

# MACHINE LEARNING FOR OPTIMIZATION

January 15, 2025

- ↳ Crash course on neural nets and automatic differentiation
- ↳ Introduction to implicit layers

# ① A crash course on neural networks

Reference: M. Hertz and B. Recht, Patterns, predictions and actions

↳ A neural network (for us) : a sequence of layers  $z^{(0)}, \dots, z^{(L)}$

Layer  $l$  : 
$$\underbrace{z^{(l+1)}}_{\text{vector}} = \underbrace{\phi}_{\text{(Typically) Nonlinear activation function}} \left( \underbrace{A^{(l)} z^{(l)} + b^{(l)}}_{\substack{\text{Linear transformation} \\ \text{of } z^{(l)}}} \right)$$

$A^{(l)}, b^{(l)}$  : parameters of the  $l^{\text{th}}$  layer (to be learned)

Ex)  $\phi : y = [y_i]_i \mapsto [\max(y_i, 0)]_i$  ReLU activation

$\phi : y = [y_i]_i \mapsto [\tanh(y_i)]_i$  Hyperbolic tangent

## Classical layer structures

\* Fully connected layer

$$\underbrace{x}_{\text{vector}} \mapsto \underbrace{z}_{\text{vector}}$$

$$\forall i, \underbrace{z_i}_{\in \mathbb{R}} = \phi \left( \sum_j A_{ij} x_j + b_i \right)$$

\* Convolutional layers

$$\underbrace{x}_{\mathbb{R}^{d_0 \times d_0 \times c_0}} \mapsto \underbrace{z}_{\mathbb{R}^{d_1 \times d_1 \times c_1}}$$

$$A \in \mathbb{R}^{d_1 \times d_1 \times c_0 \times c_1} \quad b \in \mathbb{R}^{c_1}$$



$z$  - order kernel

$$\forall (i, j, k),$$

$$z_{ijk} = \phi \left( \sum_{l, m, n} A_{l, m, n, k} x_{i-l, j-m, n} + b_k \right)$$

\* Recurrent layer

$$x \mapsto z$$

$$z = x + B(\phi(Ax))$$

$A, B$ : Network parameters

$$z^{(0)} \mapsto z^{(1)} \mapsto \dots \mapsto z^{(L)}$$

$$\forall \ell, z^{(\ell+1)} = z^{(\ell)} + B(\phi(Az^{(\ell)}))$$

Training: Optimization problem

\* Data:  $\{(x_i, y_i)\}_{i=1 \dots m}$  (ex: supervised learning, want to predict  $y_i$  given  $x_i$ )

\* Neural network:  $x \mapsto NN(x; w)$

all the parameters of the network stacked (conceptually) into a big vector in  $\mathbb{R}^d$

\* Problem (Empirical Risk Minimization)

minimize  $w \in \mathbb{R}^d$

$$\frac{1}{m} \sum_{i=1}^m \text{loss}(NN(x_i; w), y_i)$$

Desired output  $y_i$

Loss function (scalar valued)

output of neural net with input  $x_i$

$F(w)$

Typical loss:  $\frac{1}{2} \|NN(x_i, w) - y_i\|^2$  (squared Euclidean norm)

$$\forall v \in \mathbb{R}^k, \|v\|^2 = \sum_{i=1}^k v_i^2$$

↳ In practice, this problem is solved using stochastic gradient techniques (SGD, Adam, ...)

⇒ Require  $\nabla_w F(w) = \frac{1}{n} \sum_{i=1}^n \nabla_w (\text{loss}(NN(x_i; w), y_i))$

⇒ Need an efficient way of computing derivatives, especially

$D_w NN(x_i; w)$

$\in \mathbb{R}^{d \times ?}$

→ Derivative of the NN mapping with respect to  $w$

- Want a technique that is more efficient than coding the derivative by hand
- Want to be able to compute  $D_w NN(x_i; w)$  for multiple  $x_i$ s ("batch") in parallel

Automatic / Implicit differentiation

• General, abstract idea

$$z^{(0)} \mapsto z^{(1)} \mapsto \dots \mapsto z^{(L)}$$

$$z^{(1)} = f_1(z^{(0)}), \quad z^{(2)} = f_2(z^{(1)}), \quad \dots$$

$$z^{(l)} \in \mathbb{R}^{d_l} \quad , \quad f_l: \mathbb{R}^{d_{l-1}} \mapsto \mathbb{R}^{d_l}$$

Jacobian  $D_{z^{(l-1)}} f_l(z) \in \mathbb{R}^{d_{l-1} \times d_l}$  Matrix of partial derivatives

Goal: Compute  $D_{z^{(0)}} z^{(L)}$ .

## Backpropagation

① Forward pass: Evaluate  $z^{(L)}$  and build a computational graph

Input:  $x \in \mathbb{R}^{d_0}$

Set  $v_0 = x$

For  $l = 1 \dots L$

Store  $v_{l-1}$  and compute  $v_l = f_l(v_{l-1})$

Output  $v_L$

Numerical value vs  $z^{(l)}$  function

② Backward pass: Go through the computational graph backwards

Input:  $(v_0, \dots, v_L)$

Set  $\Delta_L = D_{z^{(L)}} z^{(L)}(v_L) = I_{d_L}$

For  $l = L, \dots, 1$

$\Delta_{l-1} = \Delta_l \cdot D_{z^{(l-1)}} z^{(l)}(v_{l-1})$

Output  $\Delta_0$

↓  
Computed at layer  $l$   
doesn't require the knowledge  
of the other layers

Can prove that  $\Delta_0 = D_{z^{(0)}} z^{(L)}(x)$  and  $\Delta_l = D_{z^{(0)}} z^{(l)}(x)$

Intuition (Frequently written this way in ML papers)

$$\frac{\partial z^{(l-1)}}{\partial z^{(l)}} \times \frac{\partial z^{(l)}}{\partial z^{(l-1)}} = \frac{\partial z^{(l)}}{\partial z^{(l)}}$$

"Chain rule"

$$D_{z^{(l)}} z^{(l-1)}(x) = D_{z^{(l)}} z^{(l)}(x) D_{z^{(l-1)}} z^{(l)}(x)$$

The technique above can be used to compute  $D_x NN(x; w)$  for any neural network with standard layers and activation functions (even ReLU, which doesn't have a derivative everywhere)

→ The same approach (and the same computational graph structure) can be used to compute  $D_w NN(x; w)$

$$Ex) \quad NN(x; w) = \overbrace{A_L}^{d_L \times d_{L-1}} \text{ReLU}(\overbrace{A_{L-1}}^{d_{L-1} \times d_{L-2}} \text{ReLU}(\dots (\overbrace{A_1}^{d_1 \times d_0} x) \dots))$$

$w$  concatenates all  $A_l$  matrices  
 $w \in \mathbb{R}^d, d = d_L d_{L-1} + d_{L-1} d_{L-2} + \dots + d_1 d_0$

$$x \mapsto \text{ReLU}(A_1 x) \mapsto \text{ReLU}(A_2 z^{(1)}) \mapsto \dots$$

$$D_x NN(x; w) = A_L \Delta_{L-1} A_{L-1} \Delta_{L-2} \dots \Delta_1 A_1$$

$$\text{where } \Delta_l = \text{diag}(z^{(l)})$$

For any  $l$ , can compute the derivative of  $NN(x; w)$

with respect to any column of  $A_l$  by using

$$u = \text{ReLU}(A_{l-1}(\dots A_1 x \dots)) \quad \left. \vphantom{u} \right] \text{input of layer } l$$

$$D_{A_l u} f_l(A_l u) = A_L \Delta_{L-1} \dots \Delta_l$$

$$f_l(z) = A_L \text{ReLU}(\dots A_{L+1} \text{ReLU}(z) \dots)$$

\*  $j$ ,

$$D_{[A_l]_{:j}} \text{NN}(x, w) = D_{A_l u} f_l(A_l u) \times D_{[A_l]_{:j}} A_l u$$

$\leftarrow j^{\text{th}} \text{ column of } A_l \quad \rightarrow$

$\underbrace{\hspace{10em}}_{= u_j I}$

$$A_l \in \mathbb{R}^{d_l \times d_{l-1}}$$

↳ Python libraries for automatic differentiation  
 Autograd (within PyTorch)  
 JAX

Takeaway:

• No need to compute derivatives by yourself, AD can do it...

• ... especially for neural networks with standard, explicitly defined layers

↳ Fully connected, Recurrent / Convolutional all look like  $x \mapsto \boxed{\text{compute } z = f(x)} \rightarrow z$

Q) What if the layer is now defined implicitly,  
 e.g.  $x \mapsto \underset{z}{\operatorname{argmin}} f(z; x) \text{ s.t. } z \in C(x)$

## ② Implicit layers

Explicit layer:  $x \mapsto \boxed{\text{Compute } f(x)} \rightarrow z = f(x)$

Implicit layer:  $x \mapsto \boxed{\text{Find } z \text{ such that } C(x, z) = 0} \rightarrow z$

Joint condition on the input and the output

→ The implicit layer paradigm separates what the layer should do from the implementation of the layer

$C(x, z) = 0$  can be

- An algebraic equation  
 ⇓  
 Deep Equilibrium Models

- A differential equation  
 Neural ODEs / PINNs

- Optimization problem / Optimality conditions of an optimization

→ Differentiable optimization