

Mnacho Echenim
small evolutions : Patrick Reignier

Grenoble INP-Ensimag

2025-2026

Outline

General presentation

The history of neural networks

The multilayer perceptron

What do we gain with an MLP?

About this class

- ▶ Deep learning
 - ▶ Very dynamic field of research
 - ▶ Increased computational power
 - ▶ Increased data availability

About this class

► Deep learning

- Very dynamic field of research
 - Increased computational power
 - Increased data availability
- Different kinds of neural networks, depending on applications
 - Convolutional networks for image processing
 - Recursive neural networks, large language models for natural language processing
 - Generative adversarial network for physical simulations and image processing

About this class

- ▶ Deep learning
 - ▶ Very dynamic field of research
 - ▶ Increased computational power
 - ▶ Increased data availability
 - ▶ Different kinds of neural networks, depending on applications
 - ▶ Convolutional networks for image processing
 - ▶ Recursive neural networks, large language models for natural language processing
 - ▶ Generative adversarial network for physical simulations and image processing
 - ▶ But not many impressive results obtained by using a neural network as a black box, with no knowledge on how they work or on the data
 - ▶ Other ML techniques can outperform neural networks, depending on the problem at hand

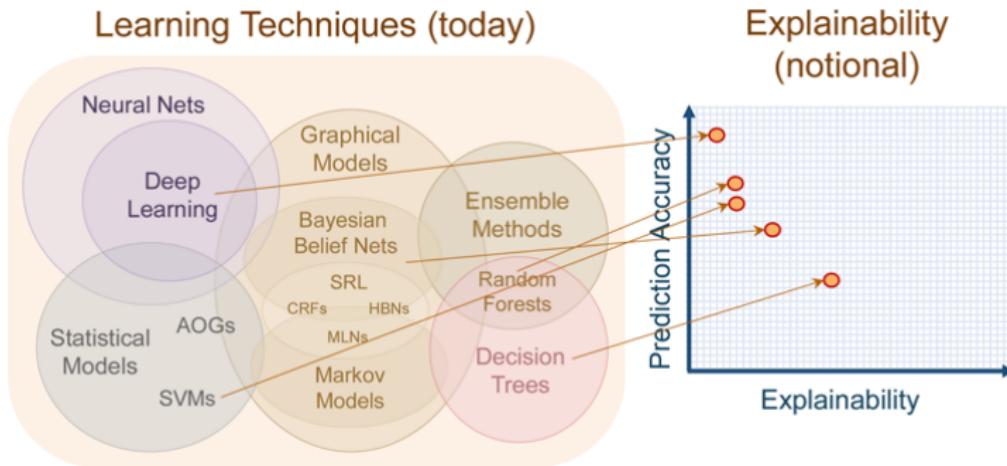
About this class

- ▶ Deep learning
 - ▶ Very dynamic field of research
 - ▶ Increased computational power
 - ▶ Increased data availability
 - ▶ Different kinds of neural networks, depending on applications
 - ▶ Convolutional networks for image processing
 - ▶ Recursive neural networks, large language models for natural language processing
 - ▶ Generative adversarial network for physical simulations and image processing
 - ▶ But not many impressive results obtained by using a neural network as a black box, with no knowledge on how they work or on the data
 - ▶ Other ML techniques can outperform neural networks, depending on the problem at hand
- ▶ Goal of this class
 - ▶ Present the fundamentals of deep learning by implementing an efficient Multilayer Perceptron (MLP) from scratch
 - ▶ Use this MLP to solve a regression problem (**supervised learning**)

About this class

- ▶ Deep learning
 - ▶ Very dynamic field of research
 - ▶ Increased computational power
 - ▶ Increased data availability
 - ▶ Different kinds of neural networks, depending on applications
 - ▶ Convolutional networks for image processing
 - ▶ Recursive neural networks, large language models for natural language processing
 - ▶ Generative adversarial network for physical simulations and image processing
 - ▶ But not many impressive results obtained by using a neural network as a black box, with no knowledge on how they work or on the data
 - ▶ Other ML techniques can outperform neural networks, depending on the problem at hand
- ▶ Goal of this class
 - ▶ Present the fundamentals of deep learning by implementing an efficient Multilayer Perceptron (MLP) from scratch
 - ▶ Use this MLP to solve a regression problem (**supervised learning**)
 - ▶ **Nb:** you may feel lost at the beginning of the hands-on labs but everything will fall into place with time

Prediction Accuracy vs Explainability



"Bornstein, Aaron M. "Is Artificial Intelligence Permanently Inscrutable?". Nautilus , no. (2016)."

Main outline

- ▶ Origins of neural networks and deep learning
- ▶ Main properties of neural networks
- ▶ Learning techniques
 - ▶ Basis (gradient descent, overfitting, ...)
 - ▶ Improvements (gradient optimization, regularization, ...)
- ▶ Application to a pricing problem
 - ▶ Question: can a neural network be used as a proxy for a costly function?

Timetable

1. Introduction, multilayer perceptrons, forward propagation
 - ▶ Hands-on: forward propagation
2. Gradient descents, backpropagation
 - ▶ Hands-on: backpropagation for stochastic gradient descent
3. Mini-batch gradient descent, gradient descent optimization
 - ▶ Hands-on: mini-batch gradient descent
4. Hands-on: mini-batch, Momentum
5. Regularization
 - ▶ Hands-on: L2 regularization
6. Feature preprocessing, batch normalization
 - ▶ Pricing project
7. Energy considerations
 - ▶ Pricing project
8. Pricing project

Some resources

- ▶ Goodfellow, Bengio, Courville. Deep Learning:
<https://www.deeplearningbook.org/>
- ▶ Aggarwal. Neural Networks and Deep Learning
- ▶ Nielsen. Neural Networks and Deep Learning:
<http://neuralnetworksanddeeplearning.com>
- ▶ Lectures by Mallat at Collège de France:
<https://www.college-de-france.fr/site/stephane-mallat/>
- ▶ Lectures by Le Cun at Collège de France:
<https://www.college-de-france.fr/site/yann-lecun/>
- ▶ Stanford School of Engineering channel:
<https://www.youtube.com/user/stanfordeng>

Outline

General presentation

The history of neural networks

The multilayer perceptron

What do we gain with an MLP?

Some landmarks

1943 W. McCulloch, W. Pitts: *A Logical Calculus of the Ideas Immanent in Nervous Activity*

- ▶ Formal model of a neuron
- ▶ Weighted inputs, binary outputs

Some landmarks

1943 W. McCulloch, W. Pitts: *A Logical Calculus of the Ideas Immanent in Nervous Activity*

- ▶ Formal model of a neuron
- ▶ Weighted inputs, binary outputs

1948 N. Wiener: *Cybernetics: Or Control and Communication in the Animal and the Machine*

- ▶ “Self-regulating mechanisms”
- ▶ Notion of feedback

Some landmarks

1943 W. McCulloch, W. Pitts: *A Logical Calculus of the Ideas Immanent in Nervous Activity*

- ▶ Formal model of a neuron
- ▶ Weighted inputs, binary outputs

1948 N. Wiener: *Cybernetics: Or Control and Communication in the Animal and the Machine*

- ▶ “Self-regulating mechanisms”
- ▶ Notion of feedback

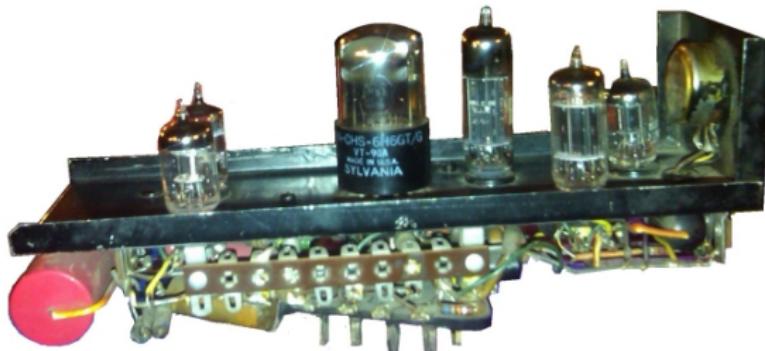
1950 A. Turing: The Turing Test

- ▶ Proposal: build a program simulating a child's mind, then educate it

Some landmarks

- 1943 W. McCulloch, W. Pitts: *A Logical Calculus of the Ideas Immanent in Nervous Activity*
- ▶ Formal model of a neuron
 - ▶ Weighted inputs, binary outputs
- 1948 N. Wiener: *Cybernetics: Or Control and Communication in the Animal and the Machine*
- ▶ “Self-regulating mechanisms”
 - ▶ Notion of feedback
- 1950 A. Turing: The Turing Test
- ▶ Proposal: build a program simulating a child's mind, then educate it
- 1951 M. Minsky and D. Edmonds: Stochastic Neural Analog Reinforcement Calculator (SNARC)
- ▶ First neural network machine

SNARC



- ▶ Image : <https://cyberneticzoo.com/mazesolvers/1951-maze-solver-minsky-edmonds-american/>
- ▶ Learning through trial and error (kind of reinforcement learning)
 - ▶ Feedback provided by the operator
- ▶ 40 neurons
- ▶ Simulate the behavior of a rat in a maze searching for food

A classification problem

1958. F. Rosenblatt: *The Perceptron: A Probabilistic Model for Information Storage and Organization in the Brain*

A classification problem

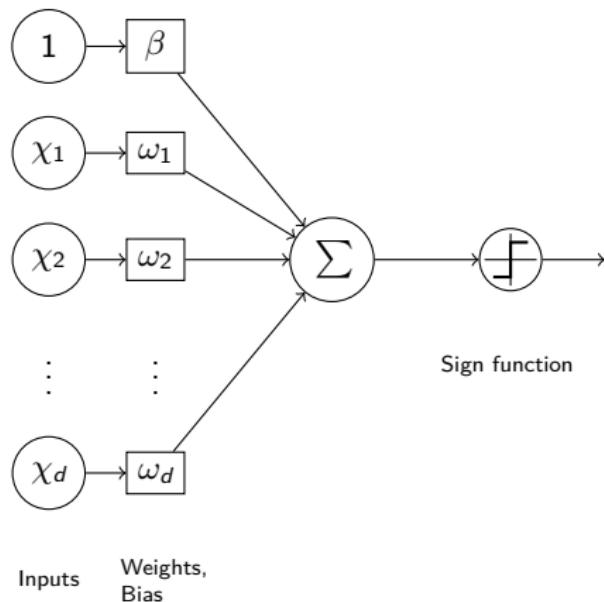
1958. F. Rosenblatt: *The Perceptron: A Probabilistic Model for Information Storage and Organization in the Brain*

- ▶ Input:
 $\{(\chi_{(i)}, \rho_{(i)}) \mid i = 1, \dots, N\}$,
where $\chi_{(i)} \in \mathbb{R}^d$ and
 $\rho_{(i)} \in \{+1, -1\}$
- ▶ Goal: construct a classifier that correctly labels the training set

A classification problem

1958. F. Rosenblatt: *The Perceptron: A Probabilistic Model for Information Storage and Organization in the Brain*

- ▶ Input:
 $\{(\chi(i), \rho(i)) \mid i = 1, \dots, N\}$,
 where $\chi(i) \in \mathbb{R}^d$ and
 $\rho(i) \in \{+1, -1\}$
- ▶ Goal: construct a classifier that correctly labels the training set



A classification problem

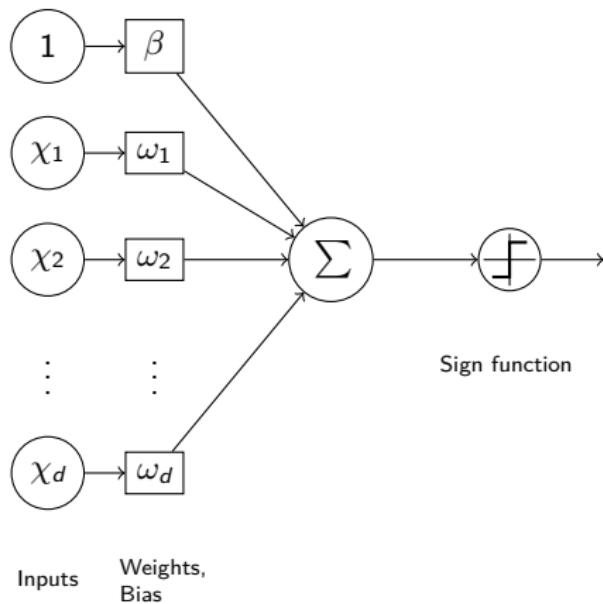
1958. F. Rosenblatt: *The Perceptron: A Probabilistic Model for Information Storage and Organization in the Brain*

► Input:

$$\{(\chi(i), \rho(i)) \mid i = 1, \dots, N\},$$

where $\chi(i) \in \mathbb{R}^d$ and
 $\rho(i) \in \{+1, -1\}$

- Goal: construct a classifier that correctly labels the training set
- Perceptron: outputs
 $\text{sign}(\omega^T \chi + \beta)$



A classification problem

1958. F. Rosenblatt: *The Perceptron: A Probabilistic Model for Information Storage and Organization in the Brain*

► Input:

$$\{(\chi_{(i)}, \rho_{(i)}) \mid i = 1, \dots, N\},$$

where $\chi_{(i)} \in \mathbb{R}^d$ and
 $\rho_{(i)} \in \{+1, -1\}$

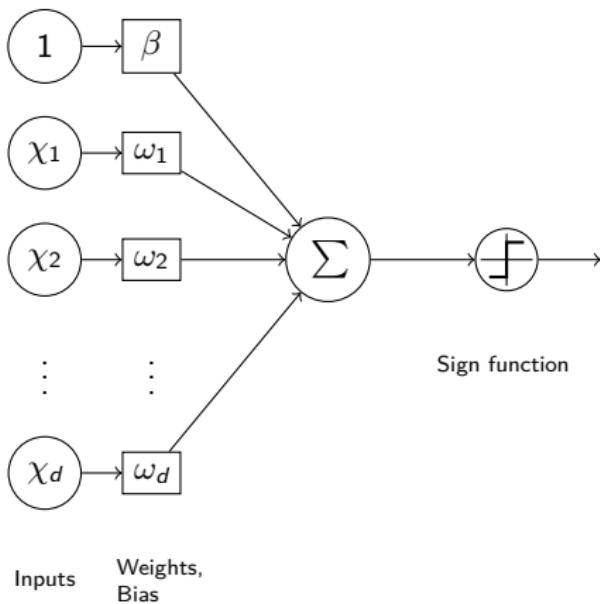
► Goal: construct a classifier that correctly labels the training set

► Perceptron: outputs
 $\text{sign}(\omega^T \chi + \beta)$

► Problem: find ω, β such that

$$\forall i = 1, \dots, N,$$

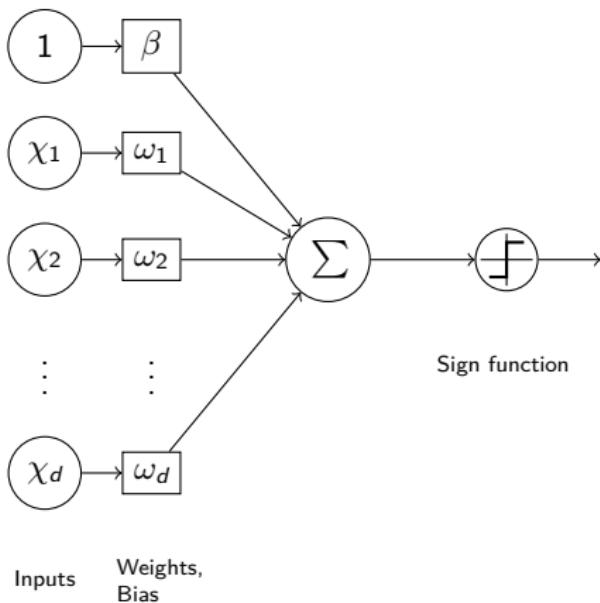
$$\text{sign}(\omega^T \chi_{(i)} + \beta) = \rho_{(i)}$$



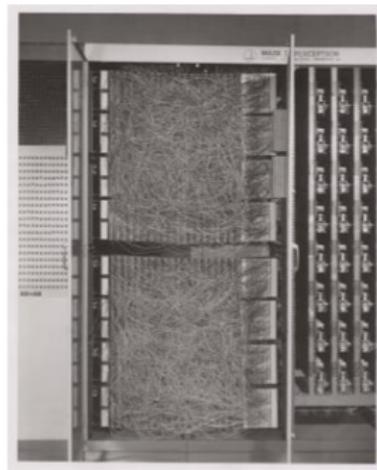
A classification problem

1958. F. Rosenblatt: *The Perceptron: A Probabilistic Model for Information Storage and Organization in the Brain*

- ▶ Input:
 $\{(\chi(i), \rho(i)) \mid i = 1, \dots, N\}$,
 where $\chi(i) \in \mathbb{R}^d$ and
 $\rho(i) \in \{+1, -1\}$
- ▶ Goal: construct a classifier that correctly labels the training set
- ▶ Perceptron: outputs
 $\text{sign}(\omega^T \chi + \beta)$
- ▶ Problem: find ω, β such that
 $\forall i = 1, \dots, N$,
 $\text{sign}(\omega^T \chi(i) + \beta) = \rho(i)$
- ▶ Reformulation: find ω, β such that
 $\forall i = 1, \dots, N$,
 $\rho(i) \cdot (\omega^T \chi(i) + \beta) > 0$



Mark1 - Perceptron



By Source (WP:NFCC#4), Fair use,
<https://en.wikipedia.org/w/index.php?curid=47541432>

Parameter update algorithm

We define:

- ▶ $\theta \stackrel{\text{def}}{=} (\beta, \omega_1, \dots, \omega_d)^T$ and $\bar{\chi} \stackrel{\text{def}}{=} (1, \chi_1, \dots, \chi_d)^T$

Goal: find θ such that $\forall i = 1, \dots, N, \rho_{(i)} \cdot (\theta^T \bar{\chi}_{(i)}) > 0$

Parameter update algorithm

We define:

- ▶ $\theta \stackrel{\text{def}}{=} (\beta, \omega_1, \dots, \omega_d)^T$ and $\bar{\chi} \stackrel{\text{def}}{=} (1, \chi_1, \dots, \chi_d)^T$

Goal: find θ such that $\forall i = 1, \dots, N, \rho_{(i)} \cdot (\theta^T \bar{\chi}_{(i)}) > 0$

Intuition:

- ▶ If for some j we have $\theta^T \bar{\chi}_{(j)} < 0$ and $\rho_{(j)} = 1$, then update θ so that the value of $\theta^T \bar{\chi}_{(j)}$ is increased
- ▶ If for some j we have $\theta^T \bar{\chi}_{(j)} > 0$ and $\rho_{(j)} = -1$, then update θ so that the value of $\theta^T \bar{\chi}_{(j)}$ is decreased

Parameter update algorithm

We define:

- ▶ $\theta \stackrel{\text{def}}{=} (\beta, \omega_1, \dots, \omega_d)^T$ and $\bar{\chi} \stackrel{\text{def}}{=} (1, \chi_1, \dots, \chi_d)^T$

Goal: find θ such that $\forall i = 1, \dots, N, \rho_{(i)} \cdot (\theta^T \bar{\chi}_{(i)}) > 0$

Intuition:

- ▶ If for some j we have $\theta^T \bar{\chi}_{(j)} < 0$ and $\rho_{(j)} = 1$, then update θ so that the value of $\theta^T \bar{\chi}_{(j)}$ is increased
- ▶ If for some j we have $\theta^T \bar{\chi}_{(j)} > 0$ and $\rho_{(j)} = -1$, then update θ so that the value of $\theta^T \bar{\chi}_{(j)}$ is decreased

Input: $\{(\chi_{(i)}, \rho_{(i)}) \mid i = 1, \dots, N\}$

- 1 $\theta \leftarrow (0, \dots, 0)^T;$
- 2 **while** $\exists j \in \{1, \dots, N\}$ such that $\rho_{(j)} \cdot \text{sign}(\theta^T \bar{\chi}_{(j)}) < 0$ **do**
- 3 | $\theta \leftarrow \theta + \rho_{(j)} \bar{\chi}_{(j)}$
- 4 **end**
- 5 **return** θ

Convergence

If there exists a separating hyperplane for the samples, then the algorithm terminates and the resulting perceptron correctly classifies all samples

Convergence

If there exists a separating hyperplane for the samples, then the algorithm terminates and the resulting perceptron correctly classifies all samples

- ▶ Separating hyperplane: $\exists \overline{\omega^*}$ such that for all $i = 1, \dots, N$,
 $\rho(i) \cdot (\overline{\omega^*}^T \chi(i)) > 0$
- ▶ Note that $\overline{\omega^*} \neq \mathbf{0}$

Convergence

If there exists a separating hyperplane for the samples, then the algorithm terminates and the resulting perceptron correctly classifies all samples

- ▶ Separating hyperplane: $\exists \overline{\omega^*}$ such that for all $i = 1, \dots, N$,
 $\rho(i) \cdot (\overline{\omega^*}^T \chi(i)) > 0$
 - ▶ Note that $\overline{\omega^*} \neq \mathbf{0}$
- ▶ Proof principle: determine an upper-bound on the number k of iterations in the **while** loop

Convergence

If there exists a separating hyperplane for the samples, then the algorithm terminates and the resulting perceptron correctly classifies all samples

- ▶ Separating hyperplane: $\exists \overline{\omega^*}$ such that for all $i = 1, \dots, N$,
 $\rho_{(i)} \cdot (\overline{\omega^*}^T \overline{\chi_{(i)}}) > 0$
 - ▶ Note that $\overline{\omega^*} \neq \mathbf{0}$
- ▶ Proof principle: determine an upper-bound on the number k of iterations in the **while** loop
- ▶ We let:
 - ▶ $\delta \stackrel{\text{def}}{=} \min_i [\rho_{(i)} \cdot (\overline{\omega^*}^T \overline{\chi_{(i)}})]$; note that $\delta > 0$

Convergence

If there exists a separating hyperplane for the samples, then the algorithm terminates and the resulting perceptron correctly classifies all samples

- ▶ Separating hyperplane: $\exists \overline{\omega^*}$ such that for all $i = 1, \dots, N$,
 $\rho_{(i)} \cdot (\overline{\omega^*}^T \overline{\chi_{(i)}}) > 0$
 - ▶ Note that $\overline{\omega^*} \neq \mathbf{0}$
- ▶ Proof principle: determine an upper-bound on the number k of iterations in the **while** loop
- ▶ We let:
 - ▶ $\delta \stackrel{\text{def}}{=} \min_i [\rho_{(i)} \cdot (\overline{\omega^*}^T \overline{\chi_{(i)}})]$; note that $\delta > 0$
 - ▶ $R \stackrel{\text{def}}{=} \max_i \|\overline{\chi_{(i)}}\|$; note that $R > 0$

Convergence

If there exists a separating hyperplane for the samples, then the algorithm terminates and the resulting perceptron correctly classifies all samples

- ▶ Separating hyperplane: $\exists \overline{\omega^*}$ such that for all $i = 1, \dots, N$,
 $\rho_{(i)} \cdot (\overline{\omega^*}^T \overline{\chi_{(i)}}) > 0$
 - ▶ Note that $\overline{\omega^*} \neq \mathbf{0}$
- ▶ Proof principle: determine an upper-bound on the number k of iterations in the **while** loop
- ▶ We let:
 - ▶ $\delta \stackrel{\text{def}}{=} \min_i [\rho_{(i)} \cdot (\overline{\omega^*}^T \overline{\chi_{(i)}})]$; note that $\delta > 0$
 - ▶ $R \stackrel{\text{def}}{=} \max_i \|\overline{\chi_{(i)}}\|$; note that $R > 0$
- ▶ We have (See, e.g., <http://www.cs.columbia.edu/~mcollins/courses/6998-2012/notes/perc.converge.pdf>):

$$k \leq \frac{R^2 \|\overline{\omega^*}\|^2}{\delta^2}$$

What if the samples admit no separating hyperplane?

Minsky & Papert (1969): *Perceptrons: an introduction to computational geometry*

What if the samples admit no separating hyperplane?

Minsky & Papert (1969): *Perceptrons: an introduction to computational geometry*

- ▶ The XOR function does not admit a separating hyperplane

What if the samples admit no separating hyperplane?

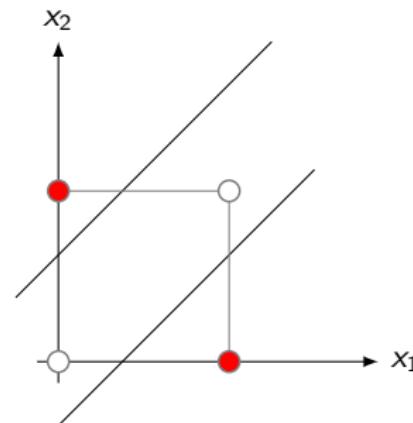
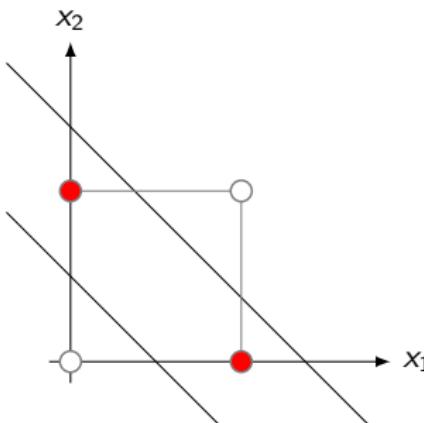
Minsky & Papert (1969): *Perceptrons: an introduction to computational geometry*

- ▶ The XOR function does not admit a separating hyperplane
- ▶ The perceptron is **incapable** of learning this function

What if the samples admit no separating hyperplane?

Minsky & Papert (1969): *Perceptrons: an introduction to computational geometry*

- ▶ The XOR function does not admit a separating hyperplane
- ▶ The perceptron is **incapable** of learning this function



Summary on the perceptron

Summary on the perceptron

- ▶ Simple learning scheme

Summary on the perceptron

- ▶ Simple learning scheme
- ▶ Convergence result when there exists a separating hyperplane
 - ▶ It is possible to evaluate the convergence speed

Summary on the perceptron

- ▶ Simple learning scheme
- ▶ Convergence result when there exists a separating hyperplane
 - ▶ It is possible to evaluate the convergence speed
- ▶ Minsky & Papert: there are simple functions that cannot be learned
 - ▶ This negative result slowed research on neural networks for a long time

Summary on the perceptron

- ▶ Simple learning scheme
- ▶ Convergence result when there exists a separating hyperplane
 - ▶ It is possible to evaluate the convergence speed
- ▶ Minsky & Papert: there are simple functions that cannot be learned
 - ▶ This negative result slowed research on neural networks for a long time
- ▶ Question: what happens if we generalize the perceptron by
 - ▶ Interconnecting several neurons together
 - ▶ Allowing activation functions other than the sign function?

Outline

General presentation

The history of neural networks

The multilayer perceptron

What do we gain with an MLP?

Definitions

A multilayer perceptron (MLP, or feedforward neural network) is a set of interconnected neurons that forms a Directed Acyclic Graph (DAG)

Definitions

A multilayer perceptron (MLP, or feedforward neural network) is a set of interconnected neurons that forms a Directed Acyclic Graph (DAG)

- ▶ The connections between neurons are **weighted**

Definitions

A multilayer perceptron (MLP, or feedforward neural network) is a set of interconnected neurons that forms a Directed Acyclic Graph (DAG)

- ▶ The connections between neurons are **weighted**
- ▶ Neurons with no incoming connection from another neuron are **input neurons**. The set of input neurons forms the **input layer**

Definitions

A multilayer perceptron (MLP, or feedforward neural network) is a set of interconnected neurons that forms a Directed Acyclic Graph (DAG)

- ▶ The connections between neurons are **weighted**
- ▶ Neurons with no incoming connection from another neuron are **input neurons**. The set of input neurons forms the **input layer**
- ▶ Neurons with no outgoing connection to another neuron are **output neurons**. The set of output neurons forms the **output layer**

Definitions

A multilayer perceptron (MLP, or feedforward neural network) is a set of interconnected neurons that forms a Directed Acyclic Graph (DAG)

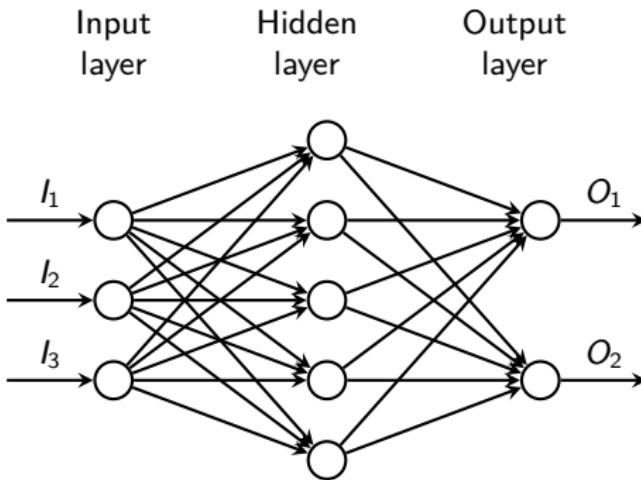
- ▶ The connections between neurons are **weighted**
- ▶ Neurons with no incoming connection from another neuron are **input neurons**. The set of input neurons forms the **input layer**
- ▶ Neurons with no outgoing connection to another neuron are **output neurons**. The set of output neurons forms the **output layer**
- ▶ All other neurons are **hidden neurons**. Hidden neurons are organized in **hidden layers**

Definitions

A multilayer perceptron (MLP, or feedforward neural network) is a set of interconnected neurons that forms a Directed Acyclic Graph (DAG)

- ▶ The connections between neurons are **weighted**
- ▶ Neurons with no incoming connection from another neuron are **input neurons**. The set of input neurons forms the **input layer**
- ▶ Neurons with no outgoing connection to another neuron are **output neurons**. The set of output neurons forms the **output layer**
- ▶ All other neurons are **hidden neurons**. Hidden neurons are organized in **hidden layers**
- ▶ Computations are performed in the hidden and output layers.

Fully connected multilayer perceptron (3 layers)



A fully connected MLP is an MLP such that all non-input neurons are connected to all the neurons of the previous layer

In what follows, all the MLPs we consider will be fully connected

Notations and assumptions

Given a fully connected MLP with L layers:

Notations and assumptions

Given a fully connected MLP with L layers:

- ▶ The neurons of layer i are denoted by $\nu_1^i, \dots, \nu_{n_i}^i$ (layer 0 denotes the input layer)

Notations and assumptions

Given a fully connected MLP with L layers:

- ▶ The neurons of layer i are denoted by $\nu_1^i, \dots, \nu_{n_i}^i$ (layer 0 denotes the input layer)
- ▶ The weight connecting ν_j^{i-1} to ν_k^i is denoted by $\omega_{j,k}^i$

Notations and assumptions

Given a fully connected MLP with L layers:

- ▶ The neurons of layer i are denoted by $\nu_1^i, \dots, \nu_{n_i}^i$ (layer 0 denotes the input layer)
- ▶ The weight connecting ν_j^{i-1} to ν_k^i is denoted by $\omega_{j,k}^i$
- ▶ The bias of ν_k^i is denoted by β_k^i

Notations and assumptions

Given a fully connected MLP with L layers:

- ▶ The neurons of layer i are denoted by $\nu_1^i, \dots, \nu_{n_i}^i$ (layer 0 denotes the input layer)
- ▶ The weight connecting ν_j^{i-1} to ν_k^i is denoted by $\omega_{j,k}^i$
- ▶ The bias of ν_k^i is denoted by β_k^i
- ▶ The sample transmitted to the input layer is denoted by $\alpha^0 = (\alpha_1^0, \dots, \alpha_{n_0}^0)^T$

Notations and assumptions

Given a fully connected MLP with L layers:

- ▶ The neurons of layer i are denoted by $\nu_1^i, \dots, \nu_{n_i}^i$ (layer 0 denotes the input layer)
- ▶ The weight connecting ν_j^{i-1} to ν_k^i is denoted by $\omega_{j,k}^i$
- ▶ The bias of ν_k^i is denoted by β_k^i
- ▶ The sample transmitted to the input layer is denoted by $\alpha^0 = (\alpha_1^0, \dots, \alpha_{n_0}^0)^T$
- ▶ The **net input** of ν_k^i is $\zeta_k^i \stackrel{\text{def}}{=} \left(\sum_{j=1}^{n_{i-1}} \omega_{j,k}^i \alpha_j^{i-1} \right) + \beta_k^i$

Notations and assumptions

Given a fully connected MLP with L layers:

- ▶ The neurons of layer i are denoted by $\nu_1^i, \dots, \nu_{n_i}^i$ (layer 0 denotes the input layer)
- ▶ The weight connecting ν_j^{i-1} to ν_k^i is denoted by $\omega_{j,k}^i$
- ▶ The bias of ν_k^i is denoted by β_k^i
- ▶ The sample transmitted to the input layer is denoted by $\alpha^0 = (\alpha_1^0, \dots, \alpha_{n_0}^0)^T$
- ▶ The **net input** of ν_k^i is $\zeta_k^i \stackrel{\text{def}}{=} \left(\sum_{j=1}^{n_{i-1}} \omega_{j,k}^i \alpha_j^{i-1} \right) + \beta_k^i$
- ▶ The **activation** of ν_k^i is denoted by α_k^i

Notations and assumptions

Given a fully connected MLP with L layers:

- ▶ The neurons of layer i are denoted by $\nu_1^i, \dots, \nu_{n_i}^i$ (layer 0 denotes the input layer)
- ▶ The weight connecting ν_j^{i-1} to ν_k^i is denoted by $\omega_{j,k}^i$
- ▶ The bias of ν_k^i is denoted by β_k^i
- ▶ The sample transmitted to the input layer is denoted by $\alpha^0 = (\alpha_1^0, \dots, \alpha_{n_0}^0)^T$
- ▶ The **net input** of ν_k^i is $\zeta_k^i \stackrel{\text{def}}{=} \left(\sum_{j=1}^{n_{i-1}} \omega_{j,k}^i \alpha_j^{i-1} \right) + \beta_k^i$
- ▶ The **activation** of ν_k^i is denoted by α_k^i
- ▶ The activation function of ν_k^i is Φ , so that $\alpha_k^i = \Phi(\zeta_k^i)$
 - ▶ Unless stated otherwise (e.g., **Dropout**), we assume it is the same for all neurons

Notations and assumptions

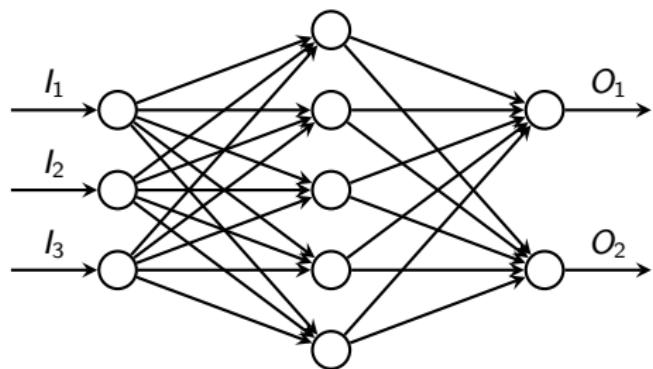
Given a fully connected MLP with L layers:

- ▶ The neurons of layer i are denoted by $\nu_1^i, \dots, \nu_{n_i}^i$ (layer 0 denotes the input layer)
- ▶ The weight connecting ν_j^{i-1} to ν_k^i is denoted by $\omega_{j,k}^i$
- ▶ The bias of ν_k^i is denoted by β_k^i
- ▶ The sample transmitted to the input layer is denoted by $\alpha^0 = (\alpha_1^0, \dots, \alpha_{n_0}^0)^T$
- ▶ The **net input** of ν_k^i is $\zeta_k^i \stackrel{\text{def}}{=} \left(\sum_{j=1}^{n_{i-1}} \omega_{j,k}^i \alpha_j^{i-1} \right) + \beta_k^i$
- ▶ The **activation** of ν_k^i is denoted by α_k^i
- ▶ The activation function of ν_k^i is Φ , so that $\alpha_k^i = \Phi(\zeta_k^i)$
 - ▶ Unless stated otherwise (e.g., **Dropout**), we assume it is the same for all neurons

Assumption

In what follows, unless specified otherwise, we will assume that the MLPs we consider have a single output node

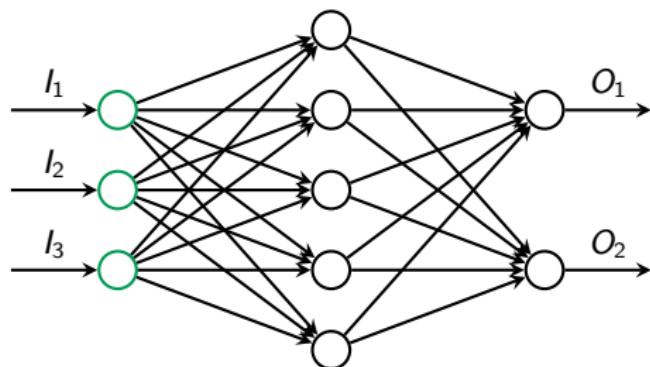
Forward propagation: an illustration



Forward propagation: an illustration

Neuron ν_1^1 :

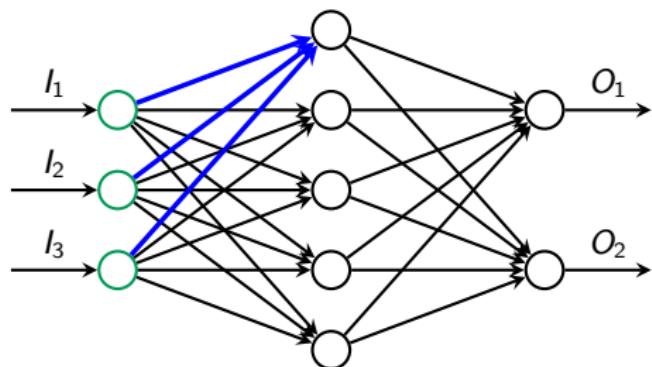
- ▶ Input: a_1^0, a_2^0, a_3^0



Forward propagation: an illustration

Neuron ν_1^1 :

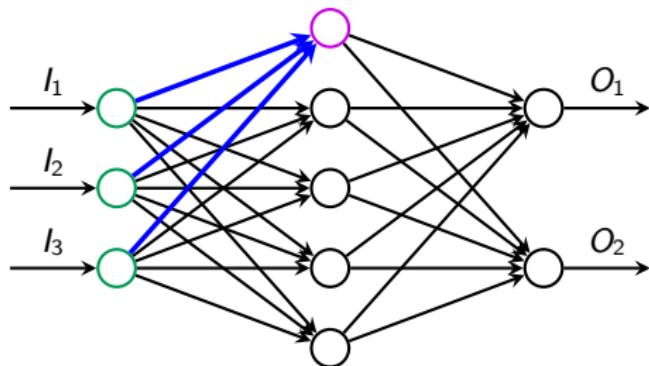
- ▶ Input: a_1^0, a_2^0, a_3^0
- ▶ Weights: $\omega_{1,1}^1, \omega_{2,1}^1, \omega_{3,1}^1$
- ▶ Bias: β_1^1



Forward propagation: an illustration

Neuron ν_1^1 :

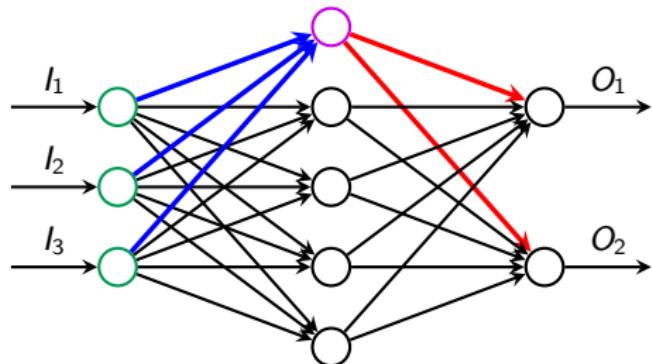
- ▶ Input: $\alpha_1^0, \alpha_2^0, \alpha_3^0$
- ▶ Weights: $\omega_{1,1}^1, \omega_{2,1}^1, \omega_{3,1}^1$
- ▶ Bias: β_1^1
- ▶ Net input: $\zeta_1^1 = \omega_{1,1}^1 \cdot \alpha_1^0 + \omega_{2,1}^1 \cdot \alpha_2^0 + \omega_{3,1}^1 \cdot \alpha_3^0 + \beta_1^1$



Forward propagation: an illustration

Neuron ν_1^1 :

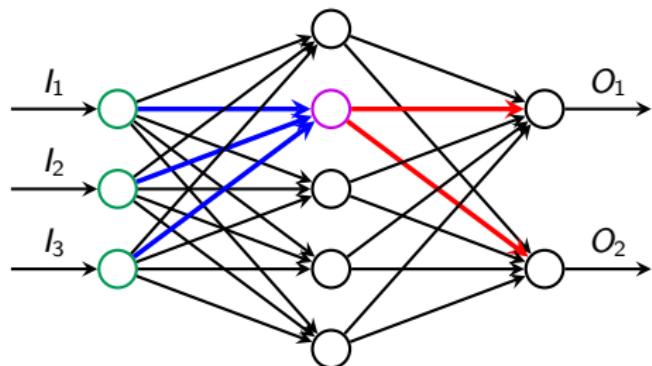
- ▶ Input: $\alpha_1^0, \alpha_2^0, \alpha_3^0$
- ▶ Weights: $\omega_{1,1}^1, \omega_{2,1}^1, \omega_{3,1}^1$
- ▶ Bias: β_1^1
- ▶ Net input: $\zeta_1^1 = \omega_{1,1}^1 \cdot \alpha_1^0 + \omega_{2,1}^1 \cdot \alpha_2^0 + \omega_{3,1}^1 \cdot \alpha_3^0 + \beta_1^1$
- ▶ Activation: $\alpha_1^1 = \Phi(\zeta_1^1)$



Forward propagation: an illustration

Neuron ν_1^1 :

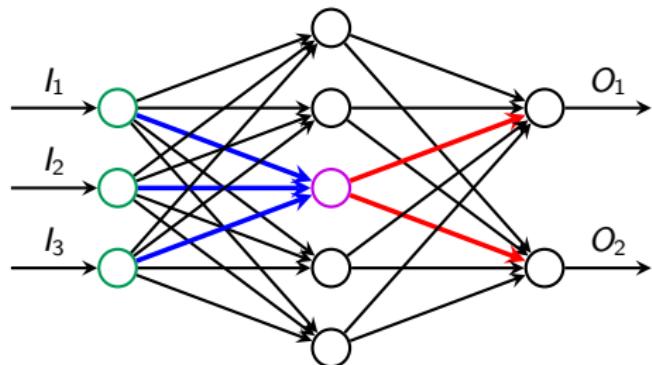
- ▶ Input: $\alpha_1^0, \alpha_2^0, \alpha_3^0$
- ▶ Weights: $\omega_{1,2}^1, \omega_{2,2}^1, \omega_{3,2}^1$
- ▶ Bias: β_2^1
- ▶ Net input: $\zeta_1^1 = \omega_{1,2}^1 \cdot \alpha_1^0 + \omega_{2,2}^1 \cdot \alpha_2^0 + \omega_{3,2}^1 \cdot \alpha_3^0 + \beta_2^1$
- ▶ Activation: $\alpha_2^1 = \Phi(\zeta_2^1)$



Forward propagation: an illustration

Neuron ν_1^1 :

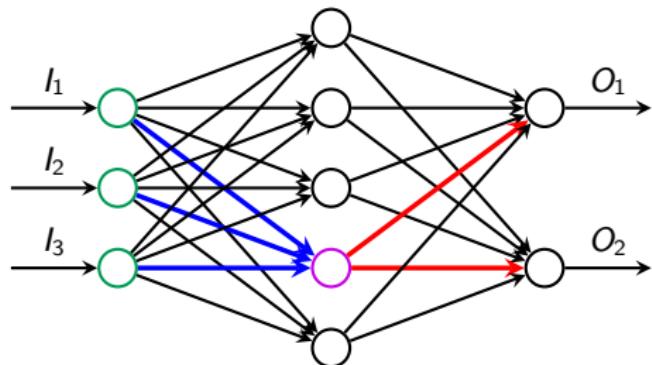
- ▶ Input: $\alpha_1^0, \alpha_2^0, \alpha_3^0$
- ▶ Weights: $\omega_{1,3}^1, \omega_{2,3}^1, \omega_{3,3}^1$
- ▶ Bias: β_3^1
- ▶ Net input: $\zeta_1^1 = \omega_{1,2}^1 \cdot \alpha_1^0 + \omega_{2,3}^1 \cdot \alpha_2^0 + \omega_{3,3}^1 \cdot \alpha_3^0 + \beta_3^1$
- ▶ Activation: $\alpha_3^1 = \Phi(\zeta_3^1)$



Forward propagation: an illustration

Neuron ν_1^1 :

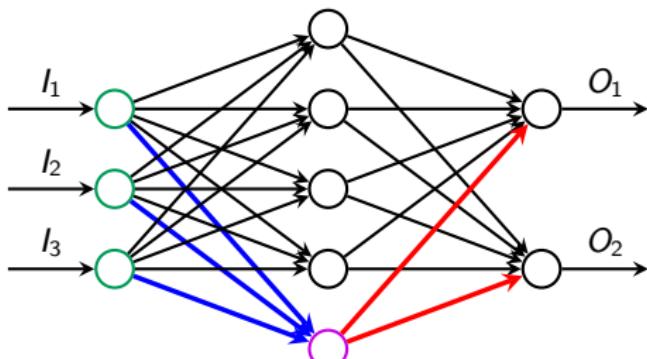
- ▶ Input: $\alpha_1^0, \alpha_2^0, \alpha_3^0$
- ▶ Weights: $\omega_{1,4}^1, \omega_{2,4}^1, \omega_{3,4}^1$
- ▶ Bias: β_4^1
- ▶ Net input: $\zeta_1^1 = \omega_{1,2}^1 \cdot \alpha_1^0 + \omega_{2,4}^1 \cdot \alpha_2^0 + \omega_{3,4}^1 \cdot \alpha_3^0 + \beta_4^1$
- ▶ Activation: $\alpha_4^1 = \Phi(\zeta_4^1)$



Forward propagation: an illustration

Neuron ν_1^1 :

- ▶ Input: $\alpha_1^0, \alpha_2^0, \alpha_3^0$
- ▶ Weights: $\omega_{1,5}^1, \omega_{2,5}^1, \omega_{3,5}^1$
- ▶ Bias: β_5^1
- ▶ Net input: $\zeta_1^1 = \omega_{1,2}^1 \cdot \alpha_1^0 + \omega_{2,5}^1 \cdot \alpha_2^0 + \omega_{3,5}^1 \cdot \alpha_3^0 + \beta_5^1$
- ▶ Activation: $\alpha_5^1 = \Phi(\zeta_5^1)$



A basic forward propagation algorithm

Input: A network with L layers

Input: $\alpha^0 \stackrel{\text{def}}{=} (\alpha_1^0, \dots, \alpha_{n_0}^0)^T$

```
1 for  $i \leftarrow 1$  to  $L$  do
2   for  $k \leftarrow 1$  to  $n_i$  do
3      $\zeta_k^i \leftarrow \left( \sum_{j=1}^{n_{i-1}} \omega_{j,k}^i \alpha_j^{i-1} \right) + \beta_k^i;$ 
4      $\alpha_k^i \leftarrow \Phi(\zeta_k^i);$ 
5   end
6 end
```

Algorithm 1: Forward propagation, basic version

Matrix representation

Matrix representation

- ▶ The weights at layer i are stored in a matrix $\Omega^i \in \mathbb{R}^{n_{i-1} \times n_i}$, where $(\Omega^i)_{j,k} = \omega_{j,k}^i$

Matrix representation

- ▶ The weights at layer i are stored in a matrix $\Omega^i \in \mathbb{R}^{n_{i-1} \times n_i}$, where $(\Omega^i)_{j,k} = \omega_{j,k}^i$
- ▶ The biases at layer i are stored in a vector $\beta^i \stackrel{\text{def}}{=} (\beta_1^i, \dots, \beta_{n_i}^i)^T$

Matrix representation

- ▶ The weights at layer i are stored in a matrix $\Omega^i \in \mathbb{R}^{n_{i-1} \times n_i}$, where $(\Omega^i)_{j,k} = \omega_{j,k}^i$
- ▶ The biases at layer i are stored in a vector $\beta^i \stackrel{\text{def}}{=} (\beta_1^i, \dots, \beta_{n_i}^i)^T$
- ▶ The net inputs at layer i are stored in a vector $\zeta^i \stackrel{\text{def}}{=} (\zeta_1^i, \dots, \zeta_{n_i}^i)^T$

Matrix representation

- ▶ The weights at layer i are stored in a matrix $\Omega^i \in \mathbb{R}^{n_{i-1} \times n_i}$, where $(\Omega^i)_{j,k} = \omega_{j,k}^i$
- ▶ The biases at layer i are stored in a vector $\beta^i \stackrel{\text{def}}{=} (\beta_1^i, \dots, \beta_{n_i}^i)^T$
- ▶ The net inputs at layer i are stored in a vector $\zeta^i \stackrel{\text{def}}{=} (\zeta_1^i, \dots, \zeta_{n_i}^i)^T$
- ▶ The activation function is extended to vectors: if $u = (u_1, \dots, u_n)^T$, then $\Phi(u) = (\Phi(u_1), \dots, \Phi(u_n))^T$

Matrix representation

- ▶ The weights at layer i are stored in a matrix $\Omega^i \in \mathbb{R}^{n_{i-1} \times n_i}$, where $(\Omega^i)_{j,k} = \omega_{j,k}^i$
- ▶ The biases at layer i are stored in a vector $\beta^i \stackrel{\text{def}}{=} (\beta_1^i, \dots, \beta_{n_i}^i)^T$
- ▶ The net inputs at layer i are stored in a vector $\zeta^i \stackrel{\text{def}}{=} (\zeta_1^i, \dots, \zeta_{n_i}^i)^T$
- ▶ The activation function is extended to vectors: if $u = (u_1, \dots, u_n)^T$, then $\Phi(u) = (\Phi(u_1), \dots, \Phi(u_n))^T$
- ▶ The activations of layer i are stored in a vector $\alpha^i \stackrel{\text{def}}{=} (\alpha_1^i, \dots, \alpha_{n_i}^i)^T$

Matrix representation

- ▶ The weights at layer i are stored in a matrix $\Omega^i \in \mathbb{R}^{n_{i-1} \times n_i}$, where $(\Omega^i)_{j,k} = \omega_{j,k}^i$
- ▶ The biases at layer i are stored in a vector $\beta^i \stackrel{\text{def}}{=} (\beta_1^i, \dots, \beta_{n_i}^i)^T$
- ▶ The net inputs at layer i are stored in a vector $\zeta^i \stackrel{\text{def}}{=} (\zeta_1^i, \dots, \zeta_{n_i}^i)^T$
- ▶ The activation function is extended to vectors: if $u = (u_1, \dots, u_n)^T$, then $\Phi(u) = (\Phi(u_1), \dots, \Phi(u_n))^T$
- ▶ The activations of layer i are stored in a vector $\alpha^i \stackrel{\text{def}}{=} (\alpha_1^i, \dots, \alpha_{n_i}^i)^T$

For all $i = 1, \dots, L$, we have $\zeta^i = [\Omega^i]^T \cdot \alpha^{i-1} + \beta^i$

Forward propagation algorithms

Input: A network with L layers

Input: $\alpha^0 \stackrel{\text{def}}{=} (\alpha_1^0, \dots, \alpha_{n_0}^0)^T$

```
1 for  $i \leftarrow 1$  to  $L$  do
2   for  $k \leftarrow 1$  to  $n_i$  do
3      $\zeta_k^i \leftarrow$ 
4      $\left( \sum_{j=1}^{n_{i-1}} \omega_{j,k}^i \alpha_j^{i-1} \right) + \beta_k^i;$ 
5      $\alpha_k^i \leftarrow \Phi(\zeta_k^i);$ 
6   end
7 end
```

Algorithm 2: Forward propagation,
basic version

Forward propagation algorithms

Input: A network with L layers

Input: $\alpha^0 \stackrel{\text{def}}{=} (\alpha_1^0, \dots, \alpha_{n_0}^0)^T$

1 **for** $i \leftarrow 1$ **to** L **do**

2 **for** $k \leftarrow 1$ **to** n_i **do**

3 $\zeta_k^i \leftarrow$

$$\left(\sum_{j=1}^{n_{i-1}} \omega_{j,k}^i \alpha_j^{i-1} \right) + \beta_k^i;$$

4 $\alpha_k^i \leftarrow \Phi(\zeta_k^i);$

5 **end**

6 **end**

Algorithm 4: Forward propagation,
basic version

Input: A network with L layers

Input: α^0

1 **for** $i \leftarrow 1$ **to** L **do**

2 $\zeta^i \leftarrow [\Omega^i]^T \cdot \alpha^{i-1} + \beta^i;$

3 $\alpha^i \leftarrow \Phi(\zeta^i);$

4 **end**

Algorithm 5: Forward propagation,
vectorized version

Outline

General presentation

The history of neural networks

The multilayer perceptron

What do we gain with an MLP?

Identity activation function

What happens if the activation function is the identity?

Identity activation function

What happens if the activation function is the identity?

- ▶ $\alpha^i = \zeta^i = [\Omega^i]^T \alpha^{i-1} + \beta^i$

Identity activation function

What happens if the activation function is the identity?

- ▶ $\alpha^i = \zeta^i = [\Omega^i]^T \alpha^{i-1} + \beta^i$
- ▶ Induction: there exists Γ^i, λ^i such that $\alpha^i = [\Gamma^i]^T \alpha^0 + \lambda^i$:

Identity activation function

What happens if the activation function is the identity?

- ▶ $\alpha^i = \zeta^i = [\Omega^i]^T \alpha^{i-1} + \beta^i$
- ▶ Induction: there exists Γ^i, λ^i such that $\alpha^i = [\Gamma^i]^T \alpha^0 + \lambda^i$:

$$\alpha^i = [\Omega^i]^T \alpha^{i-1} + \beta^i$$

Identity activation function

What happens if the activation function is the identity?

- ▶ $\alpha^i = \zeta^i = [\Omega^i]^T \alpha^{i-1} + \beta^i$
- ▶ Induction: there exists Γ^i, λ^i such that $\alpha^i = [\Gamma^i]^T \alpha^0 + \lambda^i$:

$$\begin{aligned}\alpha^i &= [\Omega^i]^T \alpha^{i-1} + \beta^i \\ &= [\Omega^i]^T ([\Gamma^{i-1}]^T \alpha^0 + \lambda^{i-1}) + \beta^i\end{aligned}$$

Identity activation function

What happens if the activation function is the identity?

- ▶ $\alpha^i = \zeta^i = [\Omega^i]^T \alpha^{i-1} + \beta^i$
- ▶ Induction: there exists Γ^i, λ^i such that $\alpha^i = [\Gamma^i]^T \alpha^0 + \lambda^i$:

$$\begin{aligned}\alpha^i &= [\Omega^i]^T \alpha^{i-1} + \beta^i \\ &= [\Omega^i]^T ([\Gamma^{i-1}]^T \alpha^0 + \lambda^{i-1}) + \beta^i\end{aligned}$$

- ▶ We let $\Gamma^i \stackrel{\text{def}}{=} \Gamma^{i-1} \Omega^i$ and $\lambda^i \stackrel{\text{def}}{=} [\Omega^i]^T \lambda^{i-1} + \beta^i$

Identity activation function

What happens if the activation function is the identity?

- ▶ $\alpha^i = \zeta^i = [\Omega^i]^T \alpha^{i-1} + \beta^i$
- ▶ Induction: there exists Γ^i, λ^i such that $\alpha^i = [\Gamma^i]^T \alpha^0 + \lambda^i$:

$$\begin{aligned}\alpha^i &= [\Omega^i]^T \alpha^{i-1} + \beta^i \\ &= [\Omega^i]^T ([\Gamma^{i-1}]^T \alpha^0 + \lambda^{i-1}) + \beta^i\end{aligned}$$

- ▶ We let $\Gamma^i \stackrel{\text{def}}{=} \Gamma^{i-1} \Omega^i$ and $\lambda^i \stackrel{\text{def}}{=} [\Omega^i]^T \lambda^{i-1} + \beta^i$
- ▶ At the last layer, we have $\alpha^L = [\Gamma^L]^T \alpha^0 + \lambda^L$
 - ▶ This is the perceptron without the step activation function

Linear activation function

What happens if the activation function is linear: $\Phi : z \mapsto \mu z + \lambda$?

Linear activation function

What happens if the activation function is linear: $\Phi : z \mapsto \mu z + \lambda$?

- ▶ Let $\bar{\mu} \stackrel{\text{def}}{=} \mu \cdot \mathbf{Id}$ and $\bar{\lambda} \stackrel{\text{def}}{=} (\lambda, \dots, \lambda)^T$

Linear activation function

What happens if the activation function is linear: $\Phi : z \mapsto \mu z + \lambda$?

- ▶ Let $\bar{\mu} \stackrel{\text{def}}{=} \mu \cdot \mathbf{Id}$ and $\bar{\lambda} \stackrel{\text{def}}{=} (\lambda, \dots, \lambda)^T$
- ▶ At layer i we have $\alpha^i = \bar{\mu} \zeta^i + \bar{\lambda} = \bar{\mu} [\Omega^i]^T \alpha^{i-1} + \bar{\mu} \beta^i + \bar{\lambda}$

Linear activation function

What happens if the activation function is linear: $\Phi : z \mapsto \mu z + \lambda$?

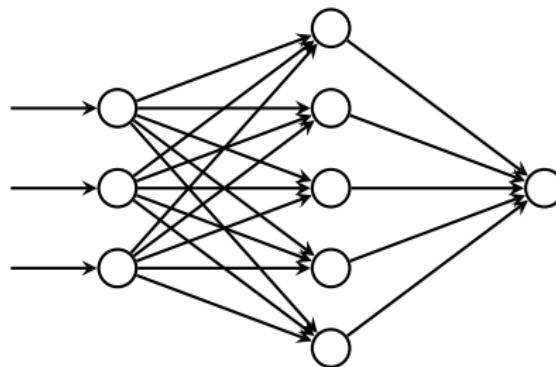
- ▶ Let $\bar{\mu} \stackrel{\text{def}}{=} \mu \cdot \mathbf{Id}$ and $\bar{\lambda} \stackrel{\text{def}}{=} (\lambda, \dots, \lambda)^T$
- ▶ At layer i we have $\alpha^i = \bar{\mu} \zeta^i + \bar{\lambda} = \bar{\mu} [\Omega^i]^T \alpha^{i-1} + \bar{\mu} \beta^i + \bar{\lambda}$
- ▶ By letting $\underline{\Omega}^i \stackrel{\text{def}}{=} \Omega^i [\bar{\mu}]^T$ and $\underline{\beta}^i \stackrel{\text{def}}{=} \bar{\mu} \beta^i + \bar{\lambda}$ we have $\alpha^i = [\underline{\Omega}^i]^T \alpha^{i-1} + \underline{\beta}^i$

Linear activation function

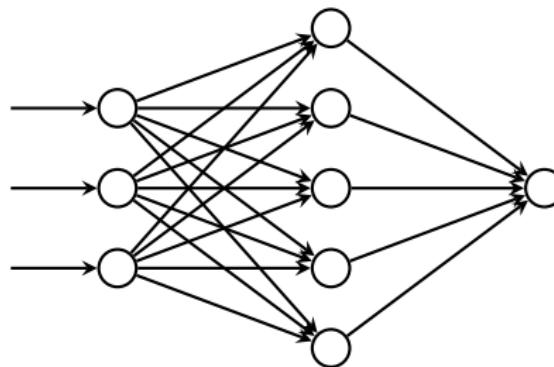
What happens if the activation function is linear: $\Phi : z \mapsto \mu z + \lambda$?

- ▶ Let $\bar{\mu} \stackrel{\text{def}}{=} \mu \cdot \mathbf{Id}$ and $\bar{\lambda} \stackrel{\text{def}}{=} (\lambda, \dots, \lambda)^T$
- ▶ At layer i we have $\alpha^i = \bar{\mu} \zeta^i + \bar{\lambda} = \bar{\mu} [\Omega^i]^T \alpha^{i-1} + \bar{\mu} \beta^i + \bar{\lambda}$
 - ▶ By letting $\underline{\Omega}^i \stackrel{\text{def}}{=} \Omega^i [\bar{\mu}]^T$ and $\underline{\beta}^i \stackrel{\text{def}}{=} \bar{\mu} \beta^i + \bar{\lambda}$ we have $\alpha^i = [\underline{\Omega}^i]^T \alpha^{i-1} + \underline{\beta}^i$
- ▶ This is the same as having an identity activation function

A specific MLP

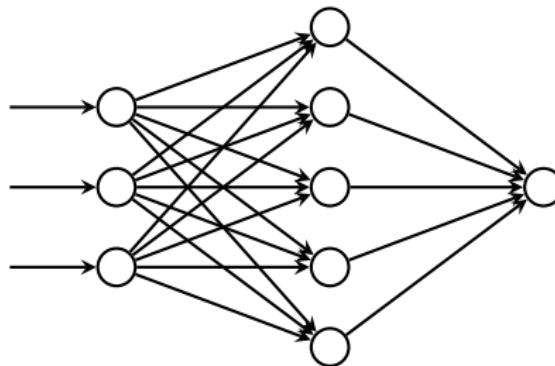


A specific MLP



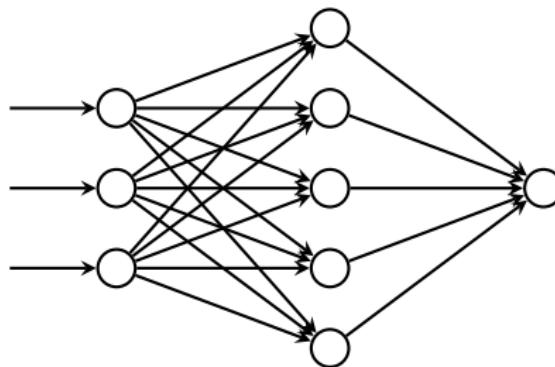
- ▶ Input layer of size d

A specific MLP



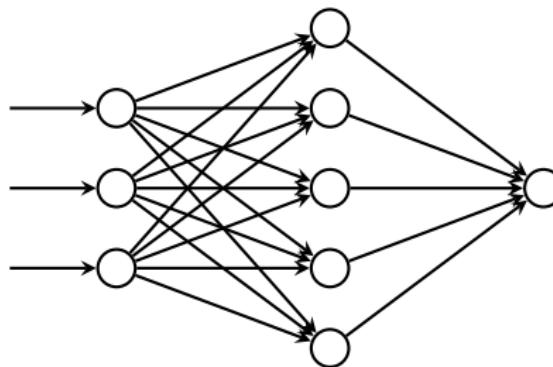
- ▶ Input layer of size d
- ▶ Single hidden layer of size K

A specific MLP



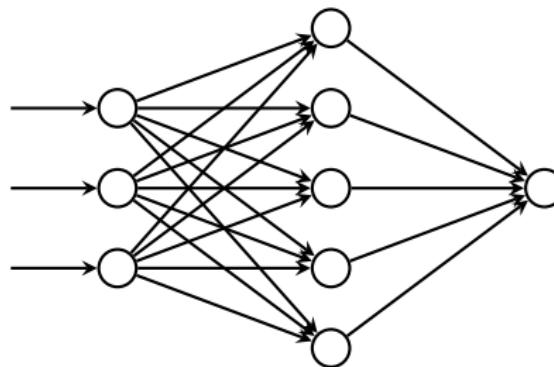
- ▶ Input layer of size d
- ▶ Single hidden layer of size K
- ▶ Single output neuron

A specific MLP



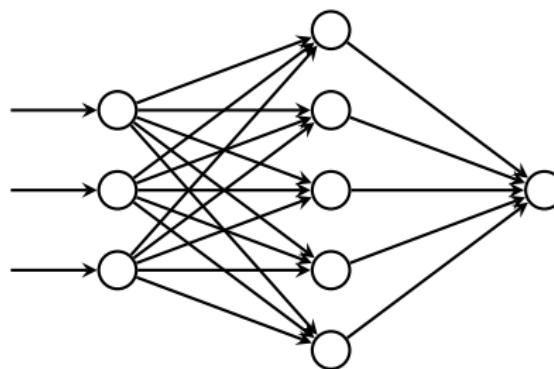
- ▶ Input layer of size d
- ▶ Single hidden layer of size K
- ▶ Single output neuron
- ▶ Activation function Φ only on hidden layer (identity on output)

A specific MLP



- ▶ Input layer of size d
- ▶ Single hidden layer of size K
- ▶ Single output neuron
- ▶ Activation function Φ only on hidden layer (identity on output)
- ▶ No bias on output layer

A specific MLP



- ▶ Input layer of size d
- ▶ Single hidden layer of size K
- ▶ Single output neuron
- ▶ Activation function Φ only on hidden layer (identity on output)
- ▶ No bias on output layer
- ▶ Activation of output neuron is $[\omega']^T \cdot \Phi([\Omega]^T \alpha^0 + \beta)$

A class of functions

- We define

$$\Lambda \stackrel{\text{def}}{=} \left\{ \langle K, \Phi, \Omega, \omega', \beta \rangle \mid K \in \mathbb{N}, \Phi : \mathbb{R} \rightarrow \mathbb{R}, \Omega \in \mathbb{R}^{d \times K}, \beta, \omega' \in \mathbb{R}^K \right\}$$

A class of functions

- ▶ We define

$$\Lambda \stackrel{\text{def}}{=} \left\{ \langle K, \Phi, \Omega, \omega', \beta \rangle \mid K \in \mathbb{N}, \Phi : \mathbb{R} \rightarrow \mathbb{R}, \Omega \in \mathbb{R}^{d \times K}, \beta, \omega' \in \mathbb{R}^K \right\}$$

- ▶ For $\lambda \in \Lambda$, we define

$$\begin{aligned} f_\lambda : \mathbb{R}^d &\rightarrow \mathbb{R} \\ x &\mapsto [\omega']^\top \cdot \Phi(\Omega^\top x + \beta) = \sum_{k=1}^K \omega'_k \Phi([\Omega_k]^\top x + \beta_k) \end{aligned}$$

A class of functions

- ▶ We define

$$\Lambda \stackrel{\text{def}}{=} \left\{ \langle K, \Phi, \Omega, \omega', \beta \rangle \mid K \in \mathbb{N}, \Phi : \mathbb{R} \rightarrow \mathbb{R}, \Omega \in \mathbb{R}^{d \times K}, \beta, \omega' \in \mathbb{R}^K \right\}$$

- ▶ For $\lambda \in \Lambda$, we define

$$\begin{aligned} f_\lambda : \mathbb{R}^d &\rightarrow \mathbb{R} \\ x &\mapsto [\omega']^\top \cdot \Phi(\Omega^\top x + \beta) = \sum_{k=1}^K \omega'_k \Phi([\Omega_k]^\top x + \beta_k) \end{aligned}$$

- ▶ We consider the set of functions $\mathcal{F}(\Lambda) \stackrel{\text{def}}{=} \{f_\lambda \mid \lambda \in \Lambda\}$

A class of functions

- We define

$$\Lambda \stackrel{\text{def}}{=} \left\{ \langle K, \Phi, \Omega, \omega', \beta \rangle \mid K \in \mathbb{N}, \Phi : \mathbb{R} \rightarrow \mathbb{R}, \Omega \in \mathbb{R}^{d \times K}, \beta, \omega' \in \mathbb{R}^K \right\}$$

- For $\lambda \in \Lambda$, we define

$$\begin{aligned} f_\lambda : \mathbb{R}^d &\rightarrow \mathbb{R} \\ x &\mapsto [\omega']^\top \cdot \Phi(\Omega^\top x + \beta) = \sum_{k=1}^K \omega'_k \Phi([\Omega_k]^\top x + \beta_k) \end{aligned}$$

- We consider the set of functions $\mathcal{F}(\Lambda) \stackrel{\text{def}}{=} \{f_\lambda \mid \lambda \in \Lambda\}$

What functions can be approximated by an element from $\mathcal{F}(\Lambda)$?

The Universal Approximation Theorem

Let $\mathcal{C}(\mathbb{R}^d)$ denote the set of continuous functions on \mathbb{R}^d and consider the set $\mathcal{M}(\Phi) \stackrel{\text{def}}{=} \text{span}\left(\left\{x \mapsto \Phi(\omega^T x + \beta) \mid \omega \in \mathbb{R}^d, \beta \in \mathbb{R}\right\}\right)$. This set is dense in $\mathcal{C}(\mathbb{R}^d)$, in the topology of uniform convergence on compact sets if and only if Φ is nonpolynomial

The Universal Approximation Theorem

Let $\mathcal{C}(\mathbb{R}^d)$ denote the set of continuous functions on \mathbb{R}^d and consider the set $\mathcal{M}(\Phi) \stackrel{\text{def}}{=} \text{span}\left(\left\{x \mapsto \Phi(\omega^T x + \beta) \mid \omega \in \mathbb{R}^d, \beta \in \mathbb{R}\right\}\right)$. This set is dense in $\mathcal{C}(\mathbb{R}^d)$, in the topology of uniform convergence on compact sets if and only if Φ is nonpolynomial

In other words: if $f \in \mathcal{C}(\mathbb{R}^d)$ and $\mathcal{D} \subseteq \mathbb{R}^d$ is a compact set, then for all $\varepsilon > 0$, there exists $\tilde{f} \in \mathcal{M}(\Phi)$ such that for all $x \in \mathcal{D}$, $|f(x) - \tilde{f}(x)| \leq \varepsilon$

The Universal Approximation Theorem

Let $\mathcal{C}(\mathbb{R}^d)$ denote the set of continuous functions on \mathbb{R}^d and consider the set $\mathcal{M}(\Phi) \stackrel{\text{def}}{=} \text{span}\left(\left\{x \mapsto \Phi(\omega^T x + \beta) \mid \omega \in \mathbb{R}^d, \beta \in \mathbb{R}\right\}\right)$. This set is dense in $\mathcal{C}(\mathbb{R}^d)$, in the topology of uniform convergence on compact sets if and only if Φ is nonpolynomial

In other words: if $f \in \mathcal{C}(\mathbb{R}^d)$ and $\mathcal{D} \subseteq \mathbb{R}^d$ is a compact set, then for all $\varepsilon > 0$, there exists $\tilde{f} \in \mathcal{M}(\Phi)$ such that for all $x \in \mathcal{D}$, $|f(x) - \tilde{f}(x)| \leq \varepsilon$

- ▶ For a given $\varepsilon > 0$, there exists $K_\varepsilon \in \mathbb{N}$ such that an MLP with a single hidden layer of size K_ε can approximate f with a precision of ε

The Universal Approximation Theorem

Let $\mathcal{C}(\mathbb{R}^d)$ denote the set of continuous functions on \mathbb{R}^d and consider the set $\mathcal{M}(\Phi) \stackrel{\text{def}}{=} \text{span}\left(\left\{x \mapsto \Phi(\omega^T x + \beta) \mid \omega \in \mathbb{R}^d, \beta \in \mathbb{R}\right\}\right)$. This set is dense in $\mathcal{C}(\mathbb{R}^d)$, in the topology of uniform convergence on compact sets if and only if Φ is nonpolynomial

In other words: if $f \in \mathcal{C}(\mathbb{R}^d)$ and $\mathcal{D} \subseteq \mathbb{R}^d$ is a compact set, then for all $\varepsilon > 0$, there exists $\tilde{f} \in \mathcal{M}(\Phi)$ such that for all $x \in \mathcal{D}$, $|f(x) - \tilde{f}(x)| \leq \varepsilon$

- ▶ For a given $\varepsilon > 0$, there exists $K_\varepsilon \in \mathbb{N}$ such that an MLP with a single hidden layer of size K_ε can approximate f with a precision of ε
- ▶ What is wrong with polynomials?
 - ▶ If Φ is a polynomial of degree m , then $\mathcal{M}(\Phi)$ only contains polynomials of degree at most m

The Universal Approximation Theorem

Let $\mathcal{C}(\mathbb{R}^d)$ denote the set of continuous functions on \mathbb{R}^d and consider the set $\mathcal{M}(\Phi) \stackrel{\text{def}}{=} \text{span}\left(\left\{x \mapsto \Phi(\omega^T x + \beta) \mid \omega \in \mathbb{R}^d, \beta \in \mathbb{R}\right\}\right)$. This set is dense in $\mathcal{C}(\mathbb{R}^d)$, in the topology of uniform convergence on compact sets if and only if Φ is nonpolynomial

In other words: if $f \in \mathcal{C}(\mathbb{R}^d)$ and $\mathcal{D} \subseteq \mathbb{R}^d$ is a compact set, then for all $\varepsilon > 0$, there exists $\tilde{f} \in \mathcal{M}(\Phi)$ such that for all $x \in \mathcal{D}$, $|f(x) - \tilde{f}(x)| \leq \varepsilon$

- ▶ For a given $\varepsilon > 0$, there exists $K_\varepsilon \in \mathbb{N}$ such that an MLP with a single hidden layer of size K_ε can approximate f with a precision of ε
- ▶ What is wrong with polynomials?
 - ▶ If Φ is a polynomial of degree m , then $\mathcal{M}(\Phi)$ only contains polynomials of degree at most m

A visual "proof" (or intuition):

<http://neuralnetworksanddeeplearning.com/chap4.html>

Issues with the Universal Approximation Theorem

- ▶ The size of the hidden layer can be extremely large

Issues with the Universal Approximation Theorem

- ▶ The size of the hidden layer can be extremely large
 - ▶ In practice, this result does not seem very useful

Issues with the Universal Approximation Theorem

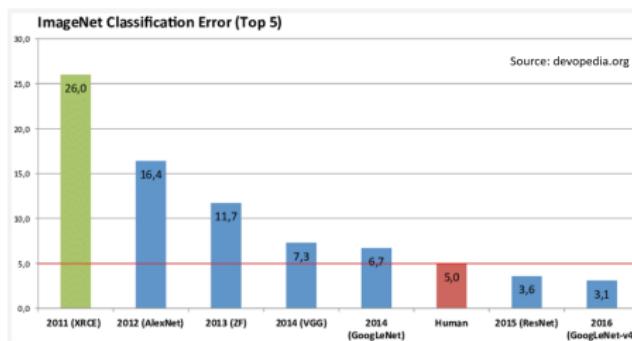
- ▶ The size of the hidden layer can be extremely large
 - ▶ In practice, this result does not seem very useful
- ▶ What happens if we increase the number of hidden layers and bound their sizes?

Issues with the Universal Approximation Theorem

- ▶ The size of the hidden layer can be extremely large
 - ▶ In practice, this result does not seem very useful
- ▶ What happens if we increase the number of hidden layers and bound their sizes?
 - ▶ For a long time, not so clear there was anything to gain

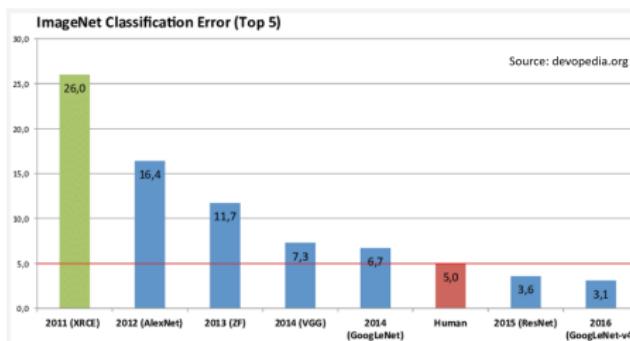
Issues with the Universal Approximation Theorem

- ▶ The size of the hidden layer can be extremely large
 - ▶ In practice, this result does not seem very useful
- ▶ What happens if we increase the number of hidden layers and bound their sizes?
 - ▶ For a long time, not so clear there was anything to gain
 - ▶ Empirically yes: impressive results using deep networks



Issues with the Universal Approximation Theorem

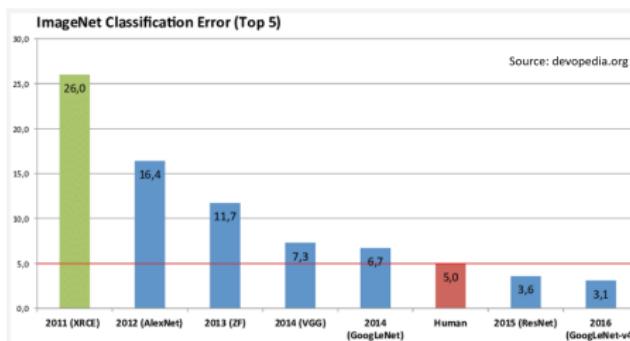
- ▶ The size of the hidden layer can be extremely large
 - ▶ In practice, this result does not seem very useful
- ▶ What happens if we increase the number of hidden layers and bound their sizes?
 - ▶ For a long time, not so clear there was anything to gain
 - ▶ Empirically yes: impressive results using deep networks



- ▶ More recent theoretical results (Eldan, Shamir. The power of Depth for Feedforward Neural Networks, 2016)

Issues with the Universal Approximation Theorem

- ▶ The size of the hidden layer can be extremely large
 - ▶ In practice, this result does not seem very useful
- ▶ What happens if we increase the number of hidden layers and bound their sizes?
 - ▶ For a long time, not so clear there was anything to gain
 - ▶ Empirically yes: impressive results using deep networks



- ▶ More recent theoretical results (Eldan, Shamir. The power of Depth for Feedforward Neural Networks, 2016)
- ▶ For the time being, it is still not possible to formally explain why things work so well