### **Deep Learning**

Lecture 2: Machine learning recap, history of neural networks and the main building blocks

Prof. Stéphane Gaïffas https://stephanegaiffas.github.io







# **Machine learning recap**

# **Supervised learning**

We have data with labeled samples  $(x_1, y_1), \ldots, (x_n, y_n)$ 

- Each  $(x_i,y_i)\in \mathcal{X} imes \mathcal{Y}$  have an unknown joint distribution P
- The  $x_i$  are input features. We have often  $\mathcal{X} = \mathbb{R}^d$ .
- The  $y_i$  are output labels. We have often  $\mathcal{Y} = \mathbb{R}$  (a regression problem, the labels do have a natural ordering) or  $\mathcal{Y} = \{1, \ldots, K\}$  (a K-class classification problem, the labels are categorical and do not have a natural ordering)

The samples are supposed independent and identically distributed (i.i.d)

The training data can be of any finite size n.

In general, we do not have any prior information about P.

In most cases,  $x_i$  is a vector, but it could be an image, a piece of text, a sample of sound, a video

## **Supervised learning**

A predictor is a measurable function

 $f:\mathcal{X}
ightarrow \widehat{\mathcal{Y}}$ 

where

$$\widehat{\mathcal{Y}} = \mathcal{Y} \quad ext{ or } \quad \widehat{\mathcal{Y}} = \Delta(\mathcal{Y}),$$

with  $\Delta(\mathcal{Y})$  the set of distributions on  $\mathcal{Y}$ . Indeed, in classification problems, we often prefer to predict the distribution of y given x

$$p(y|x) = \mathbb{P}[Y=y|X=x] \in \Delta(\mathcal{Y})$$

instead of just predicting  $y \in \mathcal{Y}$ .

#### Training

Training or learning a predictor f is done using training data using the empirical risk minimization principle

# **Empirical risk minimization**

We train the predictor by minimization of the the empirical risk or training error

$$R_n(f) = rac{1}{n}\sum_{i=1}^n \ell(y_i,f(x_i))$$

over  $f\in \mathcal{F}$  where

•  $\mathcal{F}$  is a parametrized set of functions. In linear methods we use the set of functions  $\mathcal{F} = \mathcal{F}_{ ext{linear}} = \{x \mapsto x^{ op} w : w \in \mathbb{R}^d\}.$ 

• 
$$\,\ell:\mathcal{Y} imes \widehat{\mathcal{Y}} o\mathbb{R}\,$$
 is a loss function

#### **Examples**

- For linear least-squares we take  ${\cal F}={\cal F}_{
  m linear}$  and use the square loss  $\ell(y,y')=rac{1}{2}(y-y')^2$
- For logistic regression (for binary classification where  $\mathcal{Y} = \{-1, 1\}$ ) we take  $\mathcal{F} = \mathcal{F}_{ ext{linear}}$  and the logistic loss  $\ell(y, y') = \log(1 + e^{-yy'})$

## Generalization

The aim of a predictor is generalization. We want the generalization error

$$R(f) = \mathbb{E}_{(x,y) \sim P}ig[\ell(y,f(x))ig]$$

to be as small as possible.

- The empirical risk  $R_n(f)$  is used as a way of approximating the generalization error using the training samples
- Most machine learning algorithms, including neural networks, use empirical risk minimization
- But only minimizing the training error usually leads to overfitting, most methods require cross-validation for hyper-parameters tuning

Of course, many approaches: linear methods, kernels, k-NN, ensemble methods (boosting, forests), etc. and unsupervised methods

### **Scikit-learn cheat-sheet**



### **History of neural networks**

# **Modeling neurons**

- The idea of neural networks began as a model of how neurons work in the brain.
- Called connectionism and used connected circuits to simulate intelligent behavior



Figure I ORGANIZATION OF THE MARK I PERCEPTRON

- In 1943, portrayed with a simple electrical circuit by neurophysiologist Warren McCulloch and mathematician Walter Pitts.
- Donald Hebb took the idea further by proposing that neural pathways strengthen over each successive use, especially between neurons that tend to fire at the same time.

"The organization of behavior: A neuropsychological theory", Hebb 2005

<sup>&</sup>quot;A logical calculus of the ideas immanent in nervous activity", McCulloch and Pitts 1943

## **Modeling neurons**



- Around 50', Frank Rosenblatt, a psychologist at Cornell, was working on understanding the comparatively simpler decision systems present in the eye of a fly, which underlie and determine its flee response.
- In an attempt to understand and quantify this process, he proposed the idea of a Perceptron in 1958, calling it Mark I Perceptron.
- It was a system with a simple input output relationship, modeled on a McCulloch-Pitts neuron to explain the complex decision processes in a brain using a linear threshold gate.

## **First perceptron**

A McCulloch-Pitts neuron takes in inputs, takes a weighted sum and returns 0 if the result is below threshold and 1 otherwise.

$$f(x) = egin{cases} 1 & ext{if } \sum_i w_i x_i + b \geq 0 \ 0 & ext{otherwise} \end{cases}$$

A model motivated by biology, with  $w_i$  being synaptic weights and  $x_i$  and f firing rates.



- The Mark I Perceptron was the first implementation of the perceptron algorithm.
- The machine was connected to a camera that used 20x20 photocells to produce a 400-pixel image.





The Mark I Percetron (Frank Rosenblatt).

Perceptron Research from the 50's & 60's, clip



#### The Perceptron

## **Perceptron model and algorithm**

It uses the following model

$$f(x) = \sigma \Big(\sum_j w_j x_j + b \Big)$$

where  $x_j$  are inputs,  $w_j$  are weights, b is a bias and activation  $\sigma(z) = \mathbf{1}(z > 0)$ 

• How to find (w, b)?

#### The perceptron algorithm (first iterative learning algorithm)

To ease notations, put  $w = [w_1 \cdots w_d \; b]$  and  $x_i = [x_i \; 1]$ 

- Start with  $w \leftarrow 0$
- Repeat over all samples:

1. if 
$$y_i \; x_i^ op w < 0$$
 put  $w \leftarrow w + y_i x_i$ 

2. otherwise do not modify  $\boldsymbol{w}$ 

## **Perceptron model and algorithm**

- This procedure can be seen as a stochastic gradient descent method
- Indeed, a sample i is missclassified if  $y_i x_i^ op w < 0$  (labels -1 or +1), so we can consider the goodness-of-fit

$$F(w) = \sum_{i\,:\,y_i x_i^ op w < 0} F_i(w) \quad ext{where} \quad F_i(w) = -y_i x_i^ op w$$

• But  $abla F_i(w) = -y_i x_i$  so that a stochastic gradient descent step write

$$w \leftarrow w - \eta 
abla F_i(w) = w + \eta y_i x_i$$

where we sample uniformly at sample i among  $\{i: y_i x_i^{ op} w < 0\}$ .

• We use  $\eta = 1$  in the Perceptron algorithm

### **Perceptron model and algorithm**

Let  $R = \max_i |x_i|$  and  $w^{\star}$  be the optimal hyperplane of margin

 $\gamma = \min_i y_i x_i^ op w^\star$ 

with  $|w^{\star}| = 1$ . We have the following theorem.



#### Theorem (Block 1962, Novikoff 1963)

Assume that the training dataset  $(x_1, y_1), \ldots, (x_n, y_n)$  is linearly separable (meaning that  $\gamma > 0$ ). Put  $w_0 \leftarrow 0$ , then the number of steps in the perceptron algorithm is bounded by

$$k \leq rac{1+R^2}{\gamma^2}.$$

Exercice: prove it ! (pprox 10 lines)

## History

• The Perceptron project led by Rosenblatt was funded by the US Office of Naval Research.

"The Navy revealed the embryo of an electronic computer today that it expects will be able to walk, talk, see, write, reproduce itself and be conscious of its existence. Later perceptrons will be able to recognize people and call out their names and instantly translate speech in one language to speech and writing in another language, it was predicted."

Press conference, 7 July 1958, New York Times.

• For an extensive study of the perceptron, see [Principles of neurodynamics. perceptrons and the theory of brain mechanisms, Rosenblatt 1961]

# **Moving forward: (M)ADALINE**

• In 1959 at Stanford, Bernard Widrow and Marcian Hoff developed AdaLinE (ADAptive LINear Elements) and MAdaLinE (Multiple AdaLinE), the latter being the first network successfully applied to a real world problem



Main differences with the perceptron

- The loss is the square difference between the sum of the weighted input and the output
- The optimization procedure is a gradient descent: all computations are trivial since the weighted sum is linear as a function of weights

### MADALINE

Many Adalines: network with one hidden layer composed of many Adaline units. [1]



#### Applications

- Speech and pattern recognition [2]
- Weather forecasting [3]
- Adaptive filtering and adaptive signal processing [4]

<sup>[1] &</sup>quot;Madaline Rule II: a training algorithm for neural networks", Winter and Widrow 1988

<sup>[2] &</sup>quot;Real-Time Adaptive Speech-Recognition System", Talbert et al. 1963

<sup>[3] &</sup>quot;Application of the adaline system to weather forecasting", Hu 1964

<sup>[4] &</sup>quot;Adaptive signal processing", Bernard and Samuel 1985

### **Further**

- The kernel perceptron algorithm was already introduced by [1]
- Margin bounds for the Perceptron algorithm in the general non-separable case were proven by [2] and by [3] who extended existing results and gave new L1 bounds

<sup>[1] &</sup>quot;Theoretical foundations of the potential function method in pattern recognition learning", Aizerman 1964

<sup>[2] &</sup>quot;Large margin classification using the perceptron algorithm", Freund and Schapire 1999

<sup>[3] &</sup>quot;Perceptron mistake bounds", Mohri and Rostamizadeh 2013

## **Al Winter**

1969: Minsky and Papert exhibit the fact that it was difficult for perceptron to detect parity (number of activated pixels) and connectedness (are the pixels connected?). Besides, it was known that they cannot represent simple non linear function such as XOR function, see [1]

There is no reason to suppose that any of [the virtue of perceptrons] carry over to the many-layered version. Nevertheless, we consider it to be an important research problem to elucidate (or reject) our intuitive judgement that the extension is sterile. Perhaps some powerful convergence theorem will be discovered, or some profound reason for the failure to produce an interesting "learning theorem" for the multilayered machine will be found.

Minsky, Papert

- The starting point of "AI winter": a significant decline in funding of neural network research
- In [2]: more about the controversy between Rosenblatt and Minsky, Papert

<sup>[1] &</sup>quot;Perceptrons: an introduction to computational geometry", Minsky and Papert 2017

<sup>[2] &</sup>quot;A sociological study of the official history of the perceptrons controversy", Olazaran 1996.

### **Al Winter**

Perceptron model is defined as

$$f(x) = \sigma \Big(\sum_i w_i x_i + b \geq 0 \Big)$$

where  $\sigma(z)=\mathbf{1}(z>0)$  is the heavyside activation. Assuming boolean inputs  $x_i\in\{0,1\}$  it can implement

- $or(a,b) = \mathbf{1}(a+b-0.5 \ge 0)$
- $\operatorname{and}(a,b) = \mathbf{1}(a+b-1.5 \ge 0)$
- $not(a) = \mathbf{1}(-a + 0.5 \ge 0)$

but not the XOR which is not linearly separable: xor(a, b) = 1 if a = b = 0 or if a = b = 1 and xor(a, b) = 0 otherwise.

Question: find a neural network with one hidden layer that implements the xor (it exists...)

### Main building blocks

The MLP, computational graphs, SGD and backpropagation

## **Computational graphs**



The computation of

$$f(x) = \sigma \Big(\sum_i w_i x_i + b \Big)$$

can be represented as a computational graph where

- white nodes correspond to inputs and outputs;
- red nodes correspond to model parameters;
- blue nodes correspond to intermediate operations.

# **Computational graphs**

In terms of tensor operations, f can be rewritten as

 $f(x) = \sigma(x^ op w + b),$ 

for which the corresponding computational graph of f is:



## From perceptron to logistic unit

#### The sigmoid function looks a lot like the heavyside function



• But sigmoid has a non-flat gradient (and comes from sound statistical arguments).

Important. The  $x\mapsto \sigma(x^ op w+b)$  unit is the primitive of all neural networks!



The empirical risk minimization principle leads to a goodness-of-fit of the form

$$F( heta) = rac{1}{n}\sum_{i=1}^n \ell(y_i,f(x_i; heta))$$

where f is the considered model (NN architecture) and  $\theta$  encompasses all the trainable weights. Its gradient writes

$$abla F( heta) = rac{1}{n}\sum_{i=1}^n 
abla \ell(y_i, f(x_i; heta))$$

- This is called a full gradient or batch gradient: it uses all the training data (complexity grows linearly with the size *n* of the dataset)
- This is baaaaad if *n* is large !
- Idea: this empirical risk is an approximation of the expected risk, no need to minimize it with great accuracy

Stochastic gradient descent uses instead stochastic gradients

$$rac{1}{|B|}\sum_{i\in B}
abla \ell(y_i,f(x_i; heta))$$

where  $B \subset \{1, \ldots, n\}$  is sampled uniformly at random (at each step) with a minibatch size  $m = |B| \ll n$ . One step of vanilla SGD is

$$heta_{t+1} \leftarrow heta_t - \eta_t rac{1}{|B_t|} \sum_{i \in B_t} 
abla \ell(y_i, f(x_i; heta_t))$$

where  $\eta_t$  is the learning rate and  $B_t$  is the mini-batch used in the t-th iteration

- Iteration complexity is independent of *n*.
- The stochastic process  $\{ heta_t: t=1,2,\ldots\}$  depends on the mini-batches sampled at each iteration
- In practice, shuffle  $\{1,\ldots,n\}$  and use sequentially  $B_t=\{1,\ldots,m\}$ ,  $B_{t+1}=\{m+1,\ldots,2m\}$ , etc.

$$heta_{t+1} \leftarrow heta_t - \eta_t rac{1}{|B_t|} \sum_{i \in B_t} 
abla \ell(y_i, f(x_i; heta_t))$$

- Iteration complexity is independent of *n*.
- The stochastic process  $\{ heta_t: t=1,2,\ldots\}$  depends on the mini-batches sampled at each iteration
- In practice, shuffle  $\{1,\ldots,n\}$  and use sequentially  $B_t=\{1,\ldots,m\}$ ,  $B_{t+1}=\{m+1,\ldots,2m\}$ , etc.



Batch gradient descent

Stochastic gradient descent

Why is stochastic gradient descent a good idea?

- Stochastic gradients are unbiased approximations of the batch gradient (uniform sampling of the mini-batch)
- Many guarantees on SGD and more advanced variants in the convex case
- If training is limited to single pass over the data, then SGD directly minimizes the expected risk
- Bottou and Bousquet [1]: stochastic optimization algorithms (e.g., SGD) yield the best generalization performance (in terms of excess error) despite being the worst optimization algorithms for minimizing the empirical risk

## **Multi-layer perceptron**

- So far we considered the logistic unit  $h = \sigma(w^{ op}x + b)$ , where  $h \in \mathbb{R}$ ,  $\mathbf{x} \in \mathbb{R}^d$ ,  $w \in \mathbb{R}^d$  and  $b \in \mathbb{R}$ .
- These units can be composed in parallel to form a layer with q outputs:

$$h = \sigma(\mathbf{W}^{ op}x + b)$$

where  $h \in \mathbb{R}^q$ ,  $x \in \mathbb{R}^d$ ,  $\mathbf{W} \in \mathbb{R}^{d \times q}$ ,  $b \in \mathbb{R}^q$  and where  $\sigma(\cdot)$  is applied entrywise.

• We also say that the width or number of units/neurons of the layer is q. It's the output size of the layer.



## **Multi-layer perceptron**

Similarly, layers can be composed in series, such that:

$$egin{aligned} h_0 &= x \ h_1 &= \sigma(\mathbf{W}_1^ op h_0 + b_1) \ &dots \ h_L &= \sigma(\mathbf{W}_L^ op h_{L-1} + b_L) \ f(\mathbf{x}; heta) &= \hat{y} &= h_L \end{aligned}$$

where heta denotes the model parameters  $\{\mathbf{W}_k, b_k : k=1,\ldots,L\}$ .

- This model is the multi-layer perceptron
- Also known as the fully connected feedforward network.



## **Output layer**

- For binary classification, the width of the last layer L is set to q=1, which results in a single output  $h_L\in[0,1]$  that models the probability  $\mathbb{P}[Y=1|X=x].$
- For C-class classification with C>2, the activation function  $\sigma$  in the last layer is generalized to produce a vector  $h_L \in \Delta$  where  $\Delta$  is the probability simplex in  $\mathbb{R}^C$  for the probability estimates  $\mathbb{P}[Y = c | X = x]$  for  $c = 1, \ldots, C$ .

This activation is the softmax function, where its i-th output is defined as

$$ext{softmax}(z)_i = rac{\exp(z_i)}{\sum_{c=1}^C \exp(z_c)},$$

for c = 1, ..., C.

## We need to compute gradients !

• To minimize  $F(\theta)$  with SGD, given a mini-batch B, we need to compute (stochastic) gradients

$$abla F^B( heta) = rac{1}{|B|} \sum_{i \in B} 
abla \ell(y_i, f(x_i; heta))$$

with respect to heta

• Therefore, we require the evaluation of the (total) derivatives

$$rac{\mathrm{d} F^B}{\mathrm{d} \mathbf{W}_k} \quad ext{and} \quad rac{\mathrm{d} F^B}{\mathrm{d} \mathbf{b}_k}$$

of  $F^B$  with respect to all model parameters  $\mathbf{W}_k$ ,  $\mathbf{b}_k$ , for  $k=1,\ldots,L$ 

• These derivatives can be evaluated automatically from the computational graph of F using automatic differentiation

### **Chain rule**



Let us consider a 1-dimensional output composition  $f\circ g$  , such that

$$egin{aligned} y &= f(\mathbf{u}) \ \mathbf{u} &= g(x) = (g_1(x),...,g_m(x)). \end{aligned}$$

## **Chain rule**

The chain rule states that  $(f \circ g)' = (f' \circ g)g'$ .

For the total derivative, the chain rule generalizes to

$$rac{\mathrm{d}y}{\mathrm{d}x} = \sum_{k=1}^m rac{\partial y}{\partial u_k} \underbrace{rac{\mathrm{d}u_k}{\mathrm{d}x}}_{ ext{recursive case}}$$

#### **Reverse automatic differentiation**

- Since a neural network is a composition of differentiable functions, the total derivatives of the loss can be evaluated backward, by applying the chain rule recursively over its computational graph.
- The implementation of this procedure is called reverse automatic differentiation or backpropagation

Let us consider a simplified 2-layer MLP and the following loss function:

$$F(\mathbf{W}_1,\mathbf{W}_2) = \ell(y,\hat{y}) + \lambda(\|\mathbf{W}_1\|_2^2 + \|\mathbf{W}_2\|_2^2)$$

where for  $x\in\mathbb{R}^d$  ,  $y\in\mathbb{R}$  ,  $\mathbf{W}_1\in\mathbb{R}^{d imes q}$  and  $\mathbf{W}_2\in\mathbb{R}^q$  :

$$\hat{y} = f(x; \mathbf{W}_1, \mathbf{W}_2) = \sigma(\mathbf{W}_2^ op \sigma(\mathbf{W}_1^ op x))$$

In the forward pass, many intermediate values are computed and kept in memory from inputs to outputs, which results in the annotated computational graph below:



The total derivative can be computed through a backward pass, by walking through all paths from outputs to parameters in the computational graph and accumulating the terms. For example, for  $\frac{d\ell}{d\mathbf{W}_1}$  we have:





Let us zoom in on the computation of  $\hat{y}$  and its derivative with respect to  $\mathbf{W}_1$ :



- Forward pass: values  $u_1, u_2, u_3$  and  $\hat{y}$  are computed by traversing the graph from inputs to outputs given  $x, \mathbf{W}_1$  and  $\mathbf{W}_2$ .
- Backward pass: by the chain rule we have

$$egin{aligned} rac{\mathrm{d}\hat{y}}{\mathrm{d}\mathbf{W}_1} &= rac{\partial\hat{y}}{\partial u_3}rac{\partial u_3}{\partial u_2}rac{\partial u_2}{\partial u_1}rac{\partial u_1}{\partial \mathbf{W}_1} \ &= rac{\partial\sigma(u_3)}{\partial u_3}rac{\partial\mathbf{W}_2^Tu_2}{\partial u_2}rac{\partial\sigma(u_1)}{\partial u_1}rac{\partial\mathbf{W}_1^T\mathbf{x}}{\partial \mathbf{W}_1} \end{aligned}$$

• Note how evaluating partial derivatives requires the intermediate values computed during the forward pass

- This algorithm is called automatic differentiation or backpropagation
- Often mixed up with the optimization algorithm: backpropagation does not optimize, it just helps to compute gradients
- It's the core of any deep learning library: a grammar of differentiable blocks that can be combined together
- An equivalent procedure can be defined to evaluate the derivatives in forward mode, from inputs to outputs
- Since differentiation is a linear operator, automatic differentiation can be implemented efficiently in terms of tensor operations

# **Vanishing gradients**

Training deep MLPs with many layers has for long (pre-2011) been very difficult due to the vanishing gradient problem.

- Small gradients slow down, and eventually kills "learning" with stochastic gradient descent.
- Mitigated the use of deep architectures (many layer stacked)



Backpropagated gradients normalized histograms (Glorot and Bengio, 2010). Gradients for layers far from the output vanish to zero.

Let us consider a simplified 3-layer MLP, with  $x, w_1, w_2, w_3 \in \mathbb{R}$ , such that

$$f(x;w_1,w_2,w_3)=\sigma\left(w_3\sigma\left(w_2\sigma\left(w_1x
ight)
ight)
ight).$$

Under the hood, this would be evaluated as

$$egin{aligned} & u_1 = w_1 x \ & u_2 = \sigma(u_1) \ & u_3 = w_2 u_2 \ & u_4 = \sigma(u_3) \ & u_5 = w_3 u_4 \ & \hat{y} = \sigma(u_5) \end{aligned}$$

and its derivative  $\frac{d\hat{y}}{dw_1}$  as  $\frac{d\hat{y}}{dw_1} = \frac{\partial\hat{y}}{\partial u_5} \frac{\partial u_5}{\partial u_4} \frac{\partial u_4}{\partial u_3} \frac{\partial u_3}{\partial u_2} \frac{\partial u_2}{\partial u_1} \frac{\partial u_1}{\partial w_1}$   $= \frac{\partial\sigma(u_5)}{\partial u_5} w_3 \frac{\partial\sigma(u_3)}{\partial u_3} w_2 \frac{\partial\sigma(u_1)}{\partial u_1} x$  The derivative of the sigmoid activation function  $\sigma$  is:



$$rac{\mathrm{d}\sigma}{\mathrm{d}x}(x)=\sigma(x)(1-\sigma(x))$$

Notice that  $0 \leq rac{\mathrm{d}\sigma}{\mathrm{d}x}(x) \leq rac{1}{4}$  for all x.

Assume that weights  $w_1, w_2, w_3$  are initialized randomly from a Gaussian with zero-mean and small variance, such that with high probability  $-1 \le w_i \le 1$ .

Then,



This implies that the gradient  $\frac{d\hat{y}}{dw_1}$  exponentially shrinks to zero as the number of layers in the network increases.

Hence the vanishing gradient problem.

- In general, bounded activation functions (sigmoid, tanh, etc) are prone to the vanishing gradient problem.
- Note the importance of a proper initialization scheme.

## **Rectified linear units**

Instead of the sigmoid activation function, modern neural networks are for most based on rectified linear units (ReLU) (Glorot et al, 2011):

 $\operatorname{ReLU}(x) = \max(0, x)$ 



Note that the derivative of the ReLU function is

$$rac{\mathrm{d}}{\mathrm{d}x}\mathrm{ReLU}(x) = egin{cases} 0 & ext{if } x \leq 0 \ 1 & ext{otherwise} \end{cases}$$



For x = 0, the derivative is undefined. In practice, it is set to zero.

#### Therefore,

$$rac{\mathrm{d}\hat{y}}{\mathrm{d}w_1} = \underbrace{rac{\partial\sigma(u_5)}{\partial u_5}}_{=1} w_3 \underbrace{rac{\partial\sigma(u_3)}{\partial u_3}}_{=1} w_2 \underbrace{rac{\partial\sigma(u_1)}{\partial u_1}}_{=1} x$$

This solves the vanishing gradient problem, even for deep networks! (provided proper initialization)

#### Note that:

- The ReLU unit dies when its input is negative, which might block gradient descent.
- This is actually a useful property to induce sparsity.
- This issue can also be solved using leaky ReLUs, defined as

 $\mathrm{LeakyReLU}(x) = \max(lpha x, x)$ 

for a small  $lpha \in \mathbb{R}^+$  (e.g., lpha = 0.1).

# Modern deep learning is like LEGO®

### **LEGO®** Deep Learning





People are now building a new kind of software by **assembling networks of parameterized functional blocks** and by **training them from examples using some form of gradient-based optimization**.

Yann LeCun, 2018.

### **DL** as an architectural language



### **DL** as an architectural language



The toolbox

# **LEGO®** Creator Expert



# Thank you !