

FACULTY OF ENGINEERING

E016350 - Artificial Intelligence Lecture 6

Machine learning Neural networks Part 2

Aleksandra Pizurica

Ghent University Spring 2025

Feed-forward networks - example



Feed-forward network = a parameterized family of nonlinear functions:

$$a_5 = g(w_{3,5}a_3 + w_{4,5}a_4) = g\Big(w_{3,5}g(w_{1,3}x_1 + w_{2,3}x_2) + w_{4,5}g(w_{1,4}x_1 + w_{2,4}x_2)\Big)$$

Adjusting weights changes the function: do learning this way!

Weight matrix

(For simplicity, we assume no bias)



$$a_{5} = g\left(\begin{bmatrix}w_{3,5} & w_{4,5}\end{bmatrix}\begin{bmatrix}a_{3}\\a_{4}\end{bmatrix}\right) = g\left(\underbrace{\begin{bmatrix}w_{3,5} & w_{4,5}\end{bmatrix}}_{\mathbf{W}^{(2)}}g\left(\underbrace{\begin{bmatrix}w_{1,3} & w_{2,3}\\w_{1,4} & w_{2,4}\end{bmatrix}}_{\mathbf{W}^{(1)}}\begin{bmatrix}x_{1}\\x_{2}\end{bmatrix}\right)\right)$$
$$= g\left(\mathbf{W}^{(2)}g\left(\mathbf{W}^{(1)}\mathbf{x}\right)\right)$$

 $\mathbf{W}^{(l)}$ is the weight matrix in the *l*-th layer. Its rows are the weight vectors. Omitting the layer index: $\mathbf{W}(j,:) = \mathbf{w}_j^{\top} = [w_{0,j} \dots w_{n,j}].$

A. Pizurica, E016350 Artificial Intelligence (UGent) Spring 2025

Matrix notation for multilayer neural networks

For now still assume **no bias** (What should change in the equations otherwise?)



$$a^{(1)} = g(W^{(1)}x)$$

$$a^{(2)} = g(W^{(2)}a^{(1)})$$

$$\vdots$$

$$a^{(L-1)} = g(W^{(L-1)}a^{(L-2)})$$

$$z^{OUT} = W^{(L)}a^{(L-1)}$$

Regression: $h_{\mathbf{w}}(\mathbf{x}) = z^{OUT}$; Classification: $h_{\mathbf{w}}(\mathbf{x}) = Threshold(z^{OUT})$ (or feed the output scores to sigmoid or to softmax for multiclass classification)

Matrix notation for multilayer neural networks



For m neurons in layer i with d inputs:

Matrix notation for multilayer neural networks

For the neural network on the left: ____

$$\begin{pmatrix} \mathbf{w}_{1}^{(1)} \end{pmatrix}_{2} = w_{2,1}^{(1)} \\ \mathbf{x}_{1} \\ \mathbf{x}_{2} \\ \mathbf{x}_{3} \\ \mathbf{w}_{4}^{(1)} \end{pmatrix}_{3} = w_{3,4}^{(1)} \\ \begin{pmatrix} \mathbf{w}_{1}^{(1)} \end{pmatrix}_{3} = w_{3}^{(2)} \\ \begin{pmatrix} \mathbf{w}_{1}^{(1)} \end{pmatrix}_{3} = w_{3,4}^{(1)} \\ \begin{pmatrix} \mathbf{w}_{1}^{(1)} \\ \mathbf{w}_{2}^{(1)} \\ \mathbf{w}_{1}^{(1)} \end{pmatrix}_{3} = w_{3,4}^{(1)} \\ \begin{pmatrix} \mathbf{w}_{1}^{(1)} \\ \mathbf{w}_{2}^{(1)} \\ \mathbf{w}_{2}^{(1)} \\ \mathbf{w}_{3}^{(1)} \\ \mathbf{w}_$$

 $\left(\mathbf{w}_{4}^{(1)}\right)_{3}$

Two-layer regression neural network



Hypothesis class:

$$\mathcal{H} = \{h_{\mathbf{W}^{(1)}, \mathbf{w}^{(2)}} : \mathbf{W}^{(1)} \in \mathbb{R}^{m \times d}, \mathbf{w}^{(2)} \in \mathbb{R}^{m}\}$$

Two-layer classification neural network



Hypothesis class:

$$\mathcal{H} = \{ h_{\mathbf{W}^{(1)}, \mathbf{w}^{(2)}} : \mathbf{W}^{(1)} \in \mathbb{R}^{m \times d}, \mathbf{w}^{(2)} \in \mathbb{R}^{m} \}$$

Compact matrix notation for multilayer neural networks

We can use the same compact notation when bias is not omitted



$$\mathbf{a}^{(1)} = g(\mathbf{W}^{(1)}\mathbf{x})$$
$$\mathbf{a}^{(2)} = g(\mathbf{W}^{(2)}\mathbf{a}^{(1)})$$
$$\vdots$$
$$\mathbf{a}^{(L-1)} = g(\mathbf{W}^{(L-1)}\mathbf{a}^{(L-2)})$$
$$z^{OUT} = \mathbf{W}^{(L)}\mathbf{a}^{(L-1)}$$

• We stipulate: each unit has an extra input from a dummy unit that is fixed to +1 and a weight $w_{0,j}$ for that input

Multilayer neural networks in matrix notation



Slide adapted from: M. Charikar and Koyejo: Artificial Intelligence: Principles and Techniques (Stanford).

Multilayer neural networks in matrix notation



Layers represent multiple levels of abstractions



Optimization in neural networks: A motivating example

Consider regression with a four-layer neural network.

Loss on one example:

 $Loss(\mathbf{x}, y, \mathbf{W}^{(1)}, \mathbf{W}^{(2)}, \mathbf{W}^{(3)}, \mathbf{w}^{(4)}) = (\mathbf{w}^{(4)} \cdot g(\mathbf{W}^{(3)}g(\mathbf{W}^{(2)}g(\mathbf{W}^{(1)}\mathbf{x}))) - y)^2$

(Stochastic) gradient descent:

$$\begin{split} \mathbf{W}^{(1)} &\leftarrow \mathbf{W}^{(1)} - \alpha \nabla_{\mathbf{W}^{(1)}} Loss(\mathbf{x}, y, \mathbf{W}^{(1)}, \mathbf{W}^{(2)}, \mathbf{W}^{(3)}, \mathbf{w}^{(4)}) \\ \mathbf{W}^{(2)} &\leftarrow \mathbf{W}^{(2)} - \alpha \nabla_{\mathbf{W}^{(2)}} Loss(\mathbf{x}, y, \mathbf{W}^{(1)}, \mathbf{W}^{(2)}, \mathbf{W}^{(3)}, \mathbf{w}^{(4)}) \\ \mathbf{W}^{(3)} &\leftarrow \mathbf{W}^{(3)} - \alpha \nabla_{\mathbf{W}^{(3)}} Loss(\mathbf{x}, y, \mathbf{W}^{(1)}, \mathbf{W}^{(2)}, \mathbf{W}^{(3)}, \mathbf{w}^{(4)}) \\ \mathbf{w}^{(4)} &\leftarrow \mathbf{W}^{(3)} - \alpha \nabla_{\mathbf{w}^{(4)}} Loss(\mathbf{x}, y, \mathbf{W}^{(1)}, \mathbf{W}^{(2)}, \mathbf{W}^{(3)}, \mathbf{w}^{(4)}) \end{split}$$

We learned to compute the gradients in the inner layers by backpropagation. Still, this seems very complex! Easy with automatic differentiation tools that use computation graphs.

A. Pizurica, E016350 Artificial Intelligence (UGent) Sprin

Spring 2025

Neural networks 13 / 2

Digression: How do we typically represent deep neural nets



Simpler:



or even simpler:



Sergey Levine: Backpropagation - Designing, Visualizing and Understanding Deep Neural Networks.

A. Pizurica, E016350 Artificial Intelligence (UGent) Spring 2025

Neural networks 14 /

Gradient of the loss function



$$\nabla_{\mathbf{w}} Loss = \left(\frac{\partial Loss}{\partial \mathbf{w}}\right)^{\top}$$
$$\frac{\partial Loss}{\partial \mathbf{w}^{(2)}} \in \mathbb{R}^{1 \times n} ; \quad \frac{\partial Loss}{\partial \mathbf{W}^{(1)}} \in \mathbb{R}^{m \times d}$$

Backpropagation



$$\frac{\partial Loss}{\partial \mathbf{W}^{(1)}} = \frac{\partial \mathbf{z}^{(1)}}{\partial \mathbf{W}^{(1)}} \underbrace{\frac{\partial \mathbf{a}^{(1)}}{\partial \mathbf{z}^{(1)}}}_{\text{compute first: } \delta} \underbrace{\frac{\partial \mathbf{z}^{(2)}}{\partial \mathbf{a}^{(1)}}}_{\text{update: new } \delta}$$

Why this recursion?

C

Make "cheap" computation in each step nitialize: $\frac{\partial Loss}{\partial \sigma^{(2)}} = \delta_{init}$ compute: $\frac{\partial \mathbf{z}^{(2)}}{\partial \mathbf{z}^{(1)}} \delta_{init} = \delta$ $\frac{\partial Loss}{\partial \mathbf{W}^{(1)}} = \frac{\partial \mathbf{z}^{(1)}}{\partial \mathbf{W}^{(1)}} \underbrace{\frac{\partial \mathbf{a}^{(1)}}{\partial \mathbf{z}^{(1)}} \delta}_{\mathbf{z}^{(1)}}$ $O(n^2)$ compute: $\frac{\partial \mathbf{a}^{(1)}}{\partial \mathbf{r}^{(1)}} \delta = \delta_{new}$ $\frac{\partial Loss}{\partial \mathbf{W}^{(1)}} = \frac{\partial \mathbf{z}^{(1)}}{\partial \mathbf{W}^{(1)}} \delta_{new}$ $O(n^2)$

Computation graphs



Definition (Computation graph)

A directed acyclic graph whose root node represents the final mathematical expression and each node represents intermediate subexpressions.

- Automatically compute gradients (how TensorFlow and PyTorch work)
- Gain insight into modular structure of gradient computations

Neural networks 18 / 2

Computation graphs concepts: Functions as boxes



Gradients: how much does c change if a or b changes?

M. Charikar and Koyejo: Artificial Intelligence: Principles and Techniques (Stanford).

Basic building blocks of computation graphs



 σ denotes sigmoid (logistic) function. M. Charikar & Koyejo: Artificial Intelligence: Principles and Techniques (Stanford).

Function composition



$$\frac{\partial c}{\partial a} = \frac{\partial c}{\partial b} \frac{\partial b}{\partial a} = (2b)(2a) = (2a^2)(2a) = 4a^3$$

M. Charikar and Koyejo: Artificial Intelligence: Principles and Techniques (Stanford).

Two-layer neural networks



A. Pizurica, E016350 Artificial Intelligence (UGent) Spring 2025

Neural networks 22 /

A simple backpropagation example



$$Loss(x, y, \mathbf{w}) = (\mathbf{w} \cdot \mathbf{x} - y)^2$$

 $\mathbf{w} = [3, 1], \quad \mathbf{x} = [1, 2], y = 2$ backpropagation $\nabla_{\mathbf{w}} \text{Loss}(\mathbf{x}, y, \mathbf{w}) = [6, 12]$

Forward/backward values:

Forward: f_i is value for subexpression rooted at iBackward: $b_i = \frac{\partial Loss}{\partial f_i}$ is how f_i influences loss

Backpropagation:

- Forward pass: compute each f_i (from leaves to root)
- **2** Backward pass: compute each b_i (from root to leaves)

Adapted from M. Charikar and Koyejo: Artificial Intelligence: Principles and Techniques (Stanford).

Spring 2025

Backpropagation in neural networks explained



Cross-entropy loss

In deep learning, commonly we talk about minimizing cross-entropy loss

- \bullet Cross-entropy H(P,Q) is a measure of dissimilarity between the two distributions P and Q
- General definition:

$$H(P,Q) = \mathbb{E}_{z \sim P(z)}[-\log Q(z)] = -\int P(\mathbf{z}) \log Q(\mathbf{z}) d\mathbf{z}$$

- Typically: P is the true distribution over the training examples $P^*({\bf x},y)$, and Q is the predictive hypothesis $P(y|{\bf x},{\bf w})$
 - But we don't know $P^*(\mathbf{x}, y)$. We have access to its samples though!
 - So, approximate the expectation by the sum over the samples
 - Practical approach:

$$\mathbf{w}^* = \arg\min_{\mathbf{w}} - \sum_{j=1}^N \log P(y^{(j)} | \mathbf{x}^{(j)}, \mathbf{w}) = \arg\max_{\mathbf{w}} \sum_{j=1}^N \log P(y^{(j)} | \mathbf{x}^{(j)}, \mathbf{w})$$

Residual neural networks

- A popular approach to building very deep neural networks
- Instead of learning the desired mapping $h(\mathbf{x})$, the stacked nonlinear layers fit the residual $\mathcal{F}(\mathbf{x}) = h(\mathbf{x}) \mathbf{x}$. Hence, the original mapping recast to $\mathcal{F}(\mathbf{x}) + \mathbf{x}$
- It is easier to optimize the residual mapping than to optimize the original, unreferenced mapping
 - Think if an identity mapping were optimal, easier to push residual to zero than to fit identity mapping by a stack of nonlinear layers



Regularization in deep neural networks

Some common regularization approaches in deep learning include

- Weight decay: add a penalty $\lambda \sum_{i,j} w_{i,j}^2$ to the loss function
 - > Not straightforward to interpret the effect of weight decay in neural network
 - Common to use λ near 10^{-4}
- Dropout: deactivate a random chosen subset of units in each step of training



Does stochastic gradient descent (SGD) work for neural networks?



- For neural networks, optimization is hard
- In practice, SGD can work for neural nets much better than the theory predicts The gap between theory and practice not well understood yet!

Adapted from M. Charikar and Koyejo: Artificial Intelligence: Principles and Techniques (Stanford).

How to train neural networks



- Careful initialization (random noise, pre-training)
- Overparameterization (more hidden units than needed)
- Adaptive step sizes (AdaGrad, Adam)
- Don't let gradients vanish or explode!

Adapted from M. Charikar and Koyejo: Artificial Intelligence: Principles and Techniques (Stanford).

