

DEEP LEARNING FROM A STATISTICAL VIEWPOINT

SESSION 0: HOW, WHY AND WHAT IS DEEP LEARNING?

Tanguy Lefort

IMAG, Univ Montpellier, CNRS
LIRMM, Inria, Univ Montpellier, CNRS

Joseph Salmon

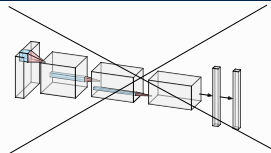
IMAG, Univ Montpellier, CNRS
Institut Universitaire de France (IUF)



UNIVERSITÉ DE
MONTPELLIER



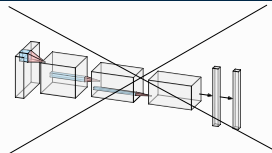
Inria



What **we are not going** to talk about:

- ▶ Zoology of neural networks (CNN, Transformers, etc.) : too many, more every other day!
- ▶ Very big neural networks (some models need thousands TPUv3-days to train⁽¹⁾ ⇒ neither practical nor theoretically well understood).

⁽¹⁾ A. Dosovitskiy et al. (2021). "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale". In: *arXiv:2010.11929 [cs]*.



What **we are not going** to talk about:

- ▶ Zoology of neural networks (CNN, Transformers, etc.) : too many, more every other day!
- ▶ Very big neural networks (some models need thousands TPUv3-days to train⁽¹⁾ \Rightarrow neither practical nor theoretically well understood).

What **we are going** to talk about:

- ▶ What is a feed-forward NN?
- ▶ Why do we use activation functions?
- ▶ Optimization with gradient methods and vanishing gradient

⁽¹⁾ A. Dosovitskiy et al. (2021). "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale". In: *arXiv:2010.11929 [cs]*.



Neural networks in parts

The perceptron
SVM and MLP

Activation functions

Optimization issues in MLP

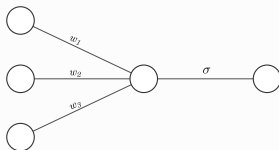
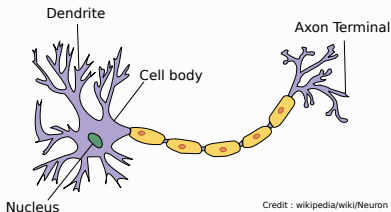
FROM THE NEURON TO THE BRAIN... ALMOST

NEURAL NETWORKS PARTS BY PARTS



The brain \simeq 86 billions neurons each connected up to 10K others.

The core element: a neuron



- ▶ Input: $x \in \mathbb{R}^d$
- ▶ Hidden layer: $\langle w, x \rangle + b$
- ▶ Output: $\theta = (w, b)$
 $\hat{f}(x, \theta) = \sigma(\langle w, x \rangle) \in \mathbb{R}^{out}$
- ▶ Training:
observations: $X = [x_1^T, \dots, x_n^T]^T$
binary labels (± 1): $y = (y_1, \dots, y_n)^T$



Neural networks in parts

The perceptron

SVM and MLP

Activation functions

Optimization issues in MLP

THE PERCEPTRON⁽²⁾

AN ELEMENTARY BRICK FOR NNs



```
def perceptron(X, y, w, b, n_iter, rho=1):
    for k in range(n_iter):
        for i in range(len(y)):
            if (X[i].T @ w + b) * y[i] <= 0: # error made
                w += rho * X[i] * y[i]
                b += rho * y[i]
    return w, b
```

Theorem: convergence⁽³⁾

The algorithm stops in a finite number of steps if the dataset (X, y) is separable, i.e., if $\exists \gamma > 0, w^{\text{sep}} \in \mathbb{R}^d$ such that $\langle w^{\text{sep}}, x_i \rangle + b > \gamma$.

Rem: Roughly # steps $\propto \max(\|x_i\|)^2 \cdot \|w^{\text{sep}}\|^2 / \gamma^2$

⁽²⁾ F. Rosenblatt (1958). *The perceptron: a theory of statistical separability in cognitive systems (Project Para)*. Cornell Aeronautical Laboratory

⁽³⁾ A. B. Novikoff (1963). *On convergence proofs for perceptrons*. Tech. rep. STANFORD RESEARCH INST MENLO PARK CA



What is the perceptron?

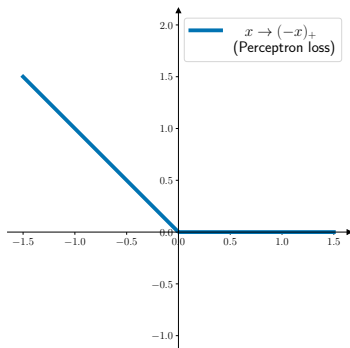
Modern answer: it is a linear binary classifier learnt with stochastic gradient descent (SGD) for the (perceptron) hinge loss and a fixed step size

(Perceptron) hinge loss:

$$\mathcal{L}(w, b) = \frac{1}{n} \sum_{i=1}^n \left(-y_i \cdot (\langle w, x_i \rangle + b) \right)_+$$

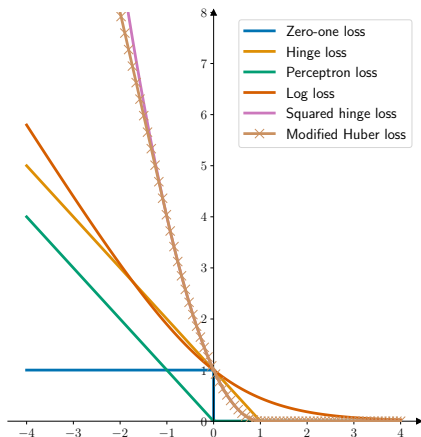
with

$$(x)_+ = \text{ReLU}(x) := \max(0, x)$$





- ▶ the Perceptron might never stop (*cf.* XOR problem)
- ▶ minimization of an 1-homogeneous function whose optimal value is trivial ($w = 0, b = 0$)
- ▶ SGD requires decreasing step size and does not converge for non-smooth case with fixed step size



Credit:

https://scikit-learn.org/dev/auto_examples/linear_model/plot_sgd_loss_functions.html



Many linear/affine predictors (classifiers) can leverage an expression:

$$\sigma(\langle w, x \rangle + b)$$

- ▶ σ the activation function
- ▶ w the weight (of the neuron)
- ▶ b the bias (of the neuron)



Many linear/affine predictors (classifiers) can leverage an expression:

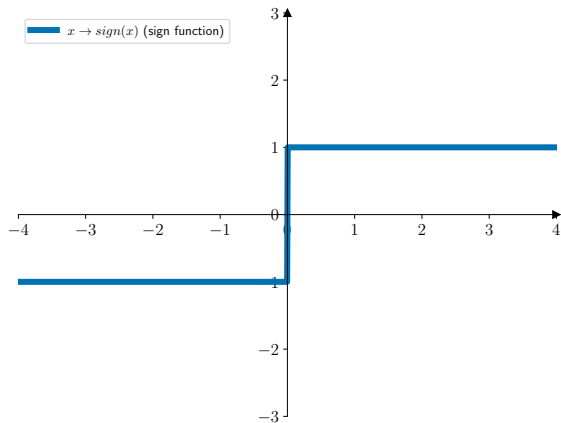
$$\sigma(\langle w, x \rangle + b)$$

- ▶ σ the activation function
- ▶ w the weight (of the neuron)
- ▶ b the bias (of the neuron)

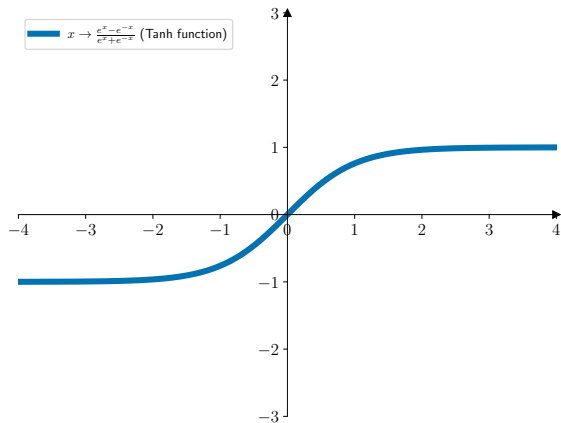
For binary classification, we can use as **activation** function σ

- ▶ the sign function: $x \mapsto \text{sign}(x)$,
- ▶ the ReLU function: $x \mapsto (x)_+$,
- ▶ the tanh function: $x \mapsto \frac{e^x - e^{-x}}{e^x + e^{-x}}$
- ▶ the sigmoid: $x \mapsto \frac{1}{1 + e^{-x}}$ (with in mind $\mathbb{P}(y_i = 1 | x_i) = \frac{1}{1 + e^{-(\langle w, x_i \rangle + b)}}$)

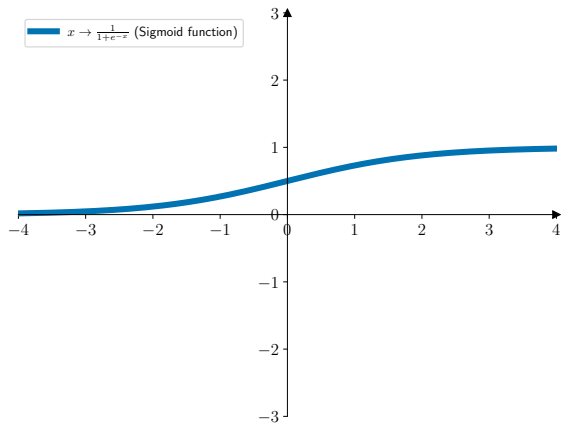
ACTIVATION FUNCTION



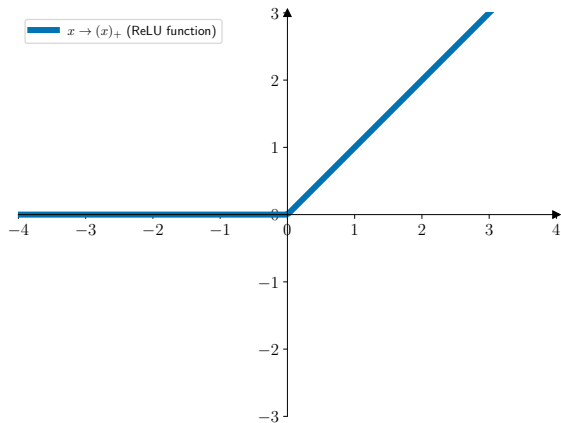
ACTIVATION FUNCTION



ACTIVATION FUNCTION



ACTIVATION FUNCTION



BACKPROPAGATION ON THE PERCEPTRON+RELU

BY HAND VS AUTOGRAD



$$\hat{L}(\theta) = \frac{1}{n} \sum_i (y_i - f(x_i, \theta))^2 = \frac{1}{n} \sum_i (y_i - (w^\top x_i + b)_+)^2$$

Denoting $u = w^\top x_i + b, v = (y_i - u)$:

$$\begin{aligned} \frac{\partial \hat{L}}{\partial w} &= \frac{\partial \hat{L}}{\partial v} \frac{\partial v}{\partial u} \frac{\partial u}{\partial w} \\ &= -\frac{2}{n} \sum_i (y_i - \hat{y}_i) x_i^\top \mathbb{1}(u > 0) \end{aligned}$$

$$\begin{aligned} \frac{\partial \hat{L}}{\partial b} &= \frac{\partial \hat{L}}{\partial v} \frac{\partial v}{\partial u} \frac{\partial u}{\partial b} \\ &= -\frac{2}{n} \sum_i (y_i - \hat{y}_i) \mathbb{1}(u > 0) \end{aligned}$$

Using AutoGrad ⁽⁴⁾:

```
model = nn.Sequential(  
    nn.Linear(2, 1), # input 2D, output 1D  
    nn.ReLU() # ReLU activation  
)  
loss = nn.MSELoss() # cost function  
y_pred = model(x) # prediction  
val_loss = loss(y, y_pred)  
val_loss.backward() # backpropagation  
model[0].weight.grad # grad wrt w  
model[0].bias.grad # grad wrt b
```

Rem:

Gradient tree can take a lot of memory

(4) A. Paszke et al. (2019). "PyTorch: An Imperative Style, High-Performance Deep Learning Library". In: arXiv:1912.01703 [cs, stat]



Neural networks in parts

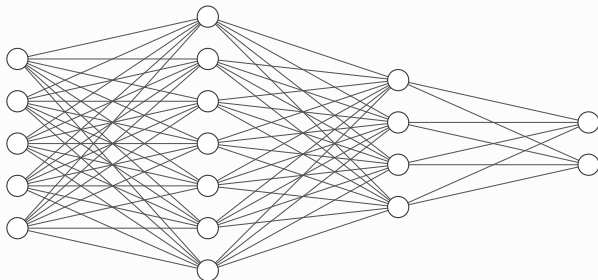
The perceptron

SVM and MLP

Activation functions

Optimization issues in MLP

- ▶ Feed-forward NN \supset MLP (Multi-Layer Perceptron)



Idea: "There is strength in numbers"

- ▶ each hidden layer : perceptrons = linear transformations + activations
- ▶ can handle multiclass cases (*e.g.*, with softmax output activation)

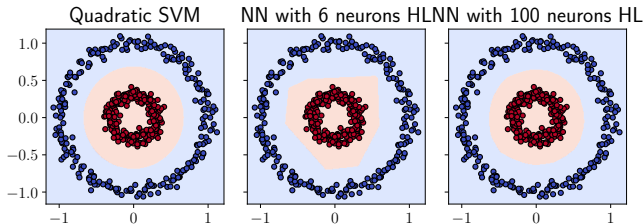
SVM AND MLP CLASSIFICATION

1 HIDDEN LAYER IN \mathbb{R}^2



Decision boundary with quadratic kernel SVM in \mathbb{R}^2 :

$$\left\{ x \mid \sum_i \alpha_i y_i k(x_i, x) + b = 0 \right\}, \quad k(x, x') = (x^\top x' + c)^2$$



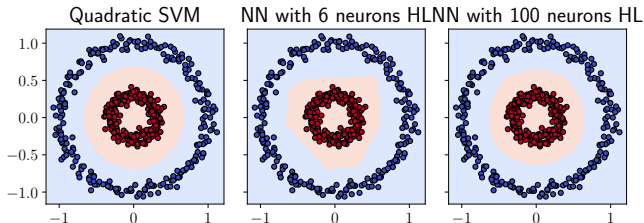
SVM AND MLP CLASSIFICATION

1 HIDDEN LAYER IN \mathbb{R}^2



Decision boundary with quadratic kernel SVM in \mathbb{R}^2 :

$$\left\{ x \mid \sum_i \alpha_i y_i k(x_i, x) + b = 0 \right\}, \quad k(x, x') = (x^\top x' + c)^2$$



Decision boundary for MLP

- ▶ MLP compute probabilities and is scalable (for labels & features)
- ▶ # linear separations = # neurons in hidden layer
- ▶ **Warning: overfitting with MLP is common**



Neural networks in parts

Activation functions

- Step function from RELU
- Approximations with NNs
- Universal theorem
- XOR problem

Optimization issues in MLP



Neural networks in parts

Activation functions

- Step function from RELU

- Approximations with NNs

- Universal theorem

- XOR problem

Optimization issues in MLP



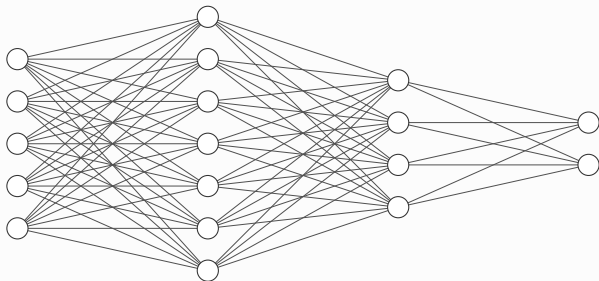
For K layers: matrices (weights) W_1, \dots, W_K and vectors (bias) b_1, \dots, b_K

- ▶ **without** activations:

$$\begin{aligned}y &= W_K x_{K-1} + b_K = W_K(W_{K-1}x_{K-2} + b_{K-1}) + b_K \\ &= W_K W_{K-1} x_{K-2} + W_K b_{K-1} + b_K \\ &= \dots \\ &= W_K W_{K-1} \dots W_1 x_0 + \mathbf{c} = \mathbf{W}x + \mathbf{c} \quad (\text{affine!})\end{aligned}$$

- ▶ **with** activations σ_K at layer k :

$$\begin{aligned}y &= \sigma_K(W_K z_{K-1} + b_K) = \sigma_K(z_K) \\ &= \sigma_K(W_K \sigma_{K-1}(W_{K-1} z_{K-2} + b_{K-1}) + b_K) \\ &= \sigma_K(W_K \sigma_{K-1}(W_{K-1} \sigma_{K-2}(\dots) + b_{K-1}) + b_K) \quad (\text{not affine!})\end{aligned}$$



Here: $\hat{f} : \mathbb{R}^5 \rightarrow \mathbb{R}^2$

$$x \mapsto \sigma_3(W_3 \sigma_2(W_2 \sigma_1(W_1 x + b_1) + b_2) + b_3)$$

with: $W_1 \in \mathbb{R}^{7 \times 5}, \quad b_1 \in \mathbb{R}^7$

$$W_2 \in \mathbb{R}^{4 \times 7}, \quad b_2 \in \mathbb{R}^4$$

$$W_3 \in \mathbb{R}^{2 \times 4}, \quad b_3 \in \mathbb{R}^2$$



Neural networks in parts

Activation functions

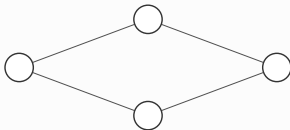
Step function from RELU

Approximations with NNs

Universal theorem

XOR problem

Optimization issues in MLP



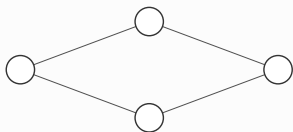
```
class Two_Perceptron(torch.nn.Module):
    def __init__(self):
        super(Two_Perceptron, self).__init__()
        self.in_fc = nn.Linear(1, 2)
        self.fc = nn.Linear(2, 1)
        self.sig = nn.Sigmoid()

    def forward(self, x):
        out = self.in_fc(x)
        out = self.sig(out)
        out = self.fc(out)
        return out
```



Simple 2 neurons-1 hidden layer and sigmoid activation:

$$x \mapsto w_3\sigma(w_1x + b_1) + w_4\sigma(w_2x + b_2)$$



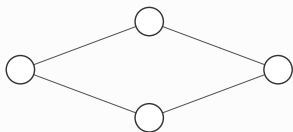
FROM THE LINEARITY TO THE NON-LINEARITY

STEP FUNCTION WITH SIGMOID

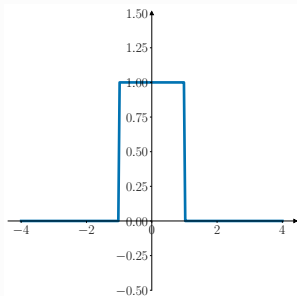


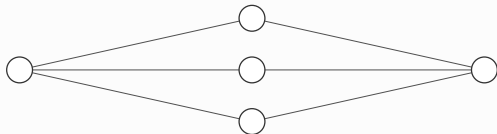
Simple 2 neurons-1 hidden layer and sigmoid activation:

$$x \mapsto w_3\sigma(w_1x + b_1) + w_4\sigma(w_2x + b_2)$$



```
model = Two_Perceptron()
model.in_fc.weight.data =
    torch.tensor([[1000.], [1000.]])
model.in_fc.bias.data =
    torch.tensor([0., -1000.])
model.fc.weight.data =
    torch.tensor([[1., -1.]])
model.fc.bias.data = torch.tensor([0.]
```



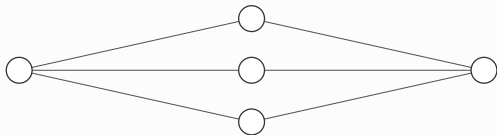


```
class Three_Perceptron(torch.nn.Module):
    def __init__(self):
        super(Three_Perceptron, self).__init__()
        self.in_fc = nn.Linear(1, 3)
        self.fc = nn.Linear(3, 1)
        self.sig = nn.ReLU()

    def forward(self, x):
        out = self.in_fc(x)
        out = self.sig(out)
        out = self.fc(out)
        return out
```

Simple 3 neurons-1 hidden layer and ReLu activation:

$$y = f(x, \theta) = \sigma(x + 1) - 2\sigma(x) + \sigma(x - 1)$$



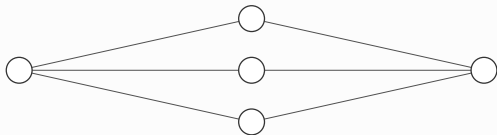
FROM THE LINEARITY TO THE NON-LINEARITY

TRIANGLE FUNCTION WITH RELU

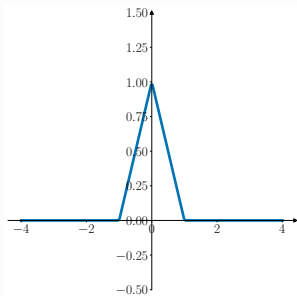


Simple 3 neurons-1 hidden layer and ReLU activation:

$$y = f(x, \theta) = \sigma(x + 1) - 2\sigma(x) + \sigma(x - 1)$$



```
model = Three_Perceptron()
model.in_fc.weight.data =
    torch.tensor([[1.], [1.], [1.]])
model.in_fc.bias.data =
    torch.tensor([1, 0., -1])
model.fc.weight.data =
    torch.tensor([[1., -2, 1.]])
model.fc.bias.data = torch.tensor([0.]
```





Neural networks in parts

Activation functions

Step function from RELU

Approximations with NNs

Universal theorem

XOR problem

Optimization issues in MLP



Notation: $\mathcal{C}([0, 1]^d) := \{f : [0, 1]^d \rightarrow \mathbb{R}, f \text{ continuous}\}$

Under the condition that σ is continuous with $\begin{cases} \lim_{t \rightarrow -\infty} \sigma(t) = 0, \\ \lim_{t \rightarrow +\infty} \sigma(t) = 1 \end{cases}$ then

Theorem⁽⁵⁾

For every function $g \in \mathcal{C}([0, 1]^d)$ and $\epsilon > 0$ there exists a feedforward neural network f such that $\|g - f\|_\infty < \epsilon$, i.e.,

$$f(x) = \sum_{j=1}^N \alpha_j \sigma(\langle w_j, x \rangle + b_j)$$

for some integer N , some $w_j \in \mathbb{R}^d$ and $b_j \in \mathbb{R}$.

⁽⁵⁾ G. Cybenko (1989). "Approximation by superpositions of a sigmoidal function". en. In: *Mathematics of Control, Signals and Systems 2.4*, pp. 303–314

⁽⁶⁾ L. Devroye, L. Györfi, and G. Lugosi (1996). *A probabilistic theory of pattern recognition*. Vol. 31. Applications of Mathematics (New York). Springer-Verlag

⁽⁷⁾ Z. Lu et al. (2017). *The Expressive Power of Neural Networks: A View from the Width*.

Notation: $\mathcal{C}([0, 1]^d) := \{f : [0, 1]^d \rightarrow \mathbb{R}, f \text{ continuous}\}$

Under the condition that σ is continuous with $\begin{cases} \lim_{t \rightarrow -\infty} \sigma(t) = 0, \\ \lim_{t \rightarrow +\infty} \sigma(t) = 1 \end{cases}$ then

Theorem⁽⁵⁾

For every function $g \in \mathcal{C}([0, 1]^d)$ and $\epsilon > 0$ there exists a feedforward neural network f such that $\|g - f\|_\infty < \epsilon$, i.e.,

$$f(x) = \sum_{j=1}^N \alpha_j \sigma(\langle w_j, x \rangle + b_j)$$

for some integer N , some $w_j \in \mathbb{R}^d$ and $b_j \in \mathbb{R}$.

Rem: one can adapt the proof⁽⁶⁾ to ReLU with the "triangle" approximation.

Rem: refined control on width/depth are possible⁽⁷⁾

⁽⁵⁾ G. Cybenko (1989). "Approximation by superpositions of a sigmoidal function". en. In: *Mathematics of Control, Signals and Systems 2.4*, pp. 303–314

⁽⁶⁾ L. Devroye, L. Györfi, and G. Lugosi (1996). *A probabilistic theory of pattern recognition*. Vol. 31. Applications of Mathematics (New York). Springer-Verlag

⁽⁷⁾ Z. Lu et al. (2017). *The Expressive Power of Neural Networks: A View from the Width*.



What the theorem says and **does not say**:

- ▶ a MLP will approximate the function **not learn (overfitting for example)**



What the theorem says and **does not say**:

- ▶ a MLP will approximate the function **not learn (overfitting for example)**
- ▶ there always exists a large NN **but might be exponentially large.**



What the theorem says and **does not say**:

- ▶ a MLP will approximate the function **not learn (overfitting for example)**
- ▶ there always exists a large NN **but might be exponentially large.**

- ▶ a single hidden layer feedforward NN can represent any function **BUT** in practice, require a large single layer (extremely time and memory consuming). **Depth is a key element in NN.**



- ▶ Widgets Runge and Gibbs phenomenon
- ▶ Show: universal approx video



Kolmogorov-Arnold theorem⁽³⁾: For any $g : [0, 1]^d \rightarrow \mathbb{R}$ continuous there exist univariate continuous functions such that:

$$g(x_1, \dots, x_d) = \sum_{q=1}^{2d} h_q \left(\sum_{p=1}^d \psi_{p,q}(x_p) \right) .$$

- ▶ *Kolmogorov's theorem is irrelevant*⁽⁴⁾: smoothness issue (learning) and h_q highly dependent of g (no parametric form).

⁽³⁾ A. N. Kolmogorov (1956). "The representation of continuous functions of several variables by superpositions of continuous functions of a smaller number of variables". In: *Doklady Akademii Nauk SSSR* 108.2, pp. 179–182.

⁽⁴⁾ F. Girosi and T. Poggio (1989). "Representation properties of networks: Kolmogorov's theorem is irrelevant". In: *Neural Computation* 1.4, pp. 465–469.

⁽⁵⁾ V. Kůrková (1991). "Kolmogorov's Theorem Is Relevant". In: *Neural Computation* 3.4, pp. 617–622.

⁽⁶⁾ J. Schmidt-Hieber (2021). "The Kolmogorov–Arnold representation theorem revisited". In: *Neural Networks* 137, pp. 119–126.



Kolmogorov-Arnold theorem⁽³⁾: For any $g : [0, 1]^d \rightarrow \mathbb{R}$ continuous there exist univariate continuous functions such that:

$$g(x_1, \dots, x_d) = \sum_{q=1}^{2d} h_q \left(\sum_{p=1}^d \psi_{p,q}(x_p) \right) .$$

- ▶ *Kolmogorov's theorem is irrelevant*⁽⁴⁾: smoothness issue (learning) and h_q highly dependent of g (no parametric form).
- ▶ *Kolmogorov's theorem is relevant*⁽⁵⁾: construction of smoother functions and reduction of the problem.

⁽³⁾ A. N. Kolmogorov (1956). "The representation of continuous functions of several variables by superpositions of continuous functions of a smaller number of variables". In: *Doklady Akademii Nauk SSSR* 108.2, pp. 179–182.

⁽⁴⁾ F. Girosi and T. Poggio (1989). "Representation properties of networks: Kolmogorov's theorem is irrelevant". In: *Neural Computation* 1.4, pp. 465–469.

⁽⁵⁾ V. Kůrková (1991). "Kolmogorov's Theorem Is Relevant". In: *Neural Computation* 3.4, pp. 617–622.

⁽⁶⁾ J. Schmidt-Hieber (2021). "The Kolmogorov–Arnold representation theorem revisited". In: *Neural Networks* 137, pp. 119–126.



Kolmogorov-Arnold theorem⁽³⁾: For any $g : [0, 1]^d \rightarrow \mathbb{R}$ continuous there exist univariate continuous functions such that:

$$g(x_1, \dots, x_d) = \sum_{q=1}^{2d} h_q \left(\sum_{p=1}^d \psi_{p,q}(x_p) \right) .$$

- ▶ *Kolmogorov's theorem is irrelevant*⁽⁴⁾: smoothness issue (learning) and h_q highly dependent of g (no parametric form).
- ▶ *Kolmogorov's theorem is relevant*⁽⁵⁾: construction of smoother functions and reduction of the problem.
- ▶ *The Kolmogorov-Arnold representation theorem revisited*⁽⁶⁾: simplify hypothesis and approximate the inner function with deep ReLU NN.

⁽³⁾ A. N. Kolmogorov (1956). "The representation of continuous functions of several variables by superpositions of continuous functions of a smaller number of variables". In: *Doklady Akademii Nauk SSSR* 108.2, pp. 179–182.

⁽⁴⁾ F. Girosi and T. Poggio (1989). "Representation properties of networks: Kolmogorov's theorem is irrelevant". In: *Neural Computation* 1.4, pp. 465–469.

⁽⁵⁾ V. Kůrková (1991). "Kolmogorov's Theorem Is Relevant". In: *Neural Computation* 3.4, pp. 617–622.

⁽⁶⁾ J. Schmidt-Hieber (2021). "The Kolmogorov–Arnold representation theorem revisited". In: *Neural Networks* 137, pp. 119–126.

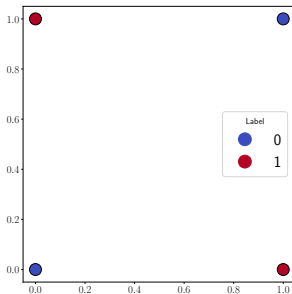


Neural networks in parts

Activation functions

- Step function from RELU
- Approximations with NNs
- Universal theorem
- XOR problem

Optimization issues in MLP



$$X = \begin{pmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{pmatrix} \quad y = \begin{pmatrix} 0 \\ 1 \\ 1 \\ 0 \end{pmatrix}$$

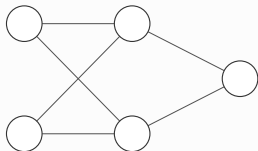
$$\hat{L}(\theta) = \frac{1}{4} \sum_i (y_i - f(x_i, \theta))^2$$

$$f(x_i, \theta) = \mathbb{1}(f(x) > 0.5)$$

Non separable case: the XOR problem

The XOR problem with this configuration can not be solved with a linear separator $x \mapsto w^\top x + b$. (To show: video xor perceptron)

1st order conditions: $\hat{w} = (0, 0)^\top$, and $\hat{b} = \frac{1}{2} \implies f(x_i, \theta) = \frac{1}{2}$.



$$\hat{L}(\theta) = \frac{1}{4} \sum_i (y_i - f(x_i, \theta))^2,$$

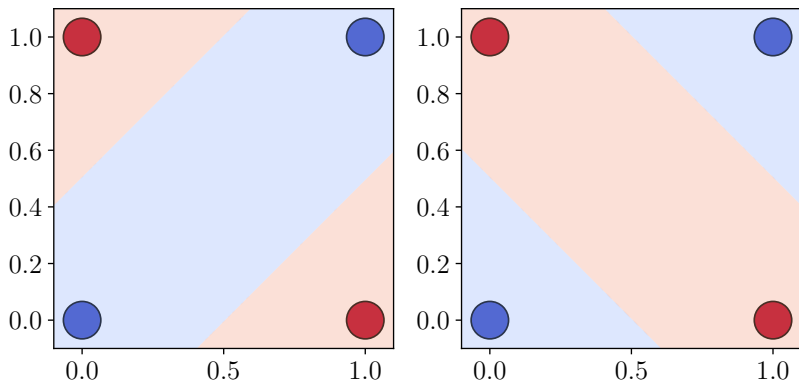
$$f(x_i, \theta) = \langle w_2, h \rangle + b_2, \quad h = (W_1 x_i + b_1)_+$$

Multiple solutions give $\hat{L}(\theta) = 0$

- ▶ if $b_1 = (0, 0)$, and $b_2 = 0$:
 - ▶ $w_2 = (1, 1)$ and $W_1 = \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix}$
 - ▶ $w_2 = (1, 1)$ and $W_1 = \begin{pmatrix} -1 & 1 \\ 1 & -1 \end{pmatrix}$
 - ▶ for $D = \text{diag}(d_1, d_2)$, $d_1, d_2 > 0$ taking $\tilde{w}_2 = D^{-1}w_2$, $\tilde{W}_1 = DW_1$ with w_2 and W_1 leading to a global minimum.
- ▶ if $b_1 = (1, -1)$ and $b_2 = 1$ and $W_1 = \begin{pmatrix} -1 & -1 \\ 1 & 1 \end{pmatrix}$ and $w_2 = (-1, -1)$
- ▶ ...

THE XOR PROBLEM

VISUALIZATION



To show: video xor1 HL2 neurons



Neural networks in parts

Activation functions

Optimization issues in MLP

Backpropagation in MLP

Vanishing gradients solutions for optimization

Introduction to resnets

Learn a neural network by minimizing your (favorite) loss \mathcal{L} on the training set:

$$\begin{aligned} \min_{\substack{W_1, \dots, W_K \\ b_1, \dots, b_K}} \mathcal{L}(\hat{f}(x_i), y_i) \\ \text{s.t. } \hat{f} = \sigma_K(W_K \sigma_{K-1}(W_{K-1} \sigma_{K-2}(\dots) + b_{K-1}) + b_K) \end{aligned}$$

Algorithm choice: use SGD and/or variants (with mini-batch) if possible with a GPU (tailored for fast matrix/vector operations)



Neural networks in parts

Activation functions

Optimization issues in MLP

Backpropagation in MLP

Vanishing gradients solutions for optimization

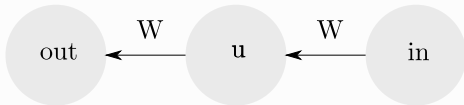
Introduction to resnets

VANISHING AND EXPLODING GRADIENT

LONG TERM DEPENDANCY⁽⁷⁾

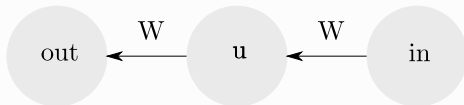


Chain rule: $\frac{\partial out}{\partial in} = \frac{\partial out}{\partial u} \frac{\partial u}{\partial in}$



⁽⁷⁾ Ian Goodfellow, Yoshua Yoshua, and Aaron Courville (2016). *Deep Learning*. MIT Press

Chain rule: $\frac{\partial out}{\partial in} = \frac{\partial out}{\partial u} \frac{\partial u}{\partial in}$



Risk with too deep NN: if $W = P \text{diag}(\lambda_i) P^{-1} \implies W^k = P \text{diag}(\lambda_i^k) P^{-1}$.

- ▶ $\lambda_i > 1 \implies$ exploding gradient,
- ▶ $\lambda_i < 1 \implies$ vanishing gradient.

Diagnose the vanishing / exploding gradient

- ▶ necessary condition (vanishing): weights distribution barely changing (**not sufficient!**)
- ▶ learning curve very unstable or not decreasing (**also not sufficient**)

⁽⁷⁾ Ian Goodfellow, Yoshua Yoshua, and Aaron Courville (2016). *Deep Learning*. MIT Press



Neural networks in parts

Activation functions

Optimization issues in MLP

Backpropagation in MLP

Vanishing gradients solutions for optimization

Introduction to resnets



Sigmoid activation induce gradient issues

Very low variations in the tails \implies flat gradient!



Sigmoid activation induce gradient issues

Very low variations in the tails \implies flat gradient!

RELU activations = our rescue?



Sigmoid activation induce gradient issues

Very low variations in the tails \implies flat gradient!

RELU activations = our rescue? **NO** (even if the internet often says so).

- ▶ too much dying RELUs \implies zero-out layers
- ▶ same for ℓ_2 and ℓ_1 penalties (but great to prevent exploding gradients instead of gradient clipping)

Currently there is no savior, just band-aids

Batch normalizations, weight initializations, leaky RELUs, ... = help

Leave MLPs and go to **Residuals networks**.



Neural networks in parts

Activation functions

Optimization issues in MLP

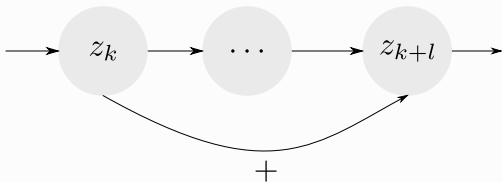
Backpropagation in MLP

Vanishing gradients solutions for optimization

Introduction to resnets

Idea: add "skip connections" to avoid flatlining gradients

$$z_{k+l} = \sigma_{k+l-1}(z_{k+l-1}) + z_k \Rightarrow \frac{\partial z_{k+l}}{\partial z_k} = \frac{\partial \sigma_{k+l-1}(z_{k+l-1})}{\partial z_k} + 1$$









Let's address the dimensionality









- ▶ dimension is an issue in CNNs (and any NN using pooling)
- ▶ authors suggested using a linear projection of z_k if needed.
- ▶ l is kept fairly small (2 or 3).



- ▶ New techniques come up very often,
- ▶ Most of the time, understanding the basics is enough to scrap the ideas and use them in your situation.
- ▶ Lots of empirical results and strategies adopted, the *why does it (not) work* is sometimes left behind (lots of research needed),
- ▶ Time vs Cost \implies essential factor in (re)using lots of NN.

We are now ready to tackle both theoretically and practically the next sessions

-  Cybenko, G. (Dec. 1989). “Approximation by superpositions of a sigmoidal function”. en. In: *Mathematics of Control, Signals and Systems* 2.4, pp. 303–314.
-  Devroye, L., L. Györfi, and G. Lugosi (1996). *A probabilistic theory of pattern recognition*. Vol. 31. Applications of Mathematics (New York). New York: Springer-Verlag.
-  Dosovitskiy, A. et al. (June 2021). “An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale”. In: *arXiv:2010.11929 [cs]*. arXiv: 2010.11929.
-  Girosi, F. and T. Poggio (1989). “Representation properties of networks: Kolmogorov’s theorem is irrelevant”. In: *Neural Computation* 1.4, pp. 465–469.
-  He, K. et al. (Dec. 2015). “Deep Residual Learning for Image Recognition”. In: *arXiv:1512.03385 [cs]*. arXiv: 1512.03385.
-  Ian Goodfellow, Yoshua Bengio, and Aaron Courville (2016). *Deep Learning*. MIT Press.

-  Kolmogorov, A. N. (1956). “The representation of continuous functions of several variables by superpositions of continuous functions of a smaller number of variables”. In: *Doklady Akademii Nauk SSSR* 108.2, pp. 179–182.
-  Kůrková, V. (Dec. 1991). “Kolmogorov’s Theorem Is Relevant”. In: *Neural Computation* 3.4, pp. 617–622.
-  Lu, Z. et al. (2017). *The Expressive Power of Neural Networks: A View from the Width*.
-  Minsky, M. L. and S. A. Papert (1969). *Perceptrons. An Introduction to Computational Geometry*. 1969, Expanded. Cambridge, MA: MIT Press.
-  Novikoff, A. B. (1963). *On convergence proofs for perceptrons*. Tech. rep. STANFORD RESEARCH INST MENLO PARK CA.
-  Paszke, A. et al. (Dec. 2019). “PyTorch: An Imperative Style, High-Performance Deep Learning Library”. In: *arXiv:1912.01703 [cs, stat]*. arXiv: 1912.01703.
-  Rosenblatt, F. (1958). *The perceptron: a theory of statistical separability in cognitive systems (Project Para)*. Cornell Aeronautical Laboratory.
-  Schmidt-Hieber, J. (2021). “The Kolmogorov–Arnold representation theorem revisited”. In: *Neural Networks* 137, pp. 119–126.