# La propagation avant en tant que matrice

Samuel Leblanc

Université de Sherbrooke

14 novembre 2024

# Références
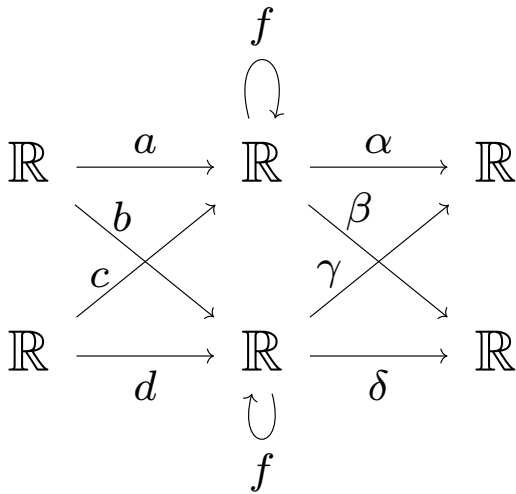
**i**. Marco ARMENTA et Pierre-Marc JODOIN. "The Representation Theory of Neural Networks". In : *Mathematics* 9.24 (2021)

**ii**. Marco ARMENTA, Thomas BRÜSTLE, Souheila HASSOUN et Markus REINEKE. "Double Framed Moduli Spaces of Quiver Representations". In : *Linear Algebra and its Applications* 650 (2022), p. 98-131

**iii**. Marco ARMENTA, Thierry JUDGE, Nathan PAINCHAUD, Youssef SKANDARANI, Carl LEMAIRE, Gabriel GIBEAU SANCHEZ, Philippe SPINO et Pierre-Marc JODOIN. "Neural Teleportation". In : *Mathematics* 11.2 (2023)

**iv**. L., Aiky RASOLOMANANA et Marco ARMENTA. *Hidden Activations Are Not Enough : A General Approach to Neural Network Predictions*. 2024. arXiv : 2409.13163 [cs.LG]
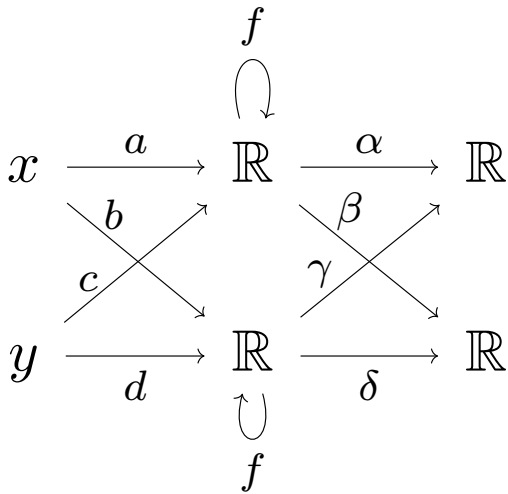
$$\mathcal{D} = \big\{(x_1, y_1), (x_2, y_2), \ldots, (x_n, y_n)\big\}$$

$$\mathcal{D} = \big\{(x_1, y_1), (x_2, y_2), \ldots, (x_n, y_n)\big\}$$

But : Trouver $\psi$ telle que $\psi(x_i) = y_i$.

$$\Psi(W, f) = \begin{pmatrix} \alpha & \gamma \\ \beta & \delta \end{pmatrix} \circ \begin{pmatrix} f \\ f \end{pmatrix} \circ \begin{pmatrix} a & c \\ b & d \end{pmatrix}$$
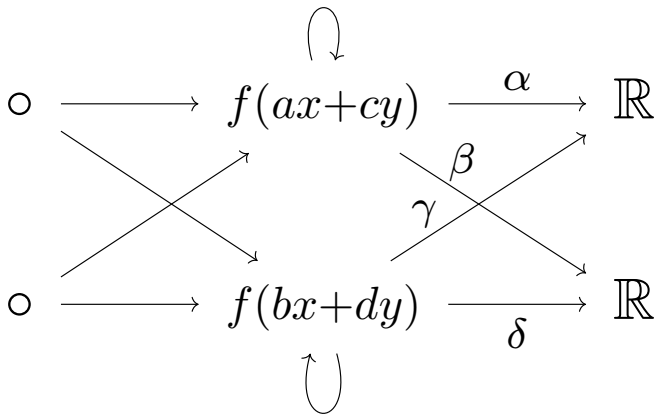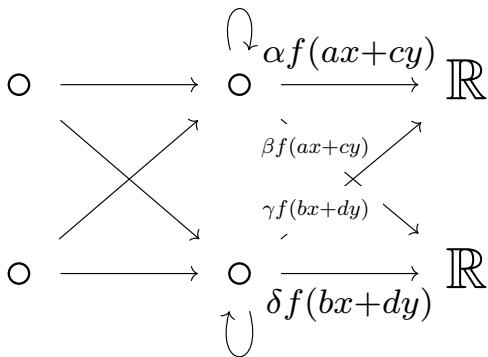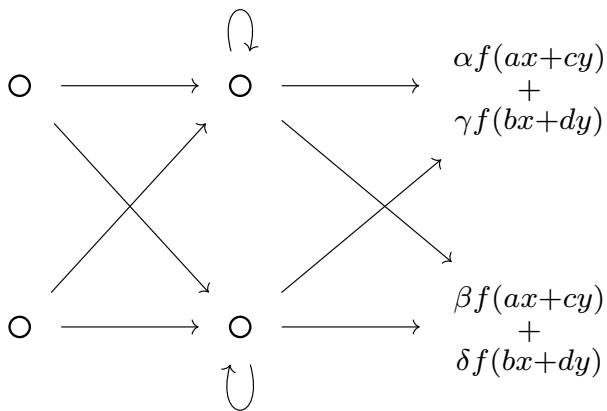
$$
\begin{array}{ccccc}
& & f & & \\
& & \circlearrowright & & \\
x & \xrightarrow{\ a\ } & \mathbb{R} & \xrightarrow{\ \alpha\ } & \mathbb{R} \\
& \begin{subarray}{c} b \\ c \end{subarray} \times & & \begin{subarray}{c} \beta \\ \gamma \end{subarray} \times & \\
y & \xrightarrow{\ d\ } & \mathbb{R} & \xrightarrow{\ \delta\ } & \mathbb{R} \\
& & \circlearrowright & & \\
& & f & &
\end{array}
$$

$$\alpha f(ax+cy)$$
$$+$$
$$\gamma f(bx+dy)$$

$$\beta f(ax+cy)$$
$$+$$
$$\delta f(bx+dy)$$

$$\mathbb{R}^2 \xrightarrow{\Psi(W,f)} \mathbb{R}^2$$

$$\mathbb{R}^2 \xrightarrow{\Psi(W,f)} \mathbb{R}^2$$

$$\varphi(W,f) \downarrow \qquad \nearrow \Psi(\cdot,1)(1)$$

$$\mathrm{Rep}Q/_{\underline{\cong}}$$

$$\begin{array}{ccc}
\mathbb{R}^2 & \xrightarrow{\ \Psi(W,f)\ } & \mathbb{R}^2 \\
{\scriptstyle \varphi(W,f)}\Big\downarrow & & \Big\uparrow{\scriptstyle \mathrm{ev}_1} \\
\mathrm{Rep}Q/_{\underset{\cong}{\cdot}} & \xrightarrow[\ \pi\ ]{} & \mathrm{Mat}_{2\times 2}(\mathbb{R})
\end{array}$$

$$\begin{array}{ccc}
\mathbb{R}^2 & \xrightarrow{\;\;\Psi(W,f)\;\;} & \mathbb{R}^2 \\
{\scriptstyle\varphi(W,f)}\downarrow & \searrow{\scriptstyle\mathbb{M}(W,f)} & \uparrow{\scriptstyle\mathrm{ev}_1} \\
\mathrm{Rep}Q/_{\underset{\cong}{\cdot}} & \xrightarrow[\pi]{} & \mathrm{Mat}_{2\times 2}(\mathbb{R})
\end{array}$$

# Attaques adversariales



(a) Image     (b) Adversarial examples with target class labels

Figure 1: Visual illustration of adversarial examples crafted by EAD (Algorithm 1). The original example is an ostrich image selected from the ImageNet dataset (Figure 1 (a)). The adversarial examples in Figure 1 (b) are classified as the target class labels by the Inception-v3 model.

Pin-Yu CHEN, Yash SHARMA, Huan ZHANG, Jinfeng YI et Cho-Jui HSIEH. "EAD : Elastic-Net Attacks to Deep Neural Networks via Adversarial Examples". In : AAAI Press, 2018. URL : https://arxiv.org/abs/1709.04114
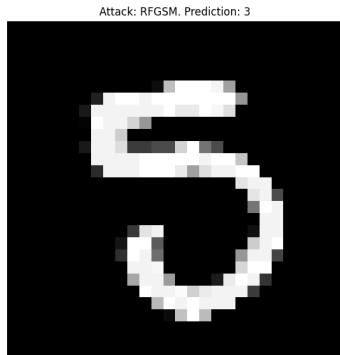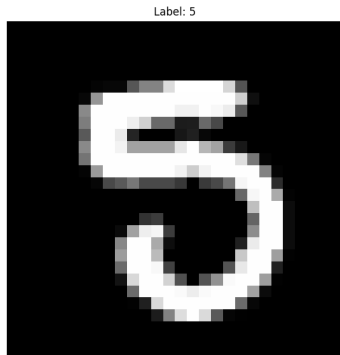
# Attaques adversariales



Figure 8. Real-life example of a backdoored stop sign near the authors'
office. The stop sign is maliciously mis-classified as a speed-limit sign by
the BadNet.

Tianyu Gu, Brendan Dolan-Gavitt et Siddharth Garg. *BadNets : Identifying Vulnerabilities in the Machine Learning Model Supply Chain*. 2019. arXiv : 1708.06733 [cs.CR]. URL : https://arxiv.org/abs/1708.06733

# MNIST — Exemple d'une détection d'attaque

Le réseau $(W, f)$ a été entrainé jusqu'à $\approx 98\%$ de précision.



Si on note par $\mathbb{M}_3$ la matrice moyenne de la classe $3$, $\Psi_3$ le vecteur de sortie moyen de la classe $3$ et $x$ l'image attaquée, alors

$$\|\mathbb{M}(W, f)(x) - \mathbb{M}_3\|_\infty \approx 84 \text{ et } \|\Psi(W, f)(x) - \Psi_3\|_\infty \approx 4.$$

CODE : *simple_adversarial_detection*. https://github.com/samueleblanc/simple_adversarial_detection. 2024