# Input relevance analysis driven by spectral learning
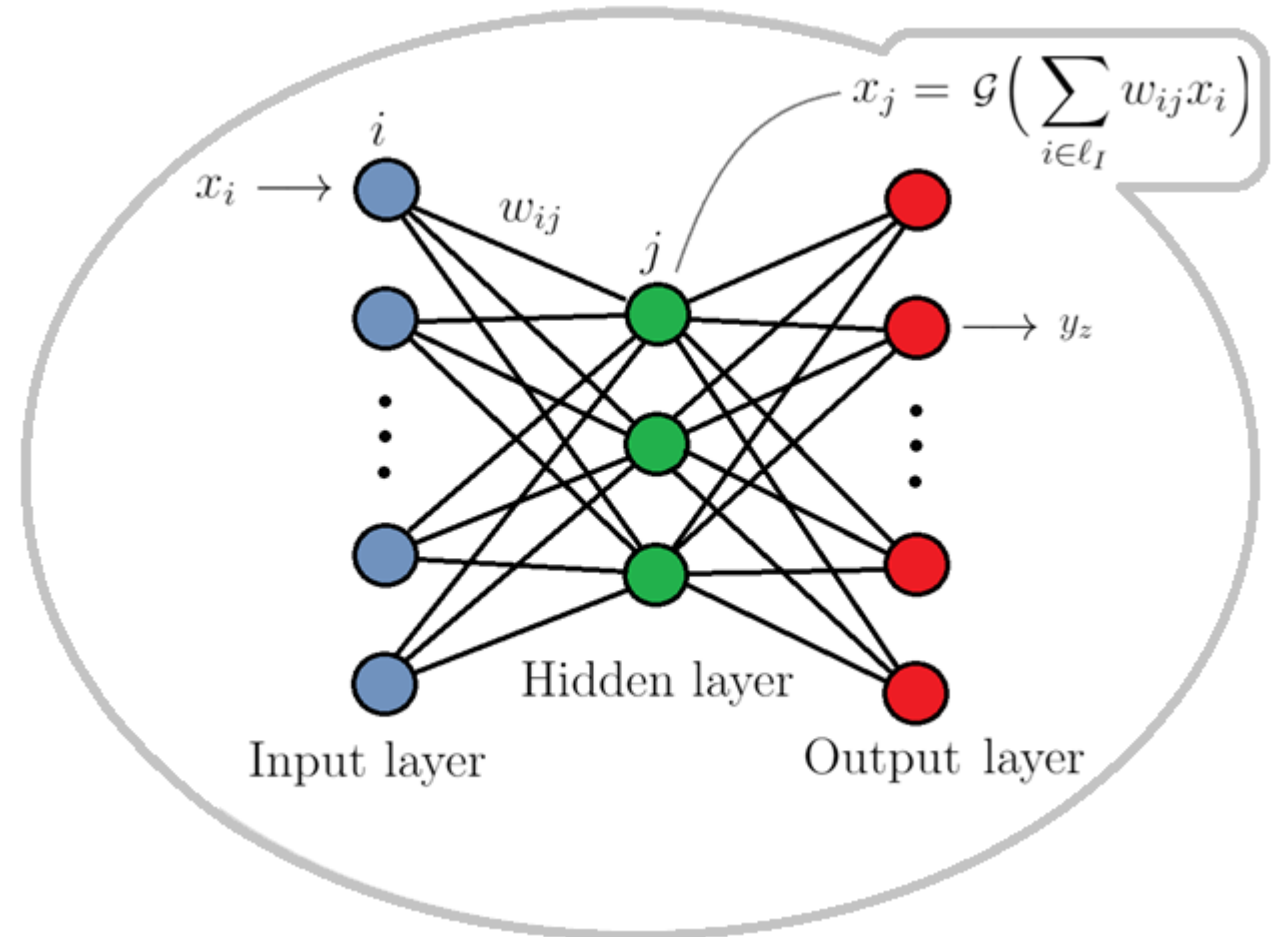
Lorenzo Chicchi

# Neural Networks

Neural networks (NNs) are models that can efficiently **approximate complex target functions** (Universal approximation theorem[1,2])

$$\vec{y} = N(\vec{x}) \approx \mathcal{F}(\vec{x})$$



$$x_j = \mathcal{G}\left(\sum_{i \in \ell_I} w_{ij} x_i\right)$$

- $\boldsymbol{G_l}$ is a non-linear function (Tanh, Sigmoid,..)
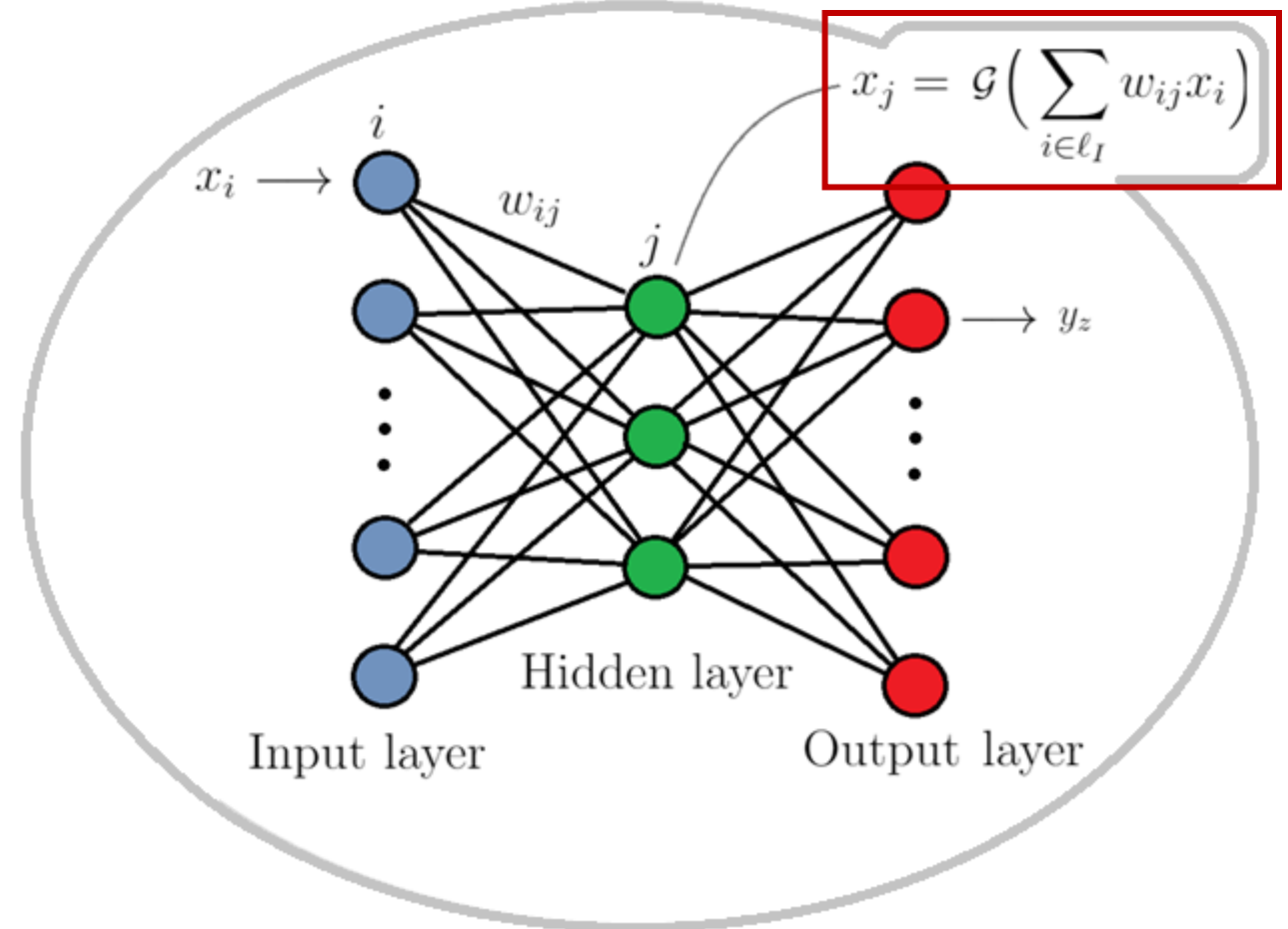- $w_{ij}$ are the network (trainable) parameters

# Neural Networks

Neural networks (NNs) are models that can efficiently **approximate complex target functions** (Universal approximation theorem[1,2])

$$\vec{y} = N(\vec{x}) \approx \mathcal{F}(\vec{x})$$

- $G_l$ is a non-linear function (Tanh, Sigmoid,..)
- $w_{ij}$ are the network (trainable) parameters

$$x_j = \mathcal{G}\left(\sum_{i \in \ell_I} w_{ij} x_i\right)$$

$x_i \longrightarrow$

$w_{ij}$

$j$

$\longrightarrow y_z$

Input layer

Hidden layer

Output layer

# Neural Networks: how to use them

Starting from a DATA SET, that is, a set of examples input-target:    $(\vec{x}, \widehat{y})_i \; i = 1, \dots N$

| $\vec{x}$ | $\widehat{y}$ |
|---|---|
|  | ⟶  5 |
|  | ⟶  1 |
|  | ⟶  1 |
|  | ⟶  9 |
|  | ⟶  2 |

# Neural Networks: how to use them

Starting from a DATA SET, that is, a set of examples input-target:     $(\vec{x}, \hat{y})_i \quad i = 1, \ldots N$

# Neural Networks: how to use them

Starting from a DATA SET, that is, a set of examples input-target: $(\vec{x}, \hat{y})_i \ i = 1, \dots N$



$\vec{x}_i \longrightarrow$
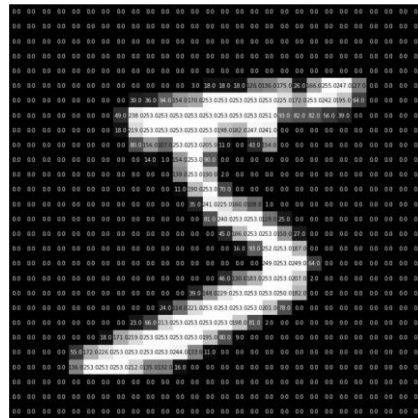
$\longrightarrow N(\vec{x}_i; \vec{w})$

# Neural Networks: how to use them
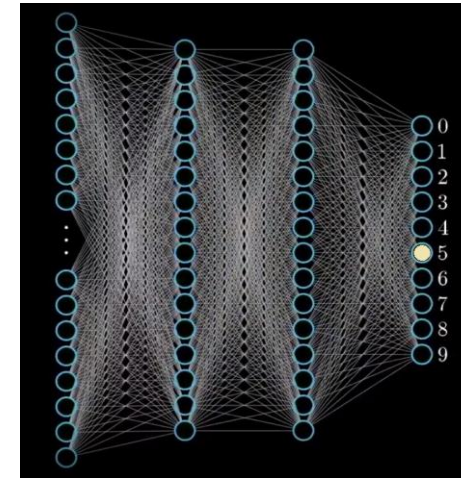
Starting from a DATA SET, that is, a set of examples input-target: $(\vec{x}, \hat{y})_i \; i = 1, \ldots N$



$\vec{x}_i \longrightarrow$

$\longrightarrow N(\vec{x}_i; \vec{w})$

**Train the network**: find the parameters $w_{ij}$ that minimize a *Loss function*:

$$\mathcal{L} \propto \sum_i |\hat{y}_i - y_i| = \sum_i |\hat{y}_i - N(\vec{x}_i; \vec{w})|$$

# Back to Neural Networks

**Feed forward architecture:**

The process consists in applying linear transformations $W_l$ alternating with non-linear functions $g_l$



$$\vec{x}_i \rightarrow \boxed{W_1(\vec{x}_i) \rightarrow g_1\big(W_1(\vec{x}_i)\big)} \rightarrow \boxed{W_2\big(g_1(W_1(\vec{x}_i))\big) \rightarrow g_2\Big(W_2\big(g_1(W_1(\vec{x}_i))\big)\Big)} \rightarrow \cdots$$

Layer 1          Layer 2

# Deep Learning in spectral domain

Giambagli et al. [1] recently introduced the idea of training neural networks by acting in the **reciprocal space**.

The main idea is to tune the spectral parameters instead of the weights of the matrix W.
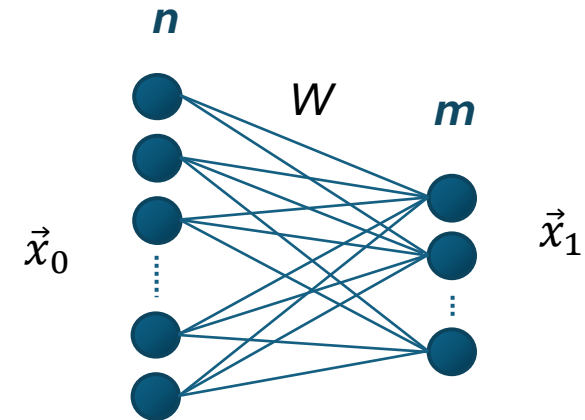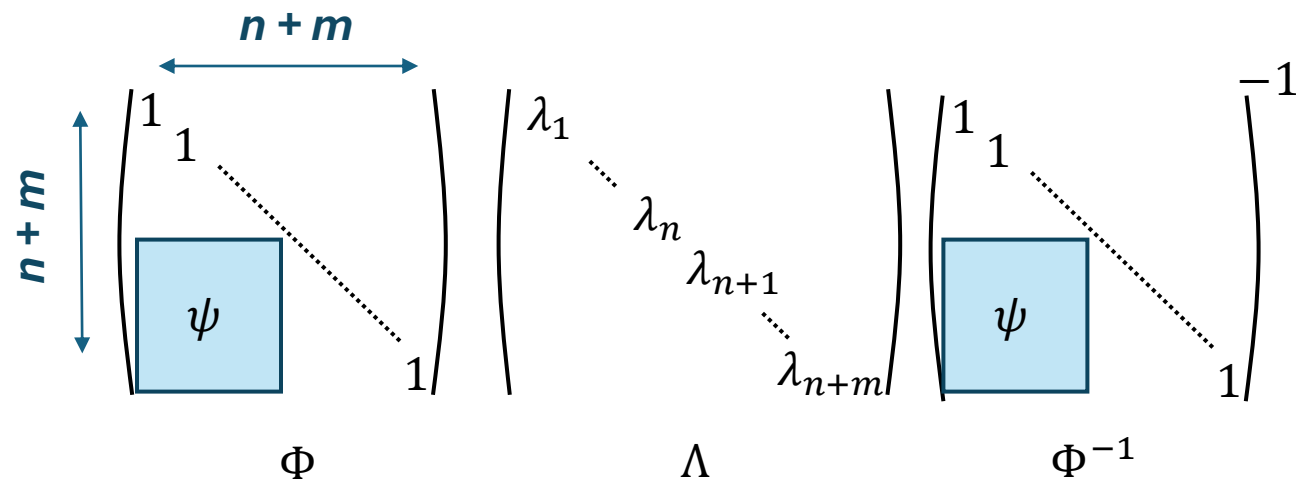
$$W = \Phi \Lambda \Phi^{-1}$$

↓

Spectral decomposition



[1] L. Giambagli, L. Buffoni, T. Carletti, W. Nocentini, and D.Fanelli, Machine learning in spectral domain, Nat. Commun. 12, 1330 (2021).
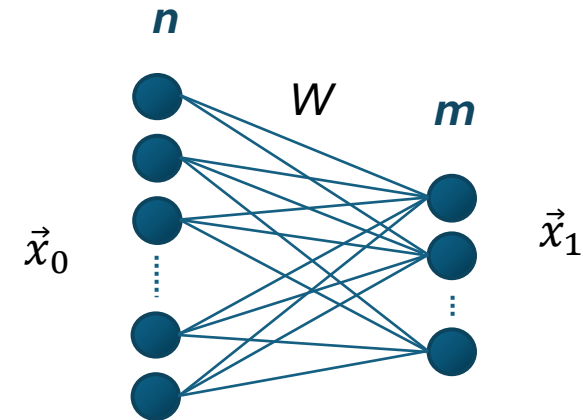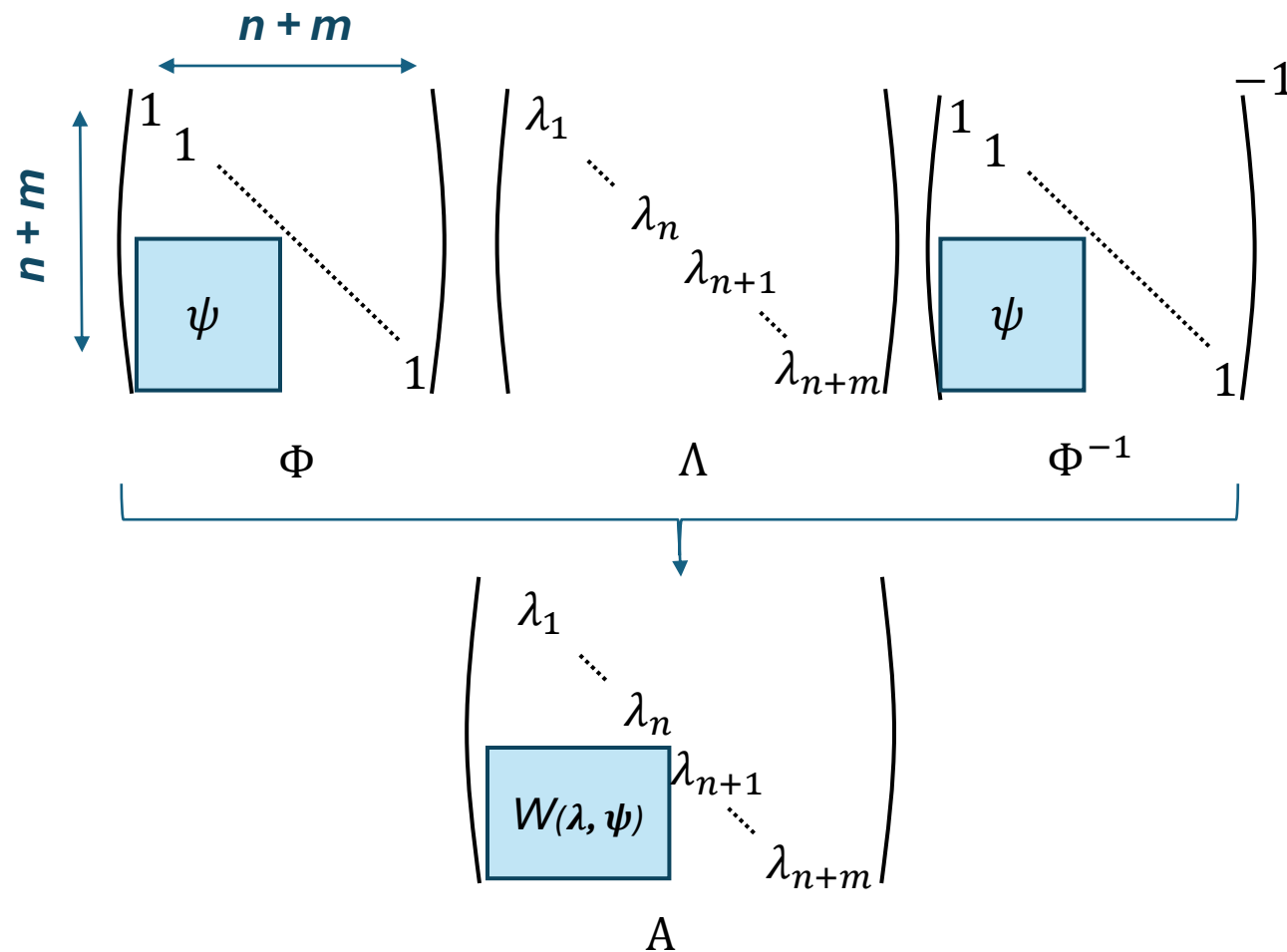
# Deep Learning in spectral domain

Starting from the reciprocal space, a suitable set of eigenvectors $\Phi$ can be chosen so as to obtain in the direct space a matrix A that describes the linear transfer between two consecutive layers.

# Deep Learning in spectral domain

Starting from the reciprocal space, a suitable set of eigenvectors $\Phi$ can be chosen so as to obtain in the direct space a matrix A that describes the linear transfer between two consecutive layers.
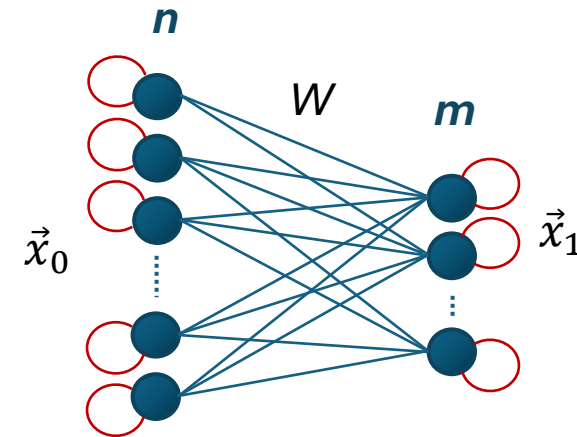
# Deep Learning in spectral domain

Starting from the reciprocal space, a suitable set of eigenvectors $\Phi$ can be chosen so as to obtain in the direct space a matrix A that describes the linear transfer between two consecutive layers.
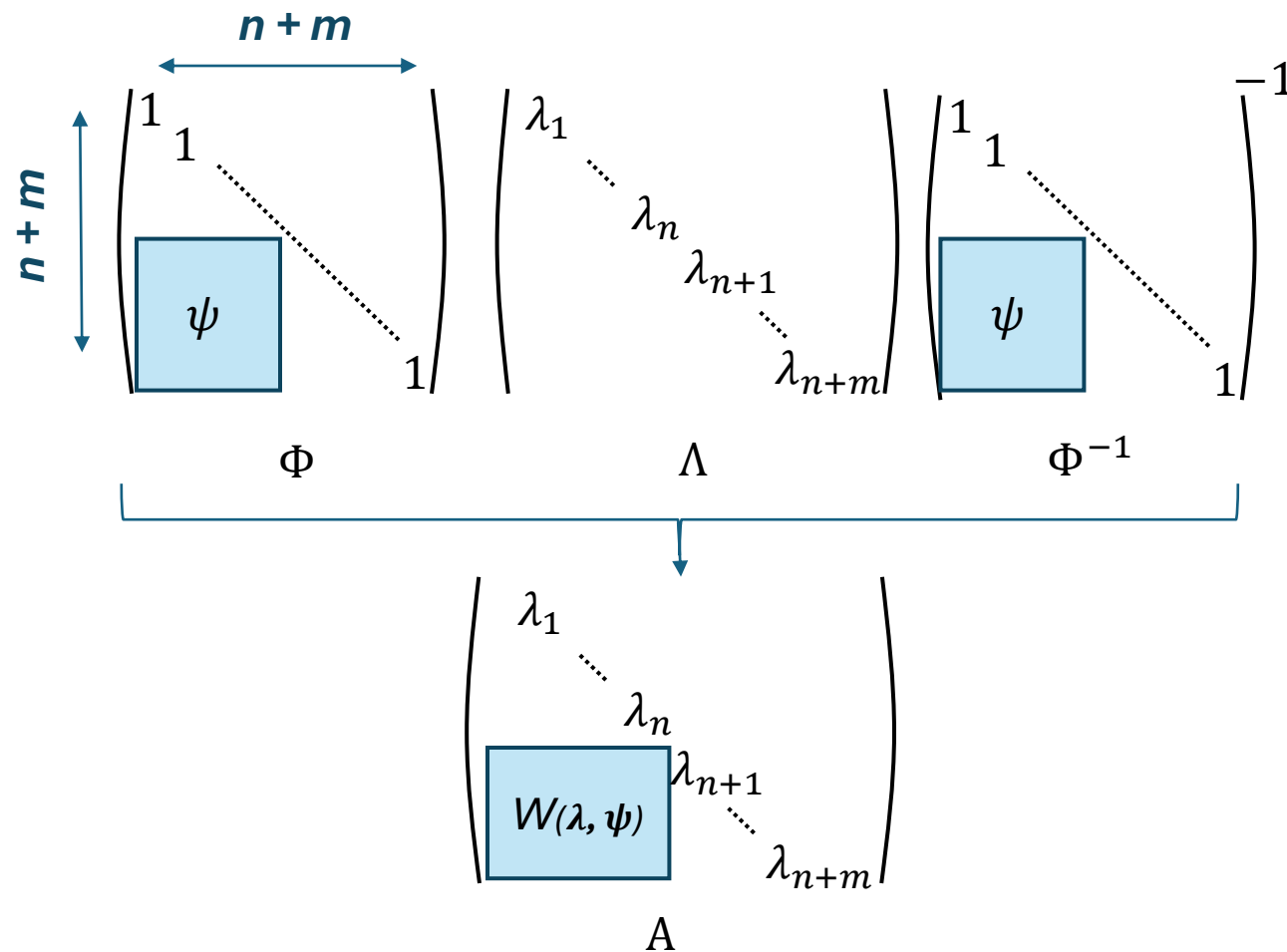
# Deep Learning in spectral domain

Starting from the reciprocal space, a suitable set of eigenvectors $\Phi$ can be chosen so as to obtain in the direct space a matrix A that describes the linear transfer between two consecutive layers.
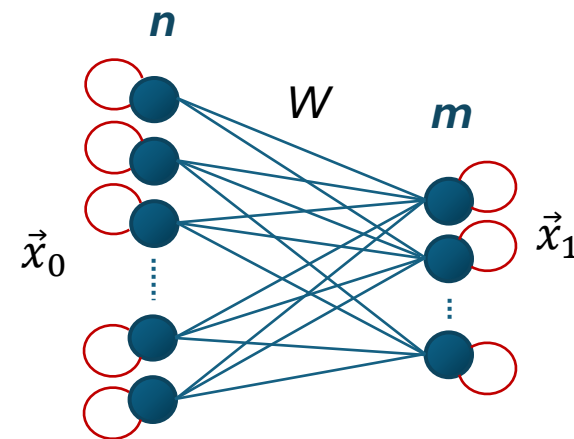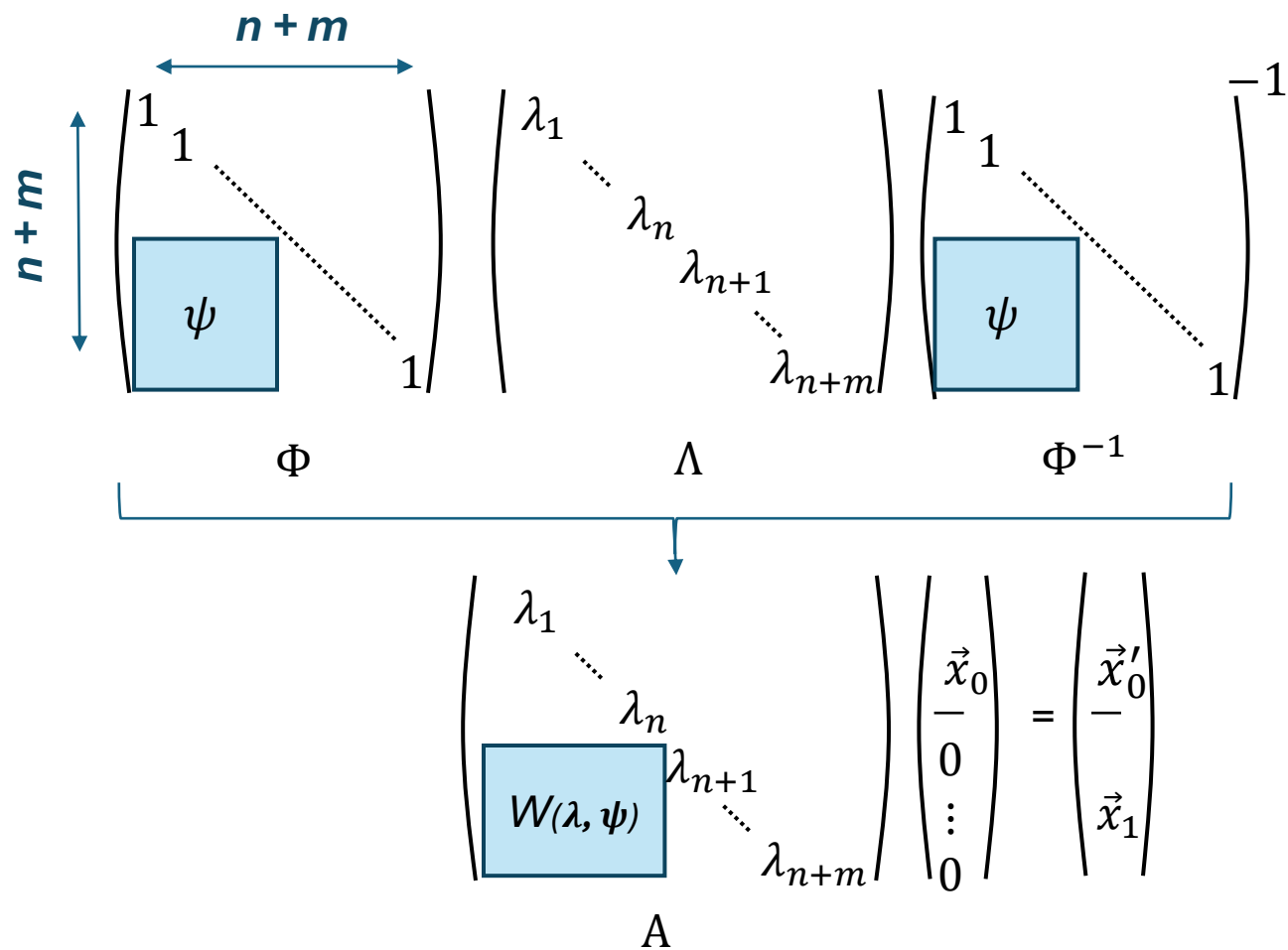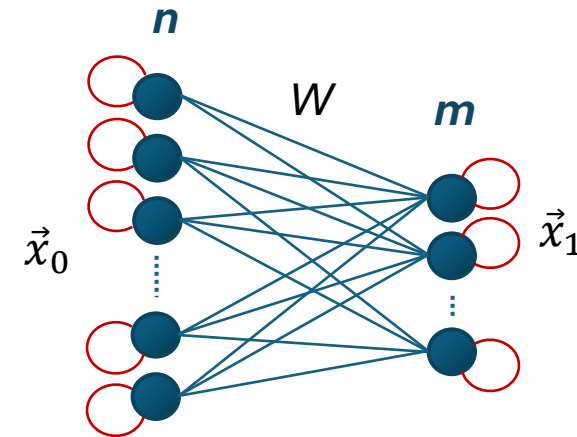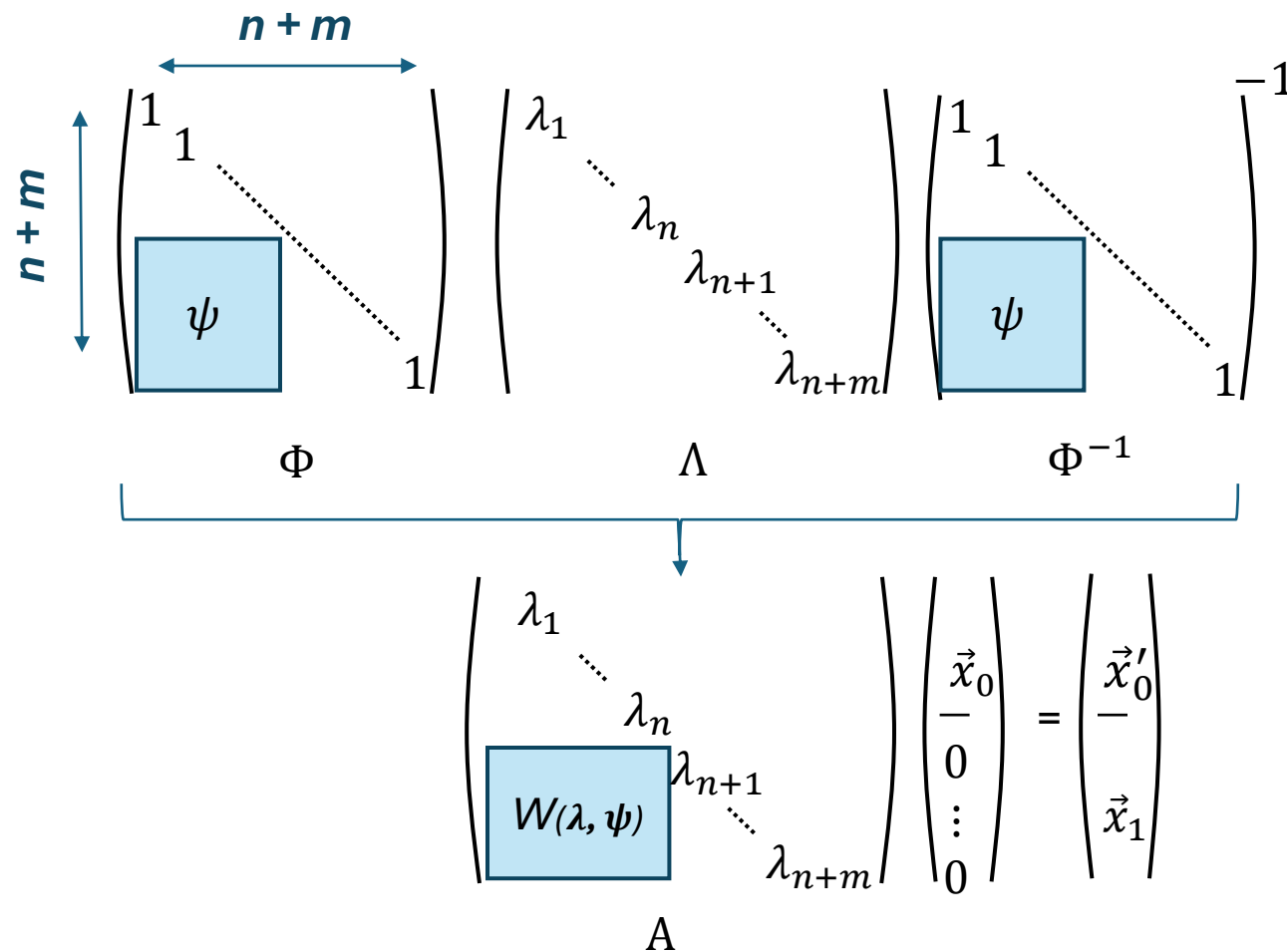
# Deep Learning in spectral domain

Starting from the reciprocal space, a suitable set of eigenvectors $\Phi$ can be chosen so as to obtain in the direct space a matrix A that describes the linear transfer between two consecutive layers.
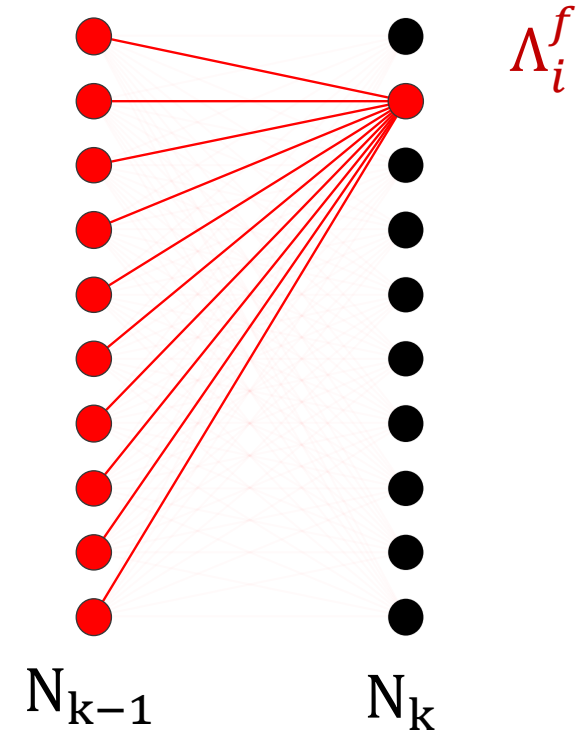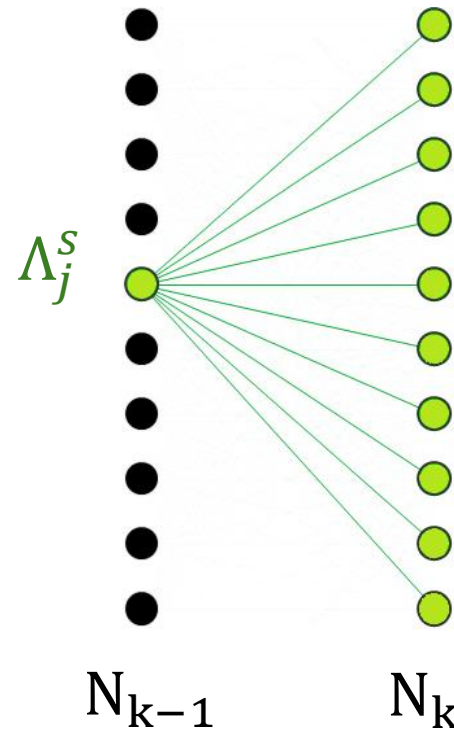


This leads to a new set of trainable parameters

$$w_{ij} = (\lambda_j - \lambda_i)\psi_{ij}$$
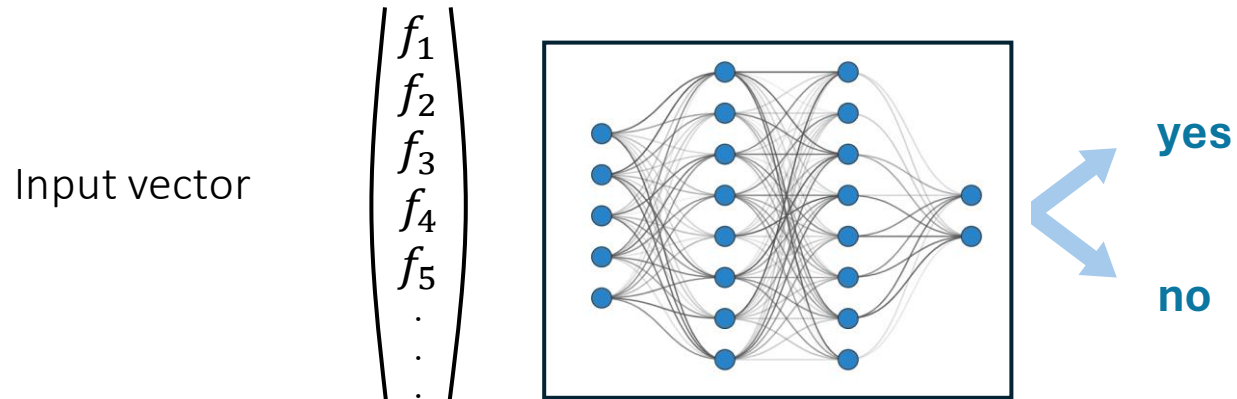
# Deep Learning in spectral domain

$$\sum_j w_{ij}^{(k)} x_j^{(k-1)} = \sum_j \psi_{ij} \Lambda_j^s x_j^{(k-1)} - \Lambda_i^f \psi_{ij} x_j^{(k-1)}$$

The eigenvalues identify the way information flows in the linear transfer

# Input feature relevance

Explainability: what does the network is looking at to take its decisions?

Input vector

$$\begin{pmatrix} f_1 \\ f_2 \\ f_3 \\ f_4 \\ f_5 \\ \cdot \\ \cdot \\ \cdot \end{pmatrix}$$
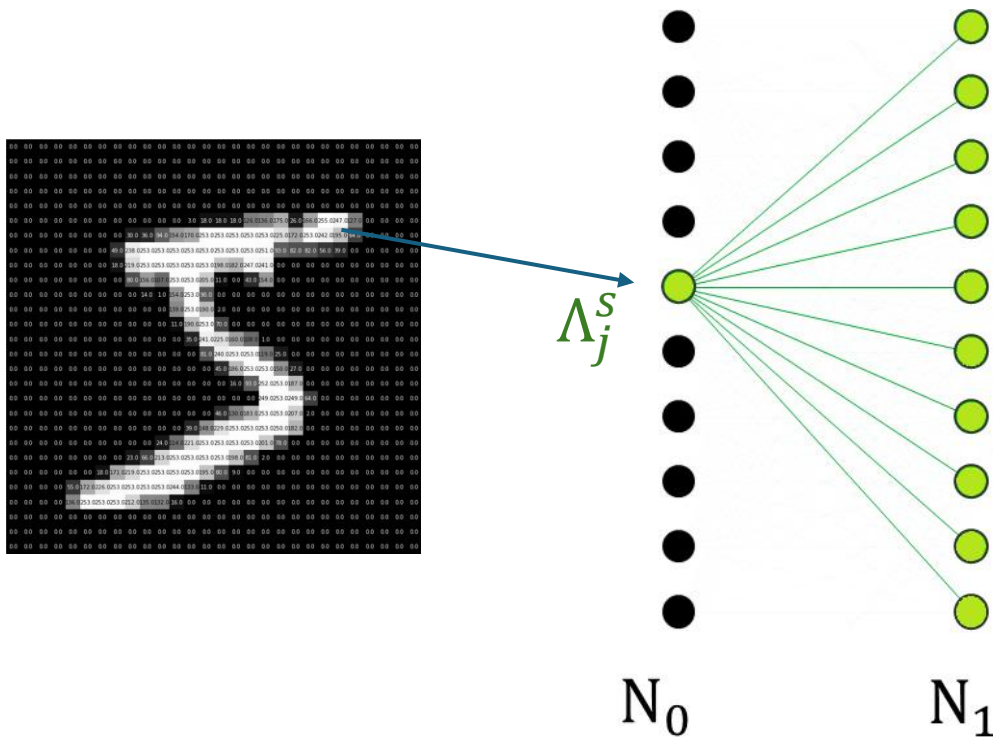


**yes**

**no**

It is crucial to define methods to identify what features are relevant!

# Input feature relevance

$$\sum_j w_{ij}^{(k)} x_j^{(k-1)} = \sum_j \psi_{ij} \Lambda_j^s x_j^{(k-1)} - \Lambda_i^f \psi_{ij} x_j^{(k-1)}$$

$$\Lambda_i^f = 0$$

$$\longrightarrow \qquad w_{ij} = \lambda_j \psi_{ij}$$

# Input feature relevance

$$\sum_j w_{ij}^{(k)} x_j^{(k-1)} = \sum_j \psi_{ij} \Lambda_j^s x_j^{(k-1)} - \Lambda_i^f \psi_{ij} x_j^{(k-1)}$$

$$\Lambda_i^f = 0$$

$$\longrightarrow \qquad w_{ij} = \lambda_j \psi_{ij}$$
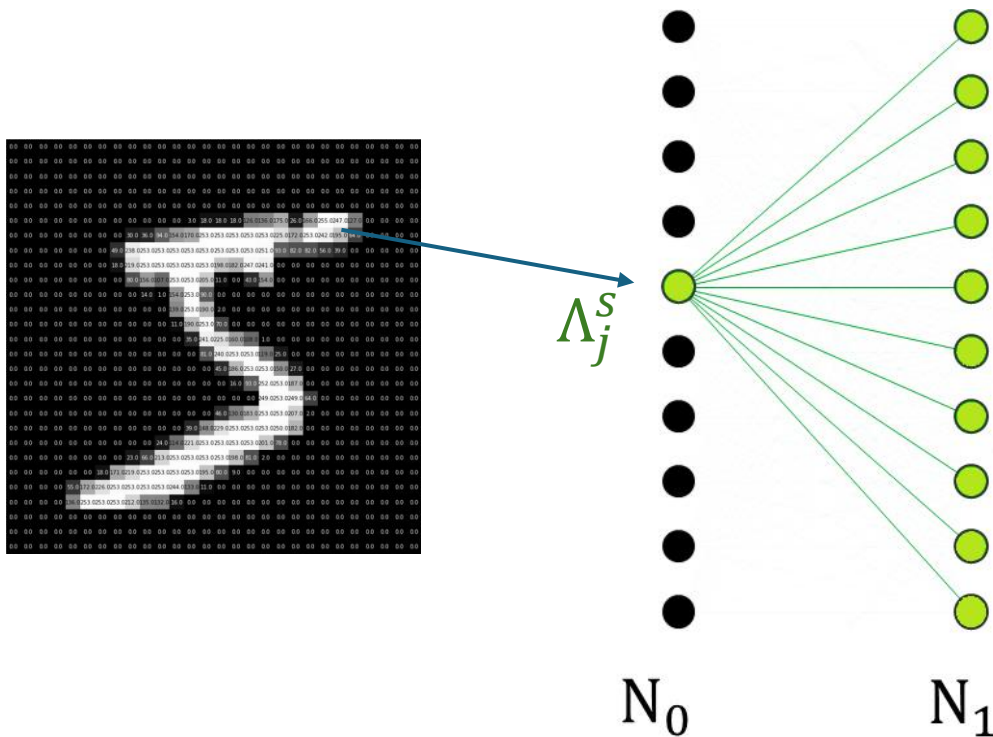


$\Lambda_j^s$

$N_0 \qquad N_1$

Not all the components of the input are important to solve the task

**Can eigenvalues be a proxy of the relevance?**

22

# Input feature relevance

$$\Lambda_i^f = 0$$

$$\sum_j w_{ij}^{(k)} x_j^{(k-1)} = \sum_j \psi_{ij} \Lambda_j^s x_j^{(k-1)} - \Lambda_i^f \psi_{ij} x_j^{(k-1)}$$

$$\longrightarrow \qquad w_{ij} = \lambda_j \psi_{ij}$$



$$\Lambda_j^s$$
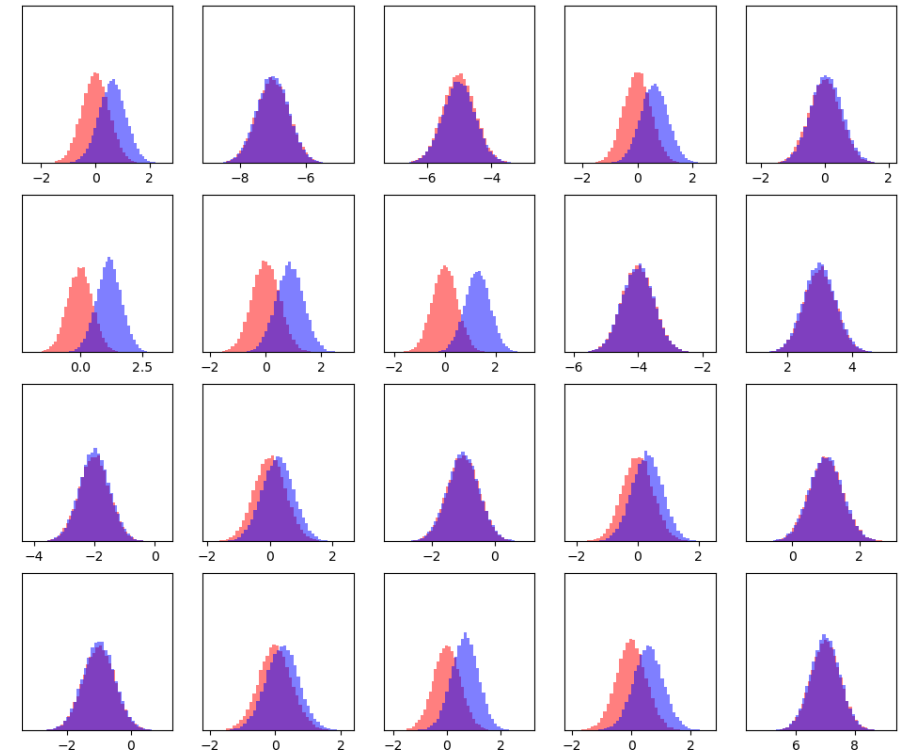
$$N_0 \qquad N_1$$

Not all the components of the input are important to solve the task

**Can eigenvalues be a proxy of the relevance?**
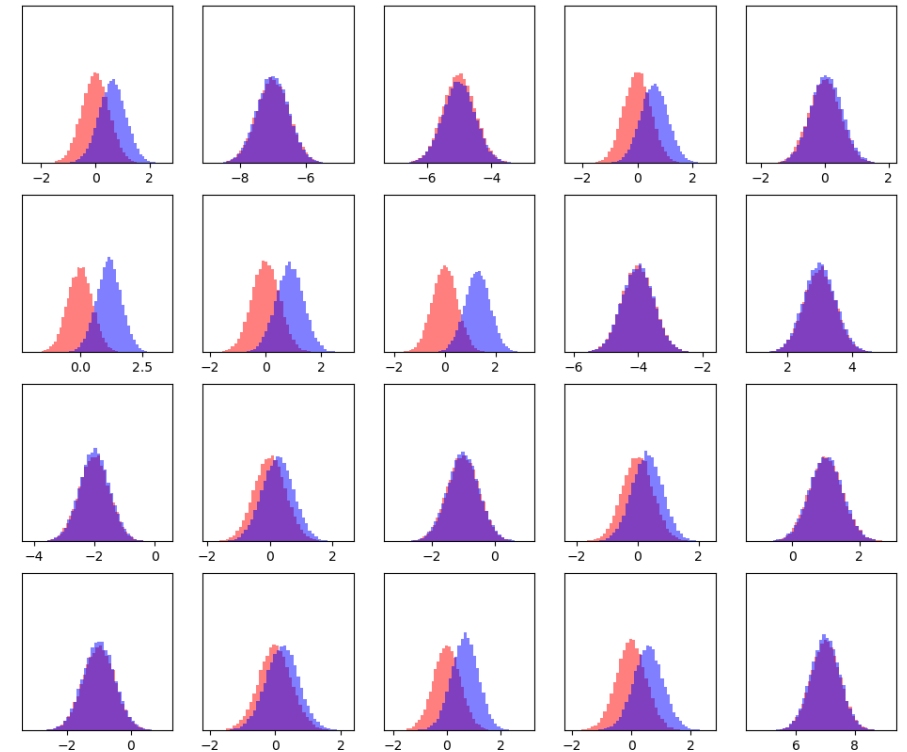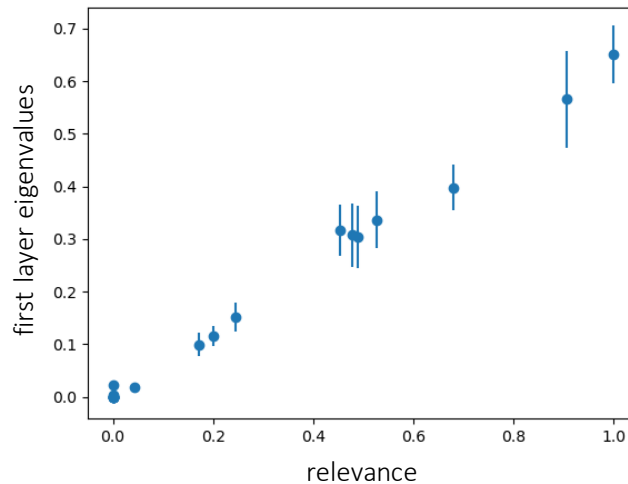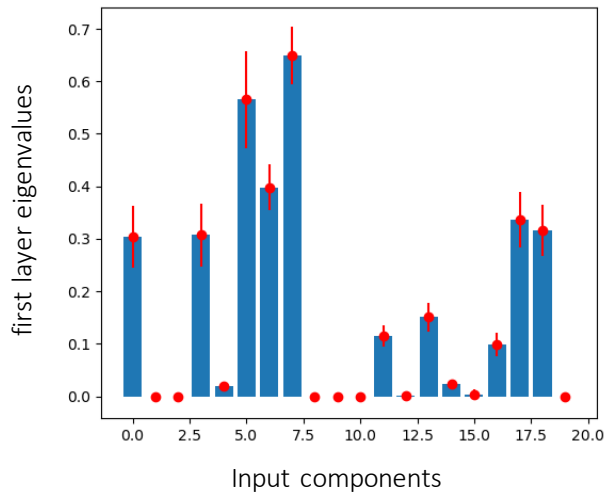
23

# A simple handmade dataset

- Classification problem with two classes $c \in [0,1]$
- $\vec{x}_0 \in R^{20}$
- $x_i$ from independent Gaussian distributions $(\mu_i^c, \sigma_i^c)$
- Some features are irrelevant: $\mu_i^0 = \mu_i^1$ and $\sigma_i^0 = \sigma_i^1$
- For some components, the two classes are partially separated. We define the relevance of the features as the *distance* between the two distributions.



distributions of the input features

# A simple handmade dataset

- Classification problem with two classes $c \in [0,1]$
- $\vec{x}_0 \in R^{20}$
- $x_i$ from independent Gaussian distributions $(\mu_i^c, \sigma_i^c)$
- Some features are irrelevant: $\mu_i^0 = \mu_i^1$ and $\sigma_i^0 = \sigma_i^1$
- For some components, the two classes are partially separated. We define the relevance of the features as the *distance* between the two distributions.



distributions of the input features

# Finding correlations

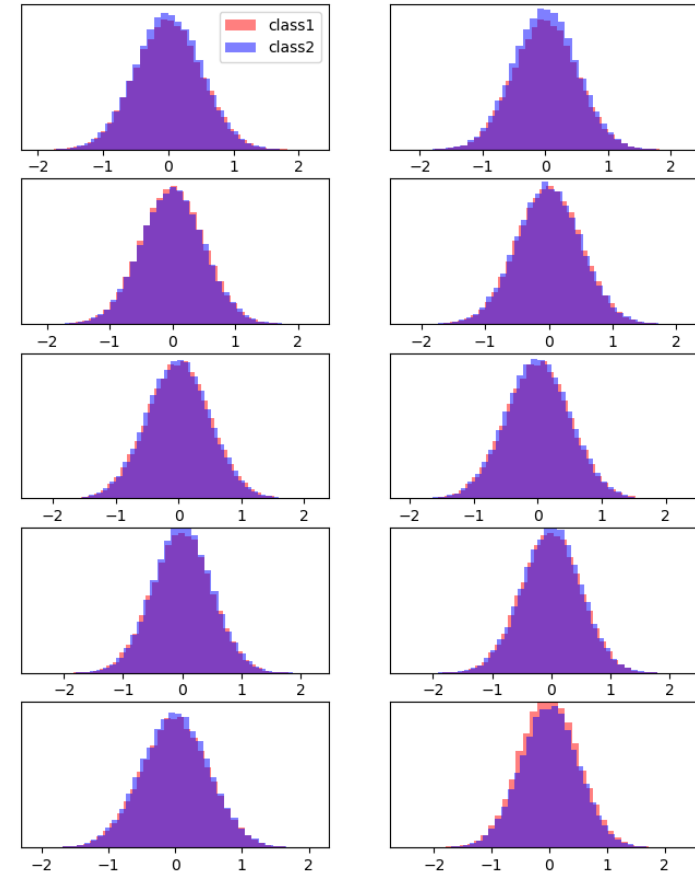- Classification problem with two classes $c \in [0,1]$

- $\vec{x}_0 \in R^{10}$

- if $c(\vec{x}) = 0$:

$$x_{2n+1} = x_{2n}$$

- if $c(\vec{x}) = 1$:

$$x_{2n+1} = \begin{cases} -x_{2n} & \text{with prob. } p \\ x_{2n} & \text{with prob. } 1-p \end{cases}$$



distributions of the input features

# Finding correlations

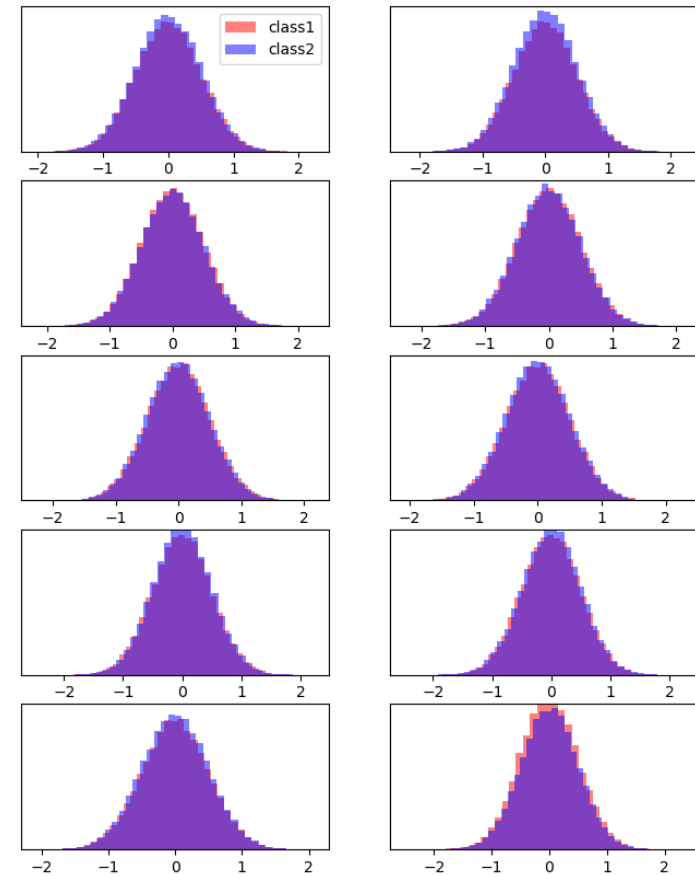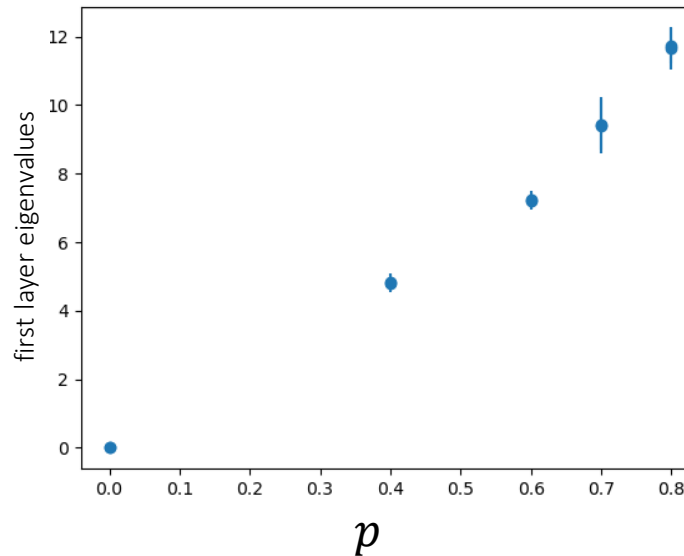- Classification problem with two classes $c \in [0,1]$
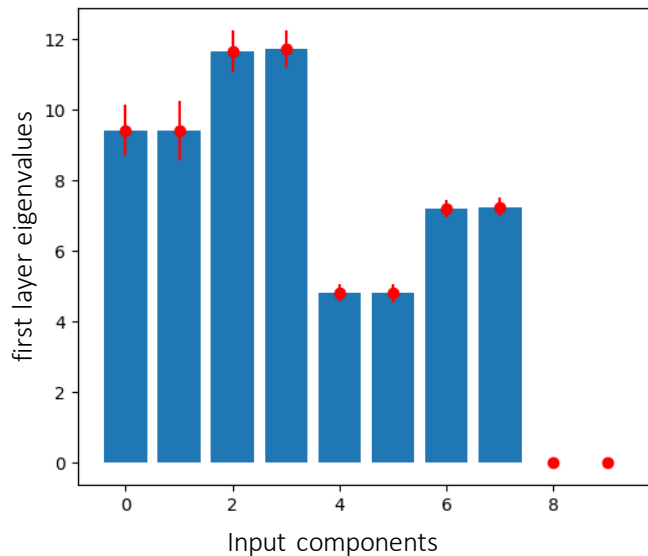- $\vec{x}_0 \in R^{10}$
- if $c(\vec{x}) = 0$:

$$x_{2n+1} = x_{2n}$$
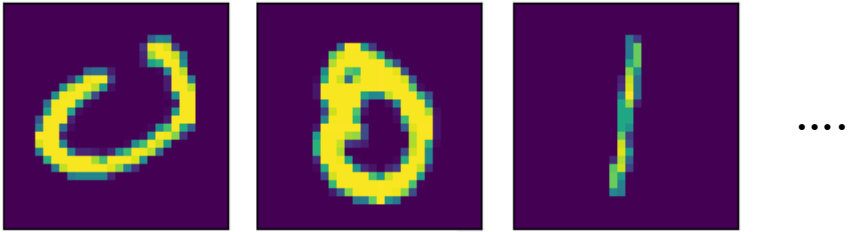
- if $c(\vec{x}) = 1$:

$$x_{2n+1} = \begin{cases} -x_{2n} & \text{with prob.} \, p \\ x_{2n} & \text{with prob.} \, 1-p \end{cases}$$



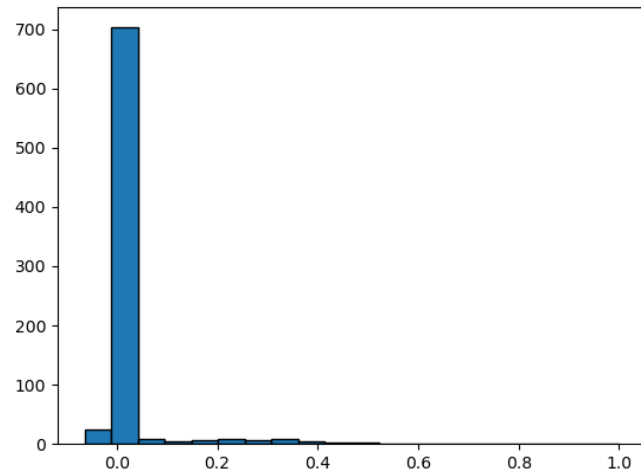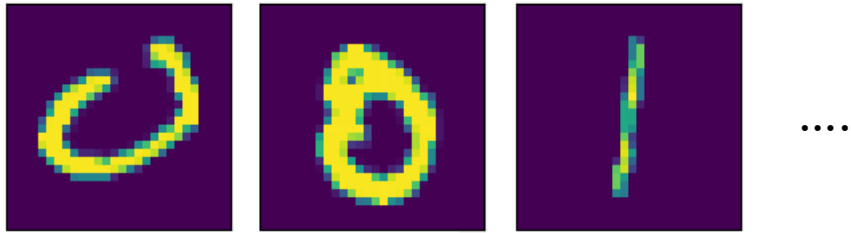

distributions of the input features

# The MNIST dataset

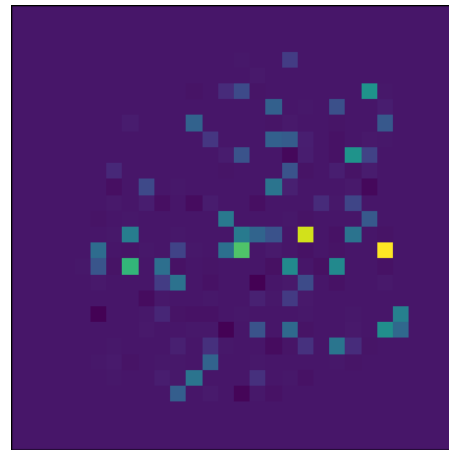A simple real dataset: images of zeros and ones $\longrightarrow \vec{x}_0 \in R^{784}$



....

# The MNIST dataset

A simple real dataset: images of zeros and ones $\longrightarrow$ $\vec{x}_0 \in R^{784}$



....



eigenvalues



Colormap: eigenvalues

# The MNIST dataset

A simple real dataset: images of zeros and ones $\longrightarrow$ $\vec{x}_0 \in R^{784}$



....


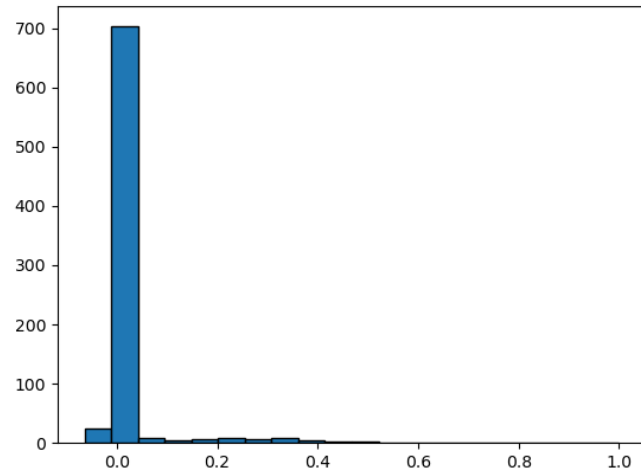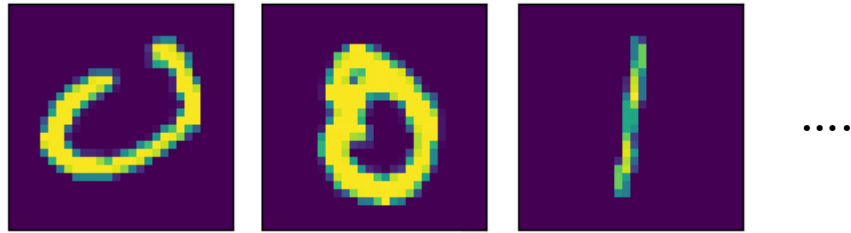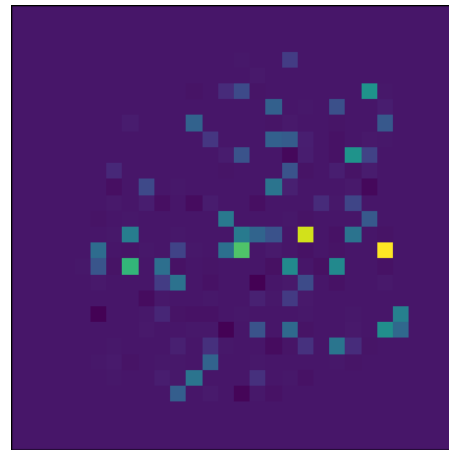
eigenvalues



Colormap: eigenvalues

# The MNIST dataset

A simple real dataset: images of zeros and ones $\longrightarrow$ $\vec{x}_0 \in R^{784}$
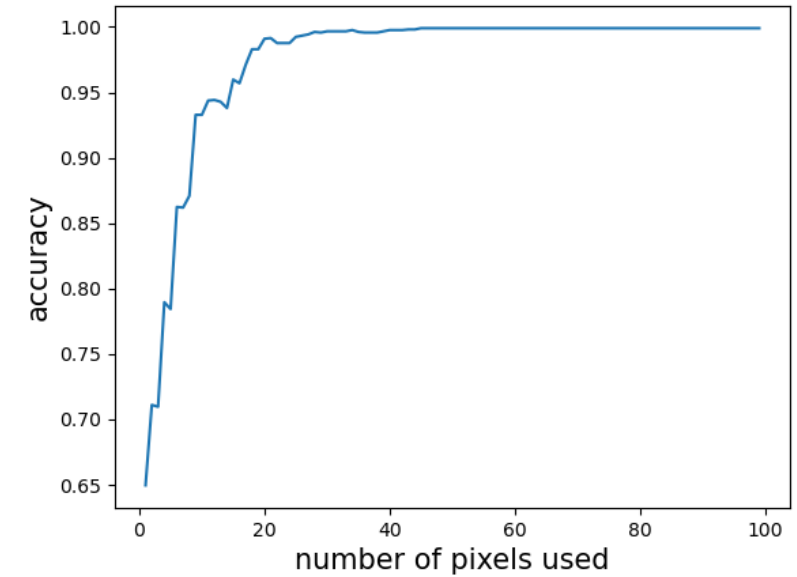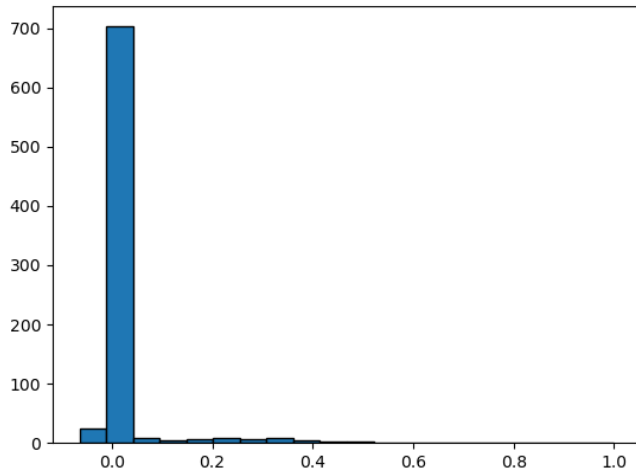




eigenvalues



Colormap: eigenvalues

# The MNIST dataset

By applying the mask to the input images, we see what the network is looking at:

# The MNIST dataset

By applying the mask to the input images, we see what the network is looking at:

# The MNIST dataset

By applying the mask to the input images, we see what the network is looking at:



To understand if the input is a "one", the network looks if this central part is active

# Conclusion

- **Spectral formulation** gives us node-related parameters by which we can rank nodes (pruning strategy)

- By applying this idea to **first-layer nodes** it is possible to identify relevant features: we tested on two handmade dataset and one very simple real dataset.

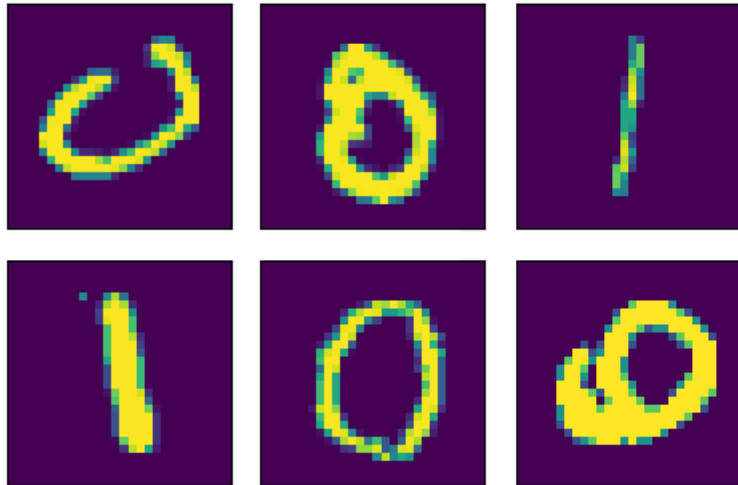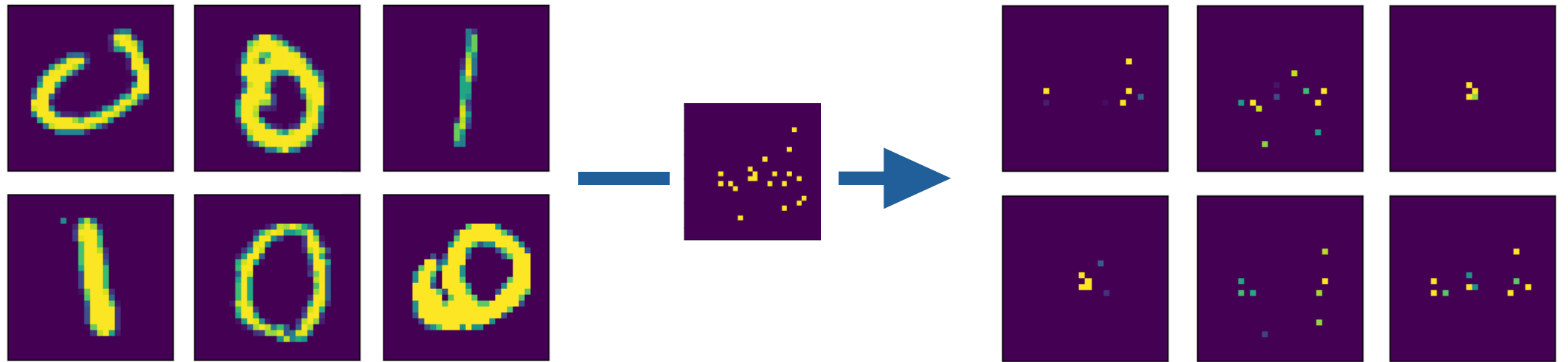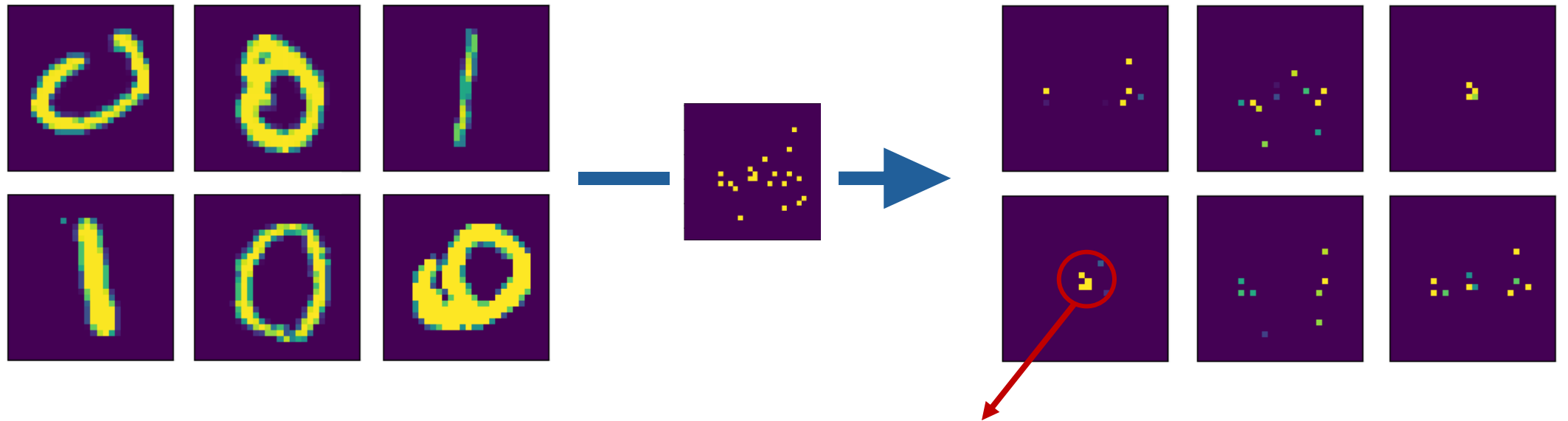- Real dataset are way more complicated, but this analysis is a first attempt to **connect post-training eigenvalues to input dimensions relevance**

---

**References:**

- Less parameters same performance  (Chicchi, Lorenzo, et al. *Physical Review E* 104.5 (2021): 054312.)

- Spectral Pruning  (Buffoni, Lorenzo, et al. *Scientific reports* 12.1 (2022): 1-9.)

- Recurrent Spectral Learning  (Chicchi, Lorenzo, et al., *Chaos, Solitons & Fractals* 168 (2023): 113128.)

- Complex Recurrent Spectral Network (Chicchi, Lorenzo, et al., *arXiv preprint arXiv:2312.07296* (2023).)

- How a student becomes a teacher: learning and forgetting through Spectral methods (Giambagli, Lorenzo, et al, *NeurIPS* 2023)

# Thank you!

# Deep Learning in spectral domain

$$\sum_j w_{ij}^{(k)} x_j^{(k-1)} = \sum_j \psi_{ij} \Lambda_j^s x_j^{(k-1)} - \Lambda_i^f \psi_{ij} x_j^{(k-1)}$$

$\Lambda_j^s = 0$



$\Lambda_i^f$

$N_{k-1}$        $N_k$

Buffoni, Lorenzo, et al. "Spectral pruning of fully connected layers." *Scientific Reports* 12.1 (2022): 11201.

# Deep Learning in spectral domain

$$\sum_j w_{ij}^{(k)} x_j^{(k-1)} = \sum_j \psi_{ij} \Lambda_j^s x_j^{(k-1)} - \Lambda_i^f \psi_{ij} x_j^{(k-1)}$$

$\Lambda_j^s = 0$

Every eigenvalues is now associated with one node : **Pruning strategy**



$\Lambda_i^f$

$N_{k-1}$      $N_k$

Buffoni, Lorenzo, et al. "Spectral pruning of fully connected layers." *Scientific Reports* 12.1 (2022): 11201.
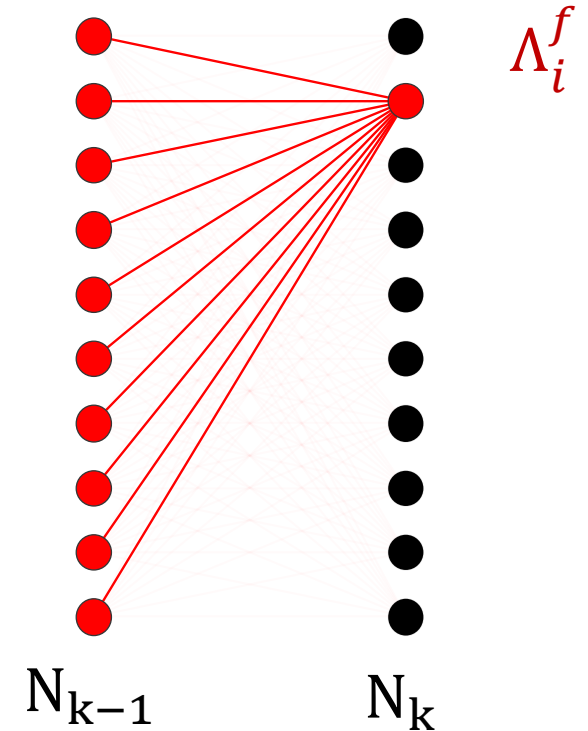
# Deep Learning in spectral domain

$$\sum_j w_{ij}^{(k)} x_j^{(k-1)} = \sum_j \psi_{ij} \Lambda_j^s x_j^{(k-1)} - \Lambda_i^f \psi_{ij} x_j^{(k-1)}$$

$\Lambda_j^s = 0$

Every eigenvalues is now associated with one node : **Pruning strategy**

$\Lambda_i^f$

$N_{k-1}$ $N_k$

Buffoni, Lorenzo, et al. "Spectral pruning of fully connected layers." *Scientific Reports* 12.1 (2022): 11201.