties of this model produce a content-addressable memory which correctly yields an entire memory from any subpart of sufficient size. The algorithm for the time evolution of the state of the system is based on asynchronous parallel processing. Additional emergent collective properties include some capacity for generalization, familiarity recognition, categorization, error correction, and time sequence retention. The collective properties are only weakly sensitive to details of the modeling or the failure of individual devices.
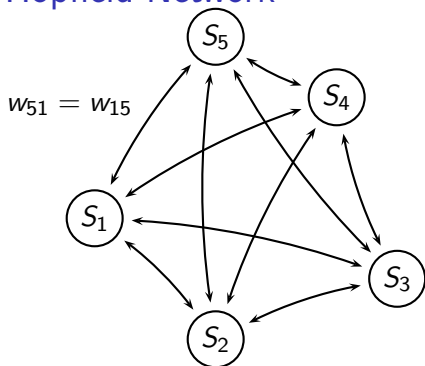
# Hopfield Network Introductory Example

recalled by the memory



- Suppose we want to store $N$ binary images in some memory.

- The memory should be *content-addressable* and insensitive to small errors.

- We present corrupted images to the memory (e.g. our brain) and recall the corresponding images.

presentation of corrupted images

# Hopfield Network



- $w_{ij}$ denotes weight connection from unit $j$ to unit $i$
- no unit has connection with itself $w_{ii} = 0$, $\forall i$
- connections are symmetric $w_{ij} = w_{ji}$, $\forall i, j$

State of unit $i$ can take values $\pm 1$ and is denoted as $S_i$. State dynamics are governed by activity rule:

$$S_i = \text{sgn}\left(\sum_j w_{ij} S_j\right), \text{ where } \text{sgn}(a) = \begin{cases} +1 & \text{if } a \geq 0, \\ -1 & \text{if } a < 0 \end{cases}$$

# Learning Rule in a Hopfield Network

Learning in Hopfield networks:

- Store a set of desired memories $\{\mathbf{x}^{(n)}\}$ in the network, where each memory is a binary pattern with $x_i \in \{-1, +1\}$.
- The weights are set using the sum of outer products

$$w_{ij} = \frac{1}{N} \sum_n x_i^{(n)} x_j^{(n)},$$

where $N$ denotes the number of units ($N$ can also be some positive constant, e.g. number of patterns). Given a $m \times 1$ column vector $\mathbf{a}$ and $1 \times n$ row vector $\mathbf{b}$. The outer product $\mathbf{a} \otimes \mathbf{b}$ is defined as the $m \times n$ matrix.

$$
\left[
\begin{array}{c}
a_1 \\
a_2 \\
a_3
\end{array}
\right]
\otimes
[b_1 \; b_2 \; b_3]
=
\left[
\begin{array}{ccc}
a_1 b_1 & a_1 b_2 & a_1 b_3 \\
a_2 b_1 & a_2 b_2 & a_2 b_3 \\
a_3 b_1 & a_3 b_2 & a_3 b_3
\end{array}
\right],
\quad m = n = 3
$$

# Learning in Hopfield Network (Example)

Suppose we want to store patterns $\mathbf{x}^{(1)} = [-1, +1, -1]$ and $\mathbf{x}^{(2)} = [+1, -1, +1]$.

$$\begin{bmatrix} -1 \\ +1 \\ -1 \end{bmatrix} \otimes [-1, +1, -1] = \begin{bmatrix} +1 & -1 & +1 \\ -1 & +1 & -1 \\ +1 & -1 & +1 \end{bmatrix}$$

$$+$$

$$\begin{bmatrix} +1 \\ -1 \\ +1 \end{bmatrix} \otimes [+1, -1, +1] = \begin{bmatrix} +1 & -1 & +1 \\ -1 & +1 & -1 \\ +1 & -1 & +1 \end{bmatrix}$$

# Learning in Hopfield Netw. (Example) (cont.)

$$\mathbf{W} = \frac{1}{3} \begin{bmatrix} 0 & -2 & +2 \\ -2 & 0 & -2 \\ +2 & -2 & 0 \end{bmatrix}$$

Recall: no unit has connection with itself.

The storage of patterns in the network can also be interpreted as constructing stable states. The condition for patterns to be stable is:

$$\text{sgn}\left(\sum_j w_{ij} x_i\right) = x_i\,, \forall i.$$

Suppose we present pattern $\mathbf{x}^{(1)}$ to the network and want to restore the corresponding pattern.

# Learning in Hopfield Netw. (Example) (cont.)

Let us assume that the network states are set as follows: $S_i = x_i, \ \forall i$. We can restore pattern $\mathbf{x}^{(1)} = [-1, +1, -1]$ as follows:

$$S_1 = \text{sgn}\left(\sum_{j=1}^{3} w_{1j} S_j\right) = -1 \qquad S_2 = \text{sgn}\left(\sum_{j=1}^{3} w_{2j} S_j\right) = +1$$

$$S_3 = \text{sgn}\left(\sum_{j=1}^{3} w_{3j} S_j\right) = -1$$

Can we also restore the original patterns by presenting "similar" patterns which are corrupted by noise?

# Updating States in a Hopfield Network

**Synchronous updates:**

- all units update their states $S_i = \text{sgn}\left(\sum_j w_{ij} S_j\right)$ simultaneously.

**Asynchronous updates:**

- one unit at a time updates its state. The sequence of selected units may be a fixed sequence or a random sequence.

Synchronously updating states can lead to oscillation (no convergence to a stable state).

# Aim of a Hopfield Network

> Our aim is that by presenting a corrupted pattern, and by applying iteratively the state update rule the Hopfield network will settle down in a stable state which corresponds to the desired pattern.

Hopfield network is a method for

- pattern completion
- error correction.

The state of a Hopfield network can be expressed in terms of the energy function (related to Ising model and spin glass theory in Physics).

$$E = -\frac{1}{2} \sum_{i,j} w_{ij} S_i S_j$$

Hopfield observed that if a state is a local minimum in the energy function, it is also a stable state for the network.

# Basin of Attraction and Stable States



Within the space the stored patterns $\mathbf{x}^{(n)}$ are acting like attractors.

# Haykin's Digit Example

Suppose we stored the following digits in the Hopfield network:



Energy = –67.73 — Pattern 0
Energy = –67.87 — Pattern 1
Energy = –82.33 — Pattern 2
Energy = –86.6 — Pattern 3
Energy = –77.73 — Pattern 4
Energy = –90.47 — Pattern 6
Energy = –83.13 — Pattern 9
Energy = –66.93 — Pattern box

# Updated States of Corrupted Digit 6

# Updated States of Corrupted Digit 6 (cont.)



Energy = −36.73 — updated unit 113
Energy = −38.4 — updated unit 57
Energy = −41.07 — updated unit 103
Energy = −42.4 — updated unit 18
Energy = −45.27 — updated unit 109

Energy = −47.6 — updated unit 83
Energy = −50.4 — updated unit 71
Energy = −52.67 — updated unit 77
Energy = −56.47 — updated unit 26
Energy = −58.4 — updated unit 15

Energy = −60.67 — updated unit 31
Energy = −63.33 — updated unit 58
Energy = −64.47 — updated unit 16
Energy = −68 — updated unit 29
Energy = −71.27 — updated unit 88

# Updated States of Corrupted Digit 6 (cont.)

The resulting pattern (stable state with energy $-90.47$) matches the desired pattern.



Energy = –73.73
updated unit 72

Energy = –77.27
updated unit 90

Energy = –81.47
updated unit 19

Energy = –84.27
updated unit 21

Energy = –87.33
updated unit 25

Energy = –90.47
updated unit 73

Energy = –90.47
Original Pattern 6

# Hopfield Networks Summary

- Learning: determine the weight matrix from the data with the outer product.
- Memory: "knowledge" is stored in the weight matrix.
- Queries to memory: apply state update rule until energy is minimized (local minimum).

John Hopfield laid the groundwork for:

- Renewed theoretical interest and connections to Physics.
- Neural network applications for optimization.
- Paved the way for recurrent networks.
- Revival of interest in multilayer networks and backpropagation.

# Backpropagation the Heart of Neural Networks

Adjust weights of connections within the network to minimize the error between the predicted and actual output.

History:

- The minimisation of errors through gradient descent (Cauchy 1847).

- ...

- Taylor Expansion of the Accumulated Rounding Error (Seppo Linnainmaa 1970 Master Thesis, backpropagation modern version).

- ...

- Learning representations by back-propagating errors (Rumelhart, Hinton and Williams, 1986).

Who Invented Backpropagation? (Excellent article by Jürgen Schmidhuber).

# Learning in Neural Networks with Backpropagation



minimize $\frac{1}{2}\|f(\mathbf{W}^{(3)}f(\mathbf{W}^{(2)}f(\mathbf{W}^{(1)}\mathbf{X} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)}) + \mathbf{b}^{(3)}) - \mathbf{Y}\|^2$

parameters to fit

$\mathbf{W}^{(3)}, \mathbf{b}^{(3)}$

$\mathbf{W}^{(2)}, \mathbf{b}^{(2)}$

$\mathbf{W}^{(1)}, \mathbf{b}^{(1)}$

Core idea:

- Calculate error of loss function and change weights and biases based on output.

- These "error" measurements for each unit can be used to calculate the partial derivatives.

- Use partial derivatives with gradient descent for updating weights and biases and minimizing loss function.

Problem: At which magnitude one shall change e.g. weight $W_{ij}^{(1)}$ based on error of $y_2$?

Input: $x_1, x_2$, output: $a_1^{(3)}, a_2^{(3)}$, target: $y_1, y_2$ and $g(\cdot)$ is activation function. NN calculates[2] $g(\mathbf{W}^{(2)} g(\mathbf{W}^{(1)} \mathbf{x}))$.

$$E(\mathbf{W}) = \frac{1}{2}\left[(a_1^{(3)} - y_1)^2 + (a_2^{(3)} - y_2)^2\right] = \frac{1}{2}\|\mathbf{a}^{(3)} - \mathbf{y}\|^2$$



$L_3$

$\mathbf{W}^{(2)}$

$L_2$

$\mathbf{W}^{(1)}$

$L_1$

$z_1^{(3)} = W_{10}^{(2)} a_0^{(2)} + W_{11}^{(2)} a_1^{(2)} + W_{12}^{(2)} a_2^{(2)} + W_{13}^{(2)} a_3^{(2)} \qquad a_1^{(3)} = g(z_1^{(2)})$

$z_2^{(3)} = W_{20}^{(2)} a_0^{(2)} + W_{21}^{(2)} a_1^{(2)} + W_{22}^{(2)} a_2^{(2)} + W_{23}^{(2)} a_3^{(2)} \qquad a_2^{(3)} = g(z_2^{(2)})$

$\underbrace{\mathbf{z}^{(3)}}_{2 \times 1} = \underbrace{\mathbf{W}^{(2)}}_{2 \times 4} \underbrace{\mathbf{a}^{(2)}}_{4 \times 1} \qquad \mathbf{a}^{(3)} = g(\mathbf{z}^{(3)})$

Forward pass

$z_1^{(2)} = W_{10}^{(1)} x_0 + W_{11}^{(1)} x_1 + W_{12}^{(1)} x_2 \qquad a_1^{(2)} = g(z_1^{(2)})$

$z_2^{(2)} = W_{20}^{(1)} x_0 + W_{21}^{(1)} x_1 + W_{22}^{(1)} x_2 \qquad a_2^{(2)} = g(z_2^{(2)})$

$z_3^{(2)} = W_{30}^{(1)} x_0 + W_{31}^{(1)} x_1 + W_{32}^{(1)} x_2 \qquad a_3^{(2)} = g(z_3^{(2)})$

$\underbrace{\mathbf{z}^{(2)}}_{3 \times 1} = \underbrace{\mathbf{W}^{(1)}}_{3 \times 3} \underbrace{\mathbf{x}}_{3 \times 1} \qquad \mathbf{a}^{(2)} = g(\mathbf{z}^{(2)})$

[2]Notation adapted from Andew Ng's slides.

# Learning in Neural Networks with Backpropagation (cont.)

For each node we calculate $\delta_j^{(l)}$, that is, error of unit $j$ in layer $l$, because $\frac{\partial}{\partial W_{ij}^{(l)}} E(\mathbf{W}) = a_j^{(l)} \delta_i^{(l+1)}$. Note $\odot$ is element wise multiplication.

$$E(\mathbf{W}) = \tfrac{1}{2}\left[(a_1^{(3)} - y_1)^2 + (a_2^{(3)} - y_2)^2\right] = \tfrac{1}{2}\|\mathbf{a}^{(3)} - \mathbf{y}\|^2$$

$L_3$

$\mathbf{W}^{(2)}$

$L_2$

$\mathbf{W}^{(1)}$

$L_1$



$\boldsymbol{\delta}^{(3)} = (\mathbf{a}^{(3)} - \mathbf{y}) \odot g'(\mathbf{z}^{(3)})$

$\boldsymbol{\delta}^{(2)} = (\mathbf{W}^{(2)})^T \delta^{(3)} \odot g'(\mathbf{z}^{(2)})$

Note $\boldsymbol{\delta}^{(1)}$ is the input, so no term.

Backward pass

# Learning in Neural Networks with Backpropagation (cont.)

*Backpropagation* = forward pass & backward pass

Given labeled training data $(\mathbf{x}_1, \mathbf{y}_1), \ldots, (\mathbf{x}_N, \mathbf{y}_N)$.

Set $\Delta_{ij}^{(l)} = 0$ for all $l, i, j$. Value $\Delta$ will be used as accumulators for computing partial derivatives.

For $n = 1$ to $N$

- Forward pass, compute $\mathbf{z}^{(2)}, \mathbf{a}^{(2)}, \mathbf{z}^{(3)}, \mathbf{a}^{(3)}, \ldots, \mathbf{z}^{(L)}, \mathbf{a}^{(L)}$
- Backward pass, compute $\boldsymbol{\delta}^{(L)}, \boldsymbol{\delta}^{(L-1)}, \ldots, \boldsymbol{\delta}^{(2)}$
- Accumulate partial derivate terms, $\boldsymbol{\Delta}^{(l)} := \boldsymbol{\Delta}^{(l)} + \boldsymbol{\delta}^{(l+1)} (\mathbf{a}^{(l)})^T$

Finally calculated partial derivatives for each parameter:
$\frac{\partial}{\partial W_{ij}^{(l)}} E(\mathbf{W}) = \frac{1}{N} \Delta_{ij}^{(l)}$ and use these in gradient descent.

See interactive demo.

# Neural Networks vs. Kernel Methods (1995 - 2012)

Support Vector Machines and Kernel Methods were favorite methods in the field of machine learning.

Neural networks suffered from:

- Slow training time: Took usually weeks and made experimentation and tuning difficult.
- Vanishing and exploding gradient problem: Especially severe with sigmoid and tanh activation functions.
- Lack of labeled data sets.
- Lack of neural network frameworks (TensorFlow, PyTorch, MXNet, etc...): Usually Matlab code.

# Era of Deep Learning Neural Networks

# Why are Deep Neural Networks so successful?



Deep Neural Networks (Backpropagation) are *universal*, that is, applicable to a large class of problems: Vision, speech, text, ... and *scale* with data. Backpropagation (forward + backward pass) is intrinsically linked to matrix multiplication (GPU's, TPU's).