# Neural Networks: Old and New

**Ju Sun**

Computer Science & Engineering

University of Minnesota, Twin Cities

September 6, 2023
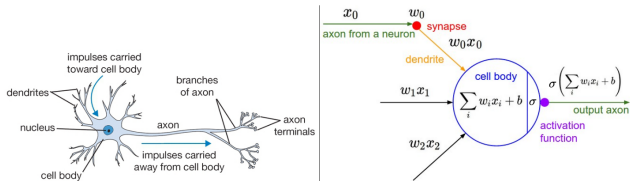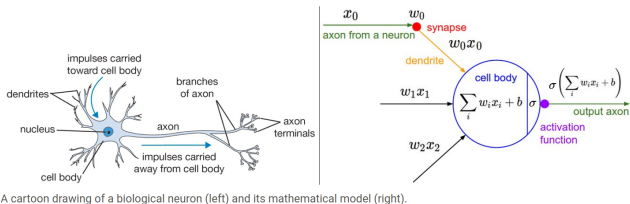
## Model of biological neurons



A cartoon drawing of a biological neuron (left) and its mathematical model (right).

Credit: Stanford CS231N

Biologically ...

– Each neuron receives signals from its **dendrites**

– Each neuron outputs signals via its single **axon**

– The axon branches out and connects via **synapese** to dendrites of other neurons

# Model of biological neurons



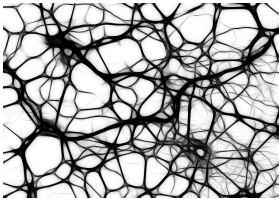A cartoon drawing of a biological neuron (left) and its mathematical model (right).

Credit: Stanford CS231N

Mathematically ...

- Each neuron receives $x_i$'s from its **dendrites**

- $x_i$'s weighted by $w_i$'s (synaptic strengths) and summed $\sum_i w_i x_i$

- The neuron fires only when the combined signal is above a certain threshold: $\sum_i w_i x_i + b$

- Fire rate is modeled by an **activation function** $\sigma$, i.e., outputting $\sigma \left( \sum_i w_i x_i + b \right)$
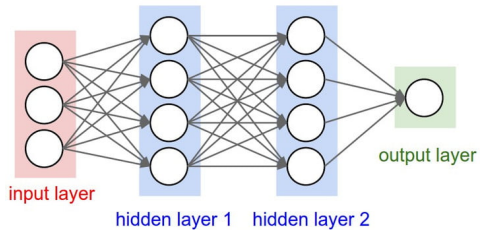
## Artificial neural networks

Brain neural networks



~ 86-billion neurons (Credit: Max Pixel)

Artificial neural networks



Why called **artificial**?

- – (Over-)simplification on neural level
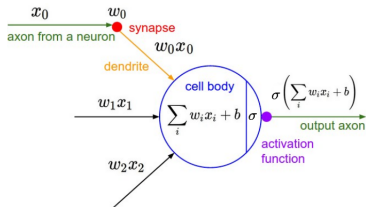- – (Over-)simplification on connection level

> In this course, neural networks are always artificial.

## Examples of activation function $\sigma$

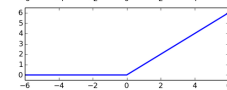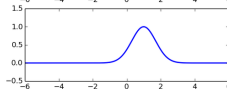

Sigmoid
$$\phi(z) = \frac{1}{1 + e^{-z}}$$

Hyperbolic Tangent
$$\phi(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

Rectified Linear
$$\phi(z) = \begin{cases} 0 & \text{if } z < 0 \\ z & \text{if } z \geq 0 \end{cases}$$

Radial Basis Function
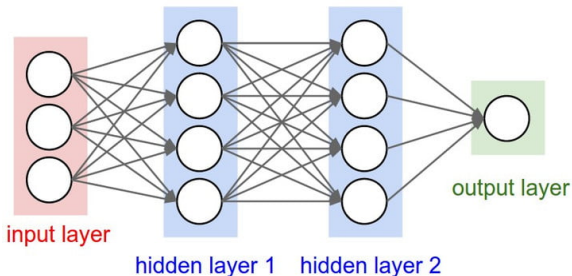$$\phi(z, c) = e^{-(\epsilon \|z - c\|)^2}$$

Credit: [Hughes and Correll, 2016]

$$\sigma \left( \sum_i w_i x_i + b \right) = \sigma \left( \boldsymbol{w}^\mathsf{T} \boldsymbol{x} + b \right)$$

## Neural networks

One neuron: $\sigma\left(\boldsymbol{w}^{\mathsf{T}}\boldsymbol{x} + b\right)$

Neural networks (NN): **structured** organization of artificial neurons



input layer

hidden layer 1    hidden layer 2

output layer

$\boldsymbol{w}$'s and $b$'s are unknown and need to be learned

Many models in machine learning **are** neural networks

**Supervised Learning**

- Gather training data $(\boldsymbol{x}_1, \boldsymbol{y}_1), \ldots, (\boldsymbol{x}_n, \boldsymbol{y}_n)$
- Choose a family of functions, e.g., $\mathcal{H}$, so that there is $f \in \mathcal{H}$ to ensure $\boldsymbol{y}_i \approx f(\boldsymbol{x}_i)$ for all $i$
- Set up a loss function $\ell$ to measure the approximation quality
- Find an $f \in \mathcal{H}$ to minimize the average loss

$$\min_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^{n} \ell(\boldsymbol{y}_i, f(\boldsymbol{x}_i))$$

... known as **empirical risk minimization** (ERM) framework in learning theory
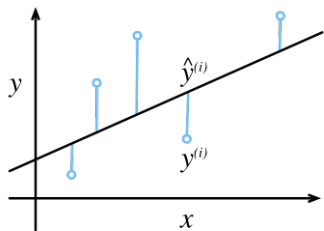
# Supervised learning meets NNs

**Supervised Learning from NN viewpoint**

- Gather training data $(\boldsymbol{x}_1, \boldsymbol{y}_1), \ldots, (\boldsymbol{x}_n, \boldsymbol{y}_n)$

- Choose a NN with $k$ neurons, so that there is a group of weights, e.g., $(\boldsymbol{w}_1, \ldots, \boldsymbol{w}_k, b_1, \ldots, b_k)$, to ensure

$$\boldsymbol{y}_i \approx \{\mathsf{NN}\,(\boldsymbol{w}_1, \ldots, \boldsymbol{w}_k, b_1, \ldots, b_k)\}\,(\boldsymbol{x}_i) \quad \forall i$$

- Set up a loss function $\ell$ to measure the approximation quality

- Find weights $(\boldsymbol{w}_1, \ldots, \boldsymbol{w}_k, b_1, \ldots, b_k)$ to minimize the average loss

$$\min_{\boldsymbol{w}'s, b's} \frac{1}{n} \sum_{i=1}^{n} \ell\,[\boldsymbol{y}_i, \{\mathsf{NN}\,(\boldsymbol{w}_1, \ldots, \boldsymbol{w}_k, b_1, \ldots, b_k)\}\,(\boldsymbol{x}_i)]$$

# Linear regression



Credit: D2L

- Data: $(\boldsymbol{x}_1, y_1), \ldots, (\boldsymbol{x}_n, y_n)$, $\boldsymbol{x}_i \in \mathbb{R}^d$
- Model: $y_i \approx \boldsymbol{w}^\intercal \boldsymbol{x}_i + b$
- Loss: $\|y - \hat{y}\|_2^2$
- Optimization:

$$\min_{\boldsymbol{w}, b} \ \frac{1}{n} \sum_{i=1}^{n} \|y_i - (\boldsymbol{w}^\intercal \boldsymbol{x}_i + b)\|_2^2$$



Credit: D2L

$\sigma$ is the identity function

# Perceptron

**Frank Rosenblatt**

(1928–1971)

– Data: $(\boldsymbol{x}_1, y_1), \ldots, (\boldsymbol{x}_n, y_n)$, $\boldsymbol{x}_i \in \mathbb{R}^d$, $y_i \in \{+1, -1\}$

– Model: $y_i \approx \sigma(\boldsymbol{w}^\mathsf{T} \boldsymbol{x}_i + b)$, $\sigma$ sign function



– Loss: $\mathbf{1}\{y \neq \hat{y}\}$

– Optimization:

$$\min_{\boldsymbol{w}, b} \frac{1}{n} \sum_{i=1}^{n} \mathbf{1}\{y_i \neq \sigma(\boldsymbol{w}^\mathsf{T} \boldsymbol{x}_i + b)\}$$

## Perceptron

Perceptron is a single artificial neuron for **binary classification**



Perceptron rule.

dominated early AI (50's – 70's)

**Logistic regression** is similar but with **sigmoid** activation



$$\phi(z) = \frac{1}{1 + e^{-z}}$$

## Softmax regression

– Data: $(\boldsymbol{x}_1, y_1), \ldots, (\boldsymbol{x}_n, y_n)$, $\boldsymbol{x}_i \in \mathbb{R}^d$, $y_i \in \{L_1, \ldots, L_p\}$, i.e., multiclass classification problem

– Data preprocessing: labels into vectors via one-hot encoding

$$L_k \Longrightarrow [\underbrace{0, \ldots, 0}_{k-1\,0's}, 1, \underbrace{0, \ldots, 0}_{p-k\,0's}]^\mathsf{T}$$

So: $y_i \Longrightarrow \boldsymbol{y}_i$

– Model: $\boldsymbol{y}_i \approx \sigma\left(\boldsymbol{W}^\mathsf{T} \boldsymbol{x}_i + \boldsymbol{b}\right)$, here $\sigma$ is the softmax function (maps vectors to vectors): for $\boldsymbol{z} \in \mathbb{R}^p$,

$$\boldsymbol{z} \mapsto \left[\frac{e^{z_1}}{\sum_j e^{z_j}}, \ldots, \frac{e^{z_p}}{\sum_j e^{z_j}}\right]^\mathsf{T}.$$

– Loss: **cross-entropy loss** $-\sum_j y_j \log \hat{y}_j$

– Optimization ...

# Softmax regression

... for multiclass classification



Credit: D2L

## Multilayer perceptrons (MLP)



Output layer — $o_1$, $o_2$, $o_3$

Hidden layer — $h_1$, $h_2$, $h_3$, $h_4$, $h_5$

Input layer — $x_1$, $x_2$, $x_3$, $x_4$

Credit: D2L

Model: $\boldsymbol{y}_i \approx \sigma_2 \left( \boldsymbol{W}_2^\mathsf{T} \sigma_1 (\boldsymbol{W}_1^\mathsf{T} \boldsymbol{x} + \boldsymbol{b}_1) + \boldsymbol{b}_2 \right)$

Also called **fully-connected networks**

Modern NNs:

- – many hidden layers: deep neural networks (DNNs)
- – refined/structured connection and/or activations
  (convolutional/recurrent/graph/... NNs)

## They're all (shallow) NNs

- Linear regression

- Perception and Logistic regression

- Softmax regression

- Multilayer perceptron (feedforward NNs)

- Support vector machines (SVM)

- PCA (autoencoder)

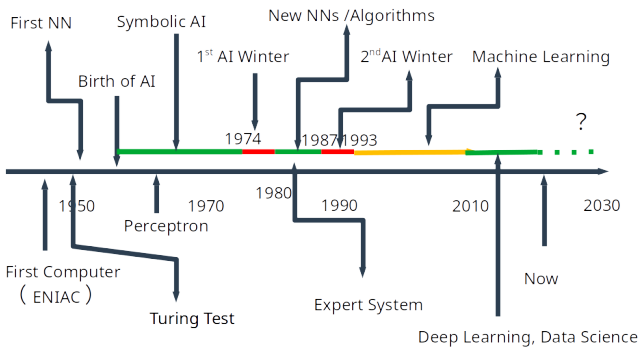- Matrix factorization

see, e.g., Chapter 2 of [Aggarwal, 2018].
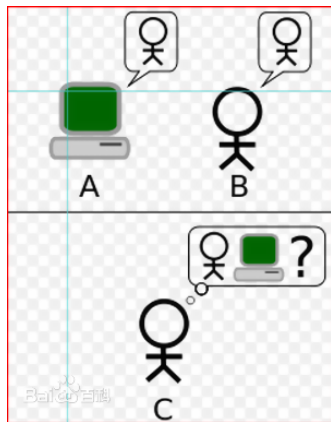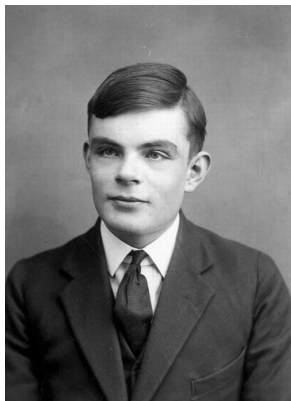
- Crucial precursors: first computer, Turing test
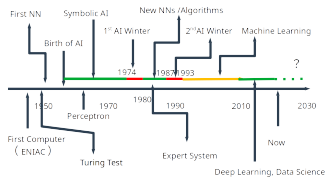- 1956: Dartmouth Artificial Intelligence Summer Research Project — Birth of AI
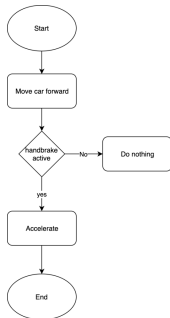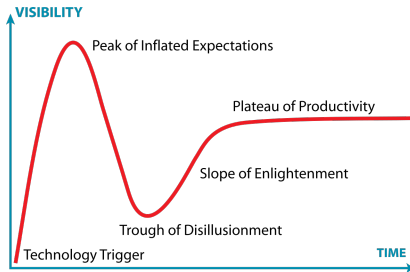
Turing Test



Alan Turing (1912–1954)

## Symbolic AI: modeling general logic and reasoning





rules for recognizing dogs?

First NN  Symbolic AI  New NNs /Algorithms

Birth of AI  1st AI Winter  2nd AI Winter  Machine Learning

1974  1987 1993  ?

1950  1970  1980  1990  2010  2030

Perceptron

First Computer
( ENIAC )  Now

Turing Test  Expert System

Deep Learning, Data Science



VISIBILITY

Peak of Inflated Expectations

Plateau of Productivity

Slope of Enlightenment

Trough of Disillusionment

Technology Trigger  TIME
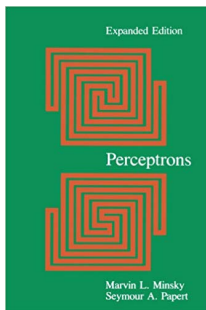
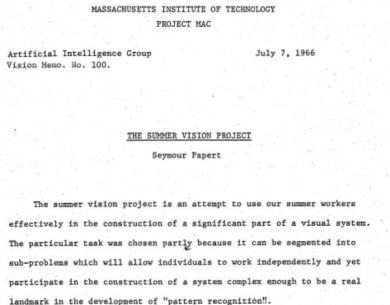Gartner hype cycle

invented 1962



written in 1969, end of
Perceptron era
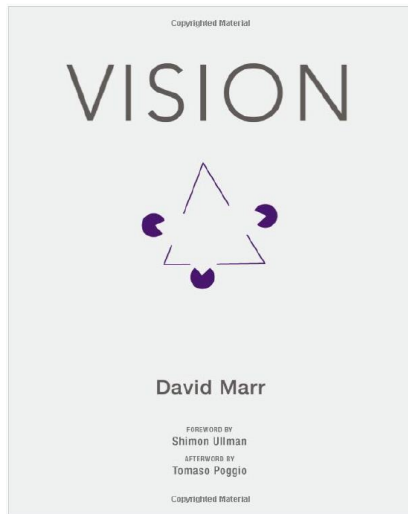


Marvin Minsky (1927–2016)

## Birth of computer vision



1966



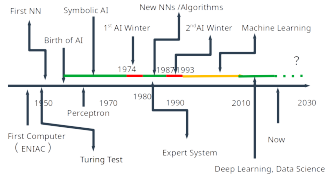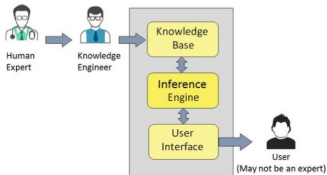around 1980

**expert system—building in domain-specific knowledge**
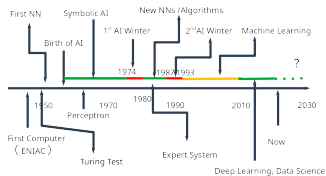




Can we build comprehensive
knowledge bases and know all
rules?

Key ingredients of DL have been in place for 25-30 years:

| Landmark | Emblem | Epoch |
|---|---|---|
| Neocognitron | Fukushima | 1980 |
| CNN | Le Cun | mid 1980s' |
| Backprop | Hinton | mid 1980's |
| SGD | Le Cun, Bengio etc | mid 1990's |
| Various | Schmidhuber | mid 1980's |
| *CTF* | *DARPA etc* | *mid 1980's* |

Machine learning takes over ...



rules learned from data, or **data-driven**

Starting 1990's

Support vector machines (SVM)

Adaboost

Decision trees and random forests

Deep learning (2010's)

...

Start from neurons

Shallow to deep neural networks

A brief history of AI

Suggested reading

## Suggested reading

- Chap 2, Neural Networks and Deep Learning.

- Chap 3–4, Dive into Deep Learning.

- Chap 1, Deep Learning with Python.

[Aggarwal, 2018] Aggarwal, C. C. (2018). **Neural Networks and Deep Learning.**
Springer International Publishing.

[Hughes and Correll, 2016] Hughes, D. and Correll, N. (2016). **Distributed machine learning in materials that couple sensing, actuation, computation and communication.** *arXiv:1606.03508.*