

CS5242 : Neural Networks and Deep Learning

Lecture 6: Multi-Layer Perceptron Depth

Semester 1 2021/22

Xavier Bresson

<https://twitter.com/xbresson>

Department of Computer Science
National University of Singapore (NUS)



Outline

- The three representations of neural networks
- The importance of non-linearities
- The importance of multiple layers
- Labs with MNIST and CIFAR

Outline

- The three representations of neural networks
- The importance of non-linearities
- The importance of multiple layers
- Labs with MNIST and CIFAR

The three representations

Computer Science

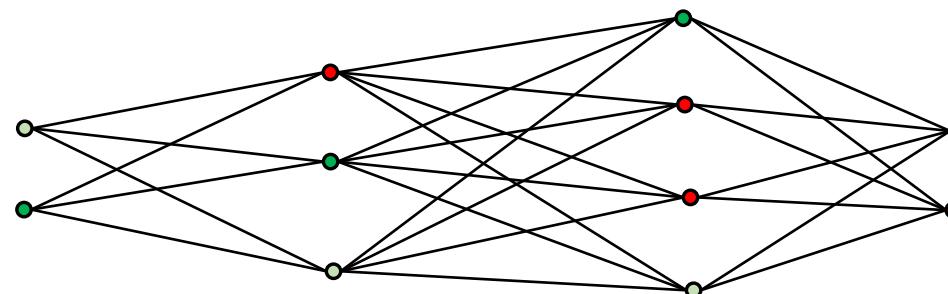
$$\begin{matrix} [0.3] \\ [0.9] \end{matrix} \boxed{\begin{bmatrix} u_{11} & u_{12} \\ u_{21} & u_{22} \\ u_{31} & u_{32} \end{bmatrix}} \begin{matrix} [-0.3] \\ [0.2] \\ [0.1] \end{matrix} \boxed{\sigma} \begin{matrix} [0.0] \\ [0.2] \\ [0.1] \end{matrix} \boxed{\begin{bmatrix} v_{11} & v_{12} & v_{13} \\ v_{21} & v_{22} & v_{23} \\ v_{31} & v_{32} & v_{33} \\ v_{41} & v_{42} & v_{43} \end{bmatrix}} \begin{matrix} [1.2] \\ [-0.7] \\ [-1.4] \\ [0.6] \end{matrix} \boxed{\sigma} \begin{matrix} [1.2] \\ [0.0] \\ [0.0] \\ [0.6] \end{matrix} \boxed{\begin{bmatrix} w_{11} & w_{12} & w_{13} & w_{14} \\ w_{21} & w_{22} & w_{23} & w_{24} \end{bmatrix}} \begin{matrix} [-2.8] \\ [-0.2] \end{matrix}$$

U σ V σ W

Mathematics

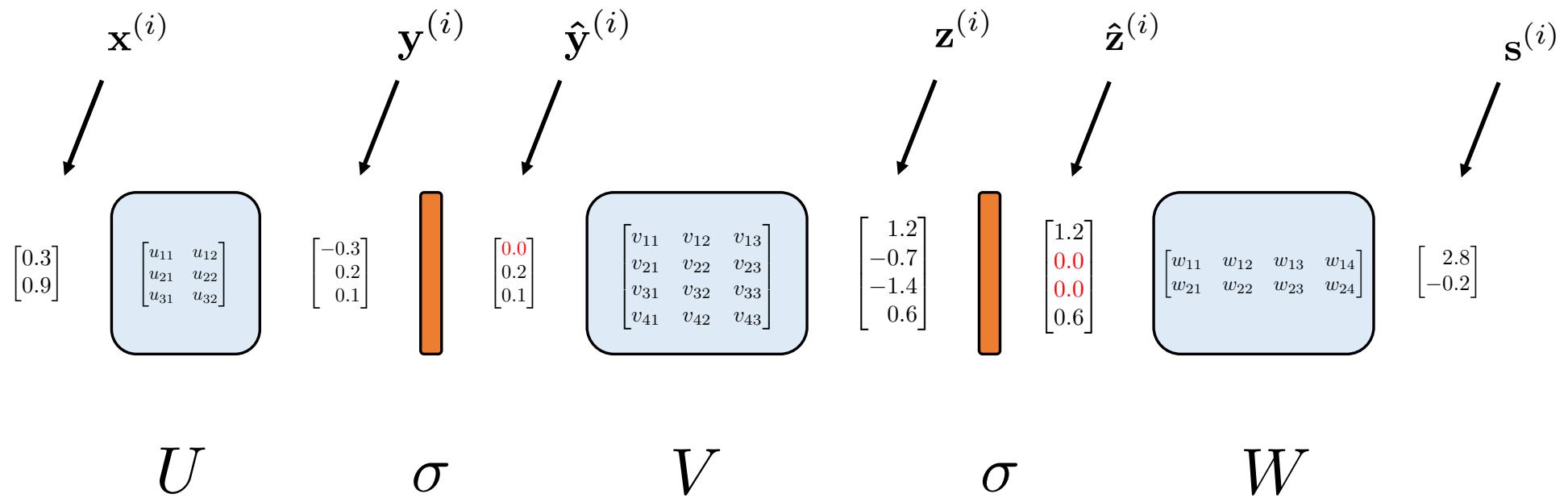
$$\mathbf{s}^{(\mathbf{i})} = W \ \sigma(V \ \sigma(U \mathbf{x}^{(\mathbf{i})}))$$

Biology



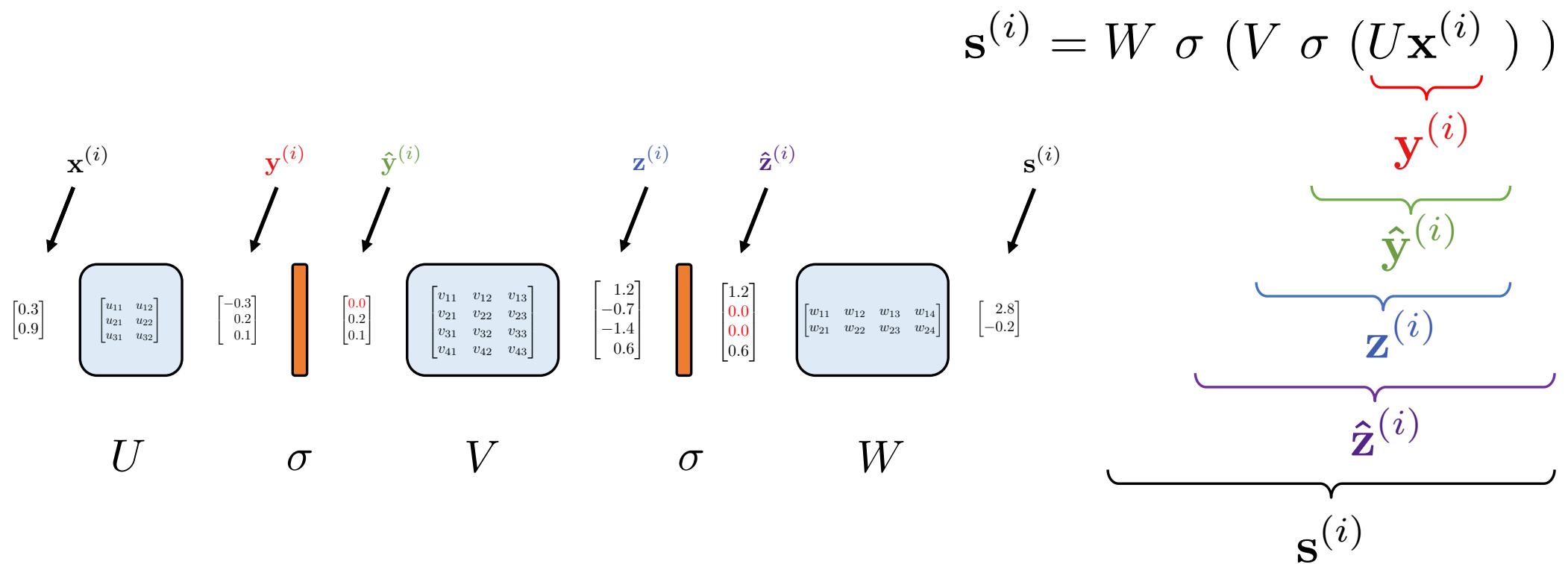
Computer science representation

- Computer science visualization of a neural network:



Mathematical representation

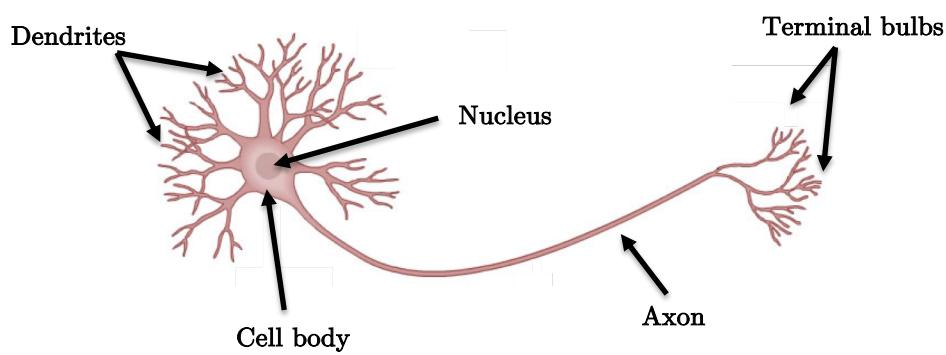
- Mathematical formulation of a neural network :



A neural net is simply a composition of linear and nonlinear functions !

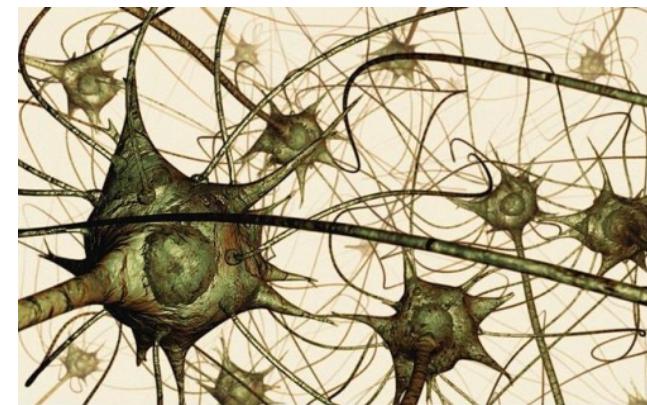
Biological representation

- A basic neuron:



- Brain = neural networks

(group of connected neurons)

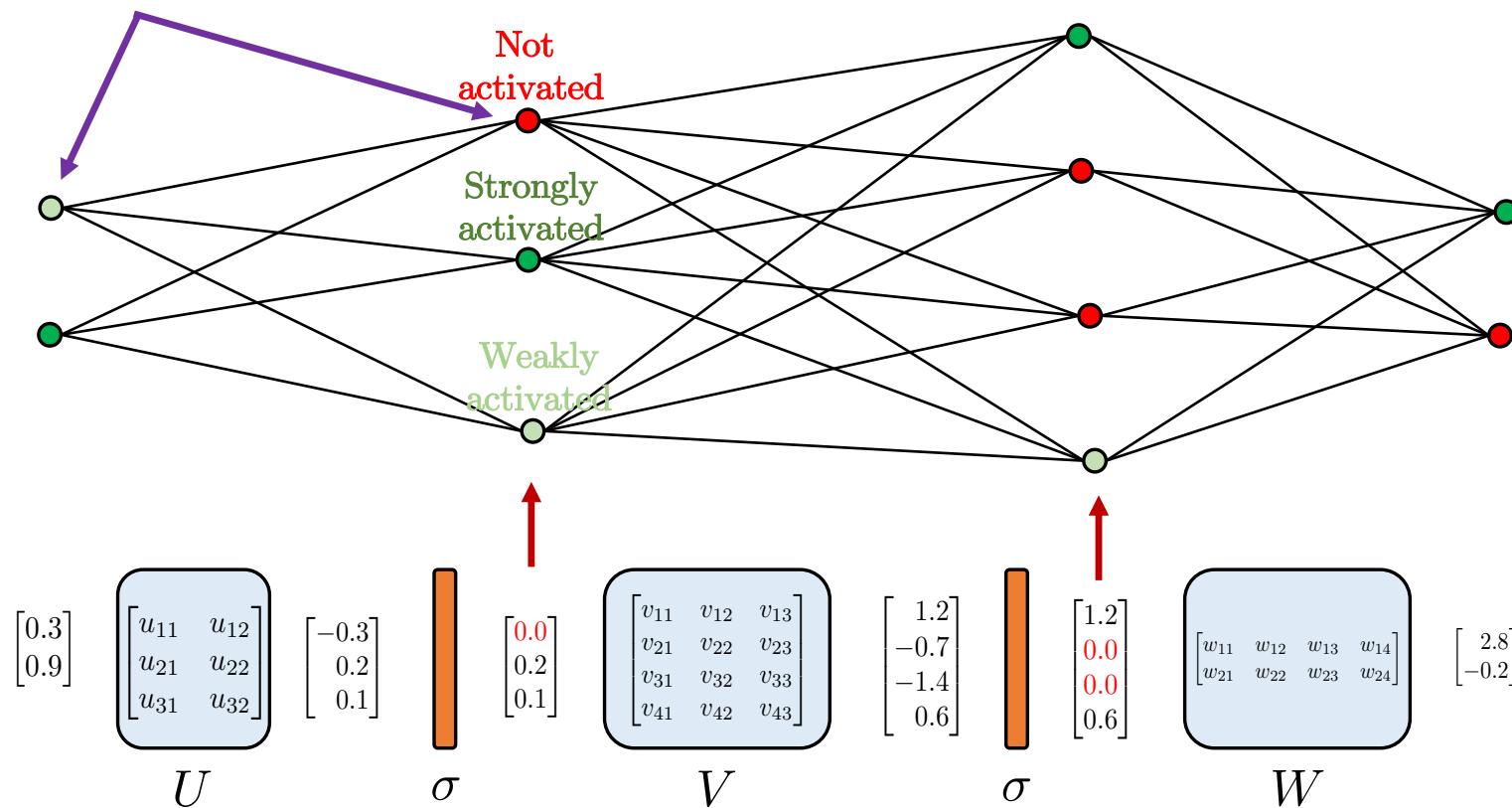


- A human brain has 10^{11} neurons and 10^{14} axons.

Biological representation

- Biological formulation of a neural network :

- Neurons



$$\begin{bmatrix} 0.3 \\ 0.9 \end{bmatrix} \quad \begin{bmatrix} u_{11} & u_{12} \\ u_{21} & u_{22} \\ u_{31} & u_{32} \end{bmatrix} \quad \begin{bmatrix} -0.3 \\ 0.2 \\ 0.1 \end{bmatrix}$$

$$\sigma \quad \begin{bmatrix} 0.0 \\ 0.2 \\ 0.1 \end{bmatrix}$$

$$\begin{bmatrix} v_{11} & v_{12} & v_{13} \\ v_{21} & v_{22} & v_{23} \\ v_{31} & v_{32} & v_{33} \\ v_{41} & v_{42} & v_{43} \end{bmatrix}$$

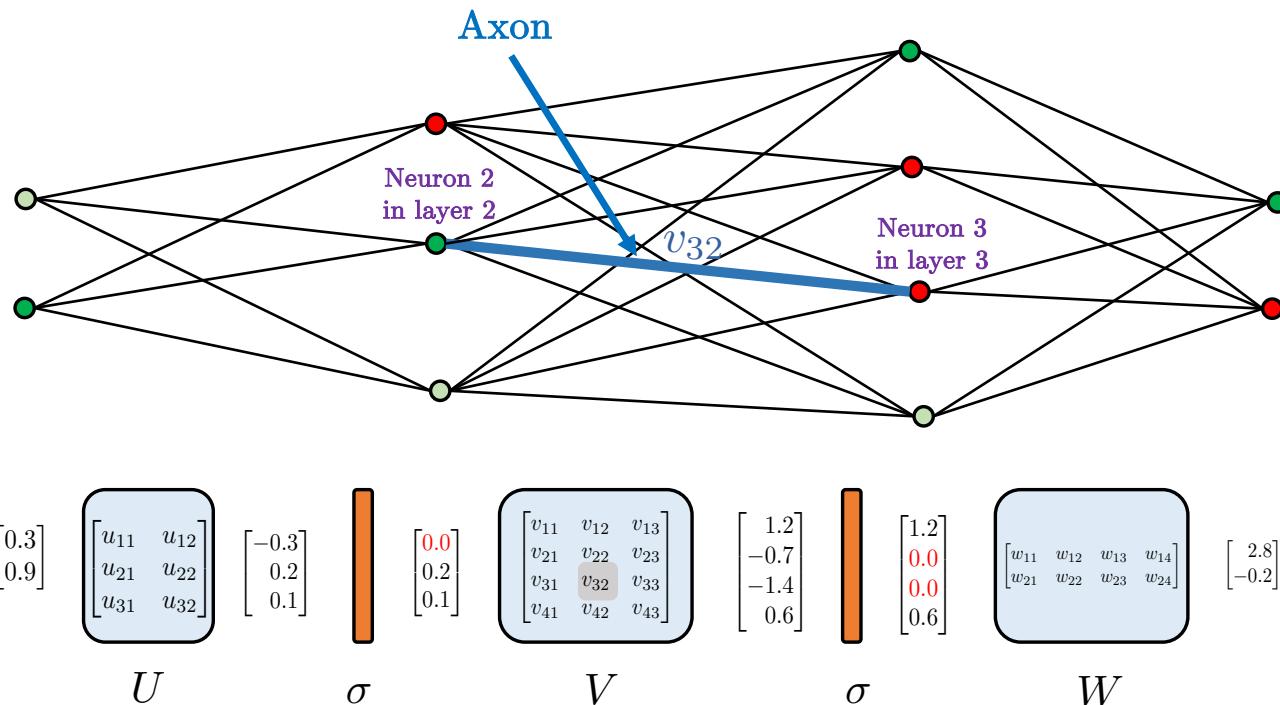
$$\begin{bmatrix} 1.2 \\ -0.7 \\ -1.4 \\ 0.6 \end{bmatrix}$$

$$\sigma \quad \begin{bmatrix} 1.2 \\ 0.0 \\ 0.0 \\ 0.6 \end{bmatrix}$$

$$\begin{bmatrix} w_{11} & w_{12} & w_{13} & w_{14} \\ w_{21} & w_{22} & w_{23} & w_{24} \end{bmatrix} \quad \begin{bmatrix} 2.8 \\ -0.2 \end{bmatrix}$$

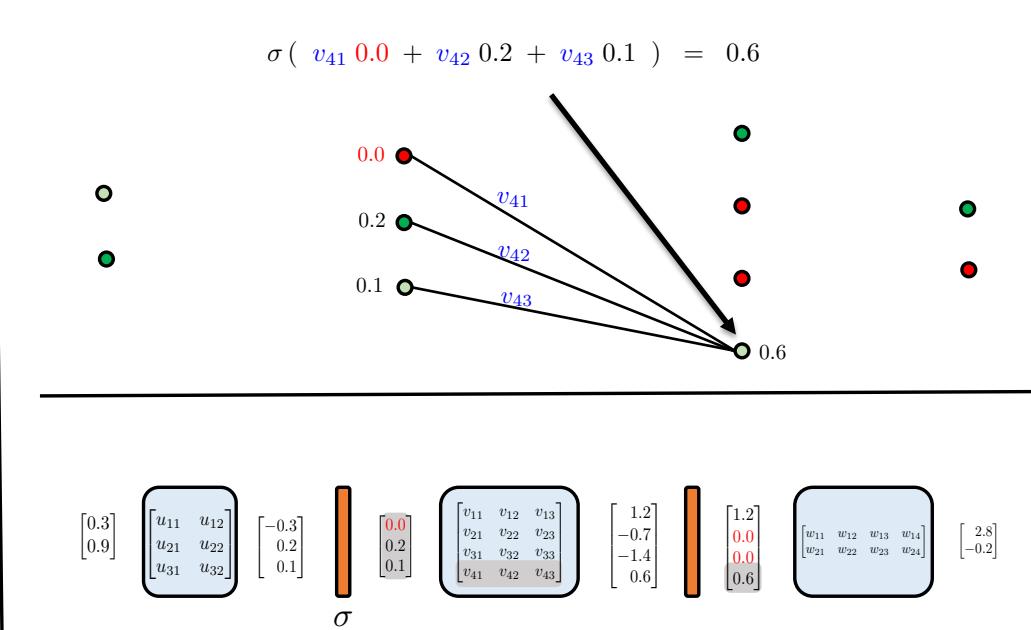
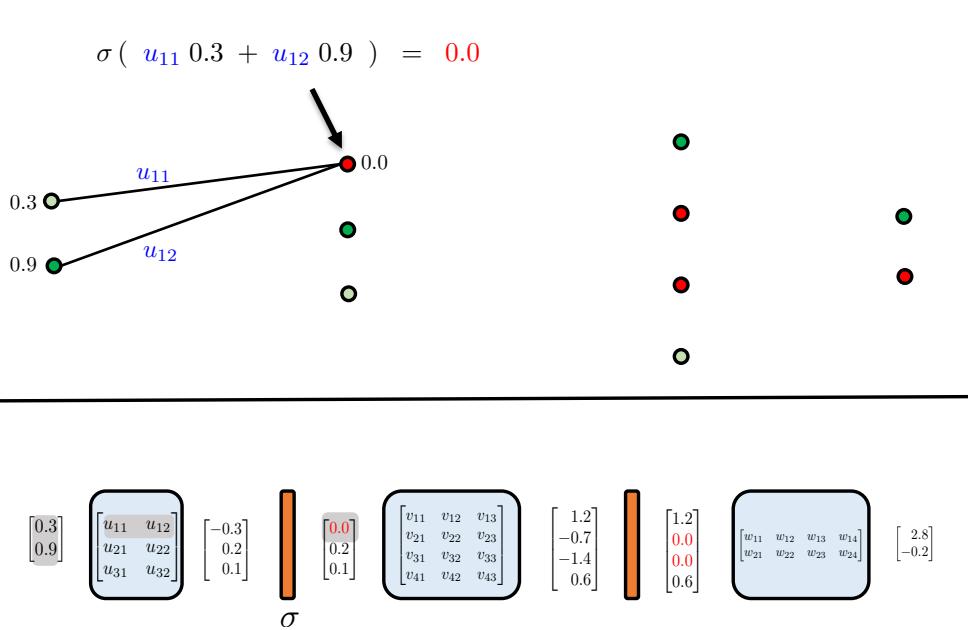
Biological representation

- Axons :
 - Each **weight** in the network can be thought as the **strength of a neuronal connection**.
 - Weight v_{32} = **strength of the connection** between neuron 2 in layer 2 and neuron 3 in layer 3.



Biological representation

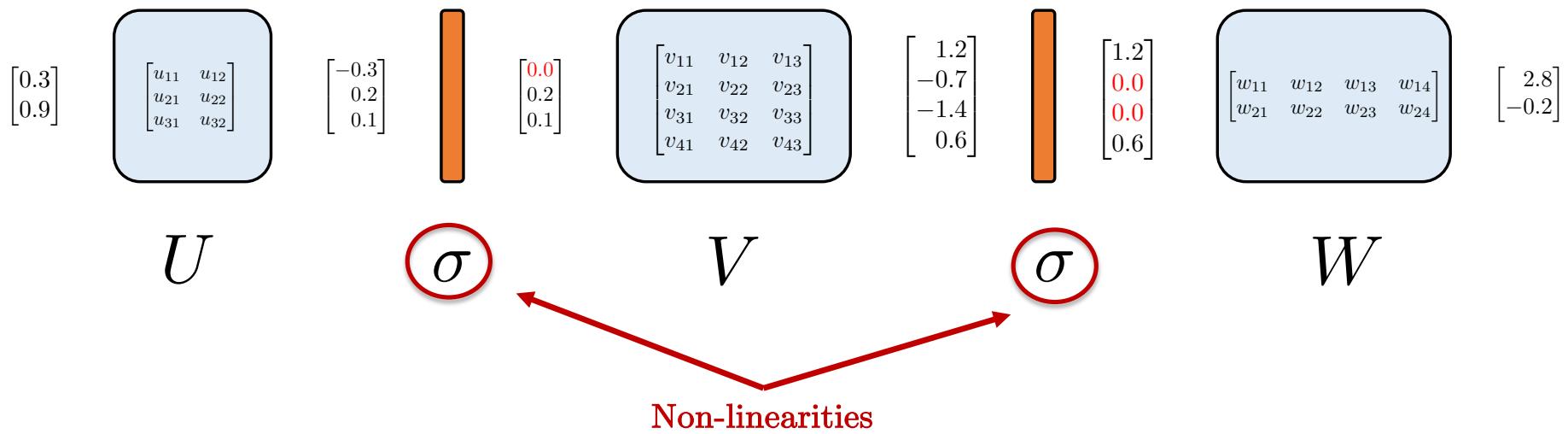
- Biological representation of an activation :



Outline

- The three representations of neural networks
- **The importance of non-linearities**
- The importance of multiple layers
- Labs with MNIST and CIFAR

The need of non-linearities



- Why non-linearities are introduced between layers?
 - Biological motivation
 - Mathematical motivation

Mathematical motivation

- Without non-linearities, the network degenerates as a linear model.

With nonlinearities:

$$\mathbf{s}^{(i)} = W\sigma(V\sigma(U\mathbf{x}^{(i)}))$$

Without nonlinearities:

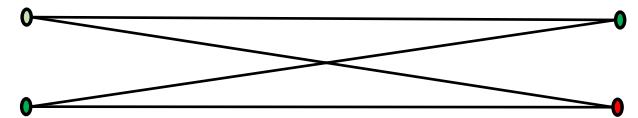
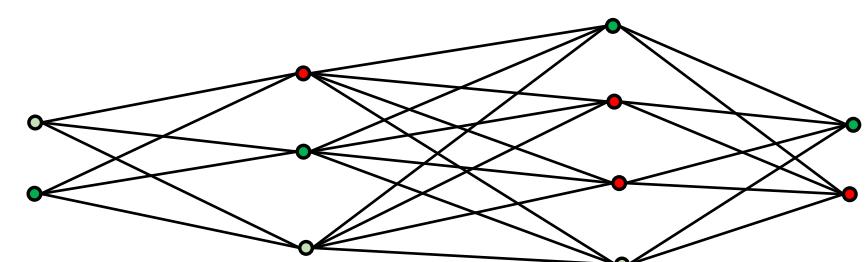
$$\mathbf{s}^{(i)} = \underbrace{WVU}_{A}\mathbf{x}^{(i)}$$

A product of matrices is just another matrix :

So the network collapses to a one-layer network !

Mathematical motivation

- Without non-linearities:



$$\begin{bmatrix} 0.3 \\ 0.9 \end{bmatrix} \begin{bmatrix} u_{11} & u_{12} \\ u_{21} & u_{22} \\ u_{31} & u_{32} \end{bmatrix} \begin{bmatrix} -0.3 \\ 0.2 \\ 0.1 \end{bmatrix} \quad \sigma \quad \begin{bmatrix} 0.0 \\ 0.2 \\ 0.1 \end{bmatrix}$$

U

$$\begin{bmatrix} v_{11} & v_{12} & v_{13} \\ v_{21} & v_{22} & v_{23} \\ v_{31} & v_{32} & v_{33} \\ v_{41} & v_{42} & v_{43} \end{bmatrix} \quad \sigma \quad \begin{bmatrix} 1.2 \\ -0.7 \\ -1.4 \\ 0.6 \end{bmatrix} \quad \begin{bmatrix} 0.0 \\ 0.0 \\ 0.6 \end{bmatrix}$$

V

$$\begin{bmatrix} w_{11} & w_{12} & w_{13} & w_{14} \\ w_{21} & w_{22} & w_{23} & w_{24} \end{bmatrix} \quad \begin{bmatrix} 2.8 \\ -0.2 \end{bmatrix}$$

W



$$\begin{bmatrix} 0.3 \\ 0.9 \end{bmatrix} \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} 2.8 \\ -0.2 \end{bmatrix}$$

A

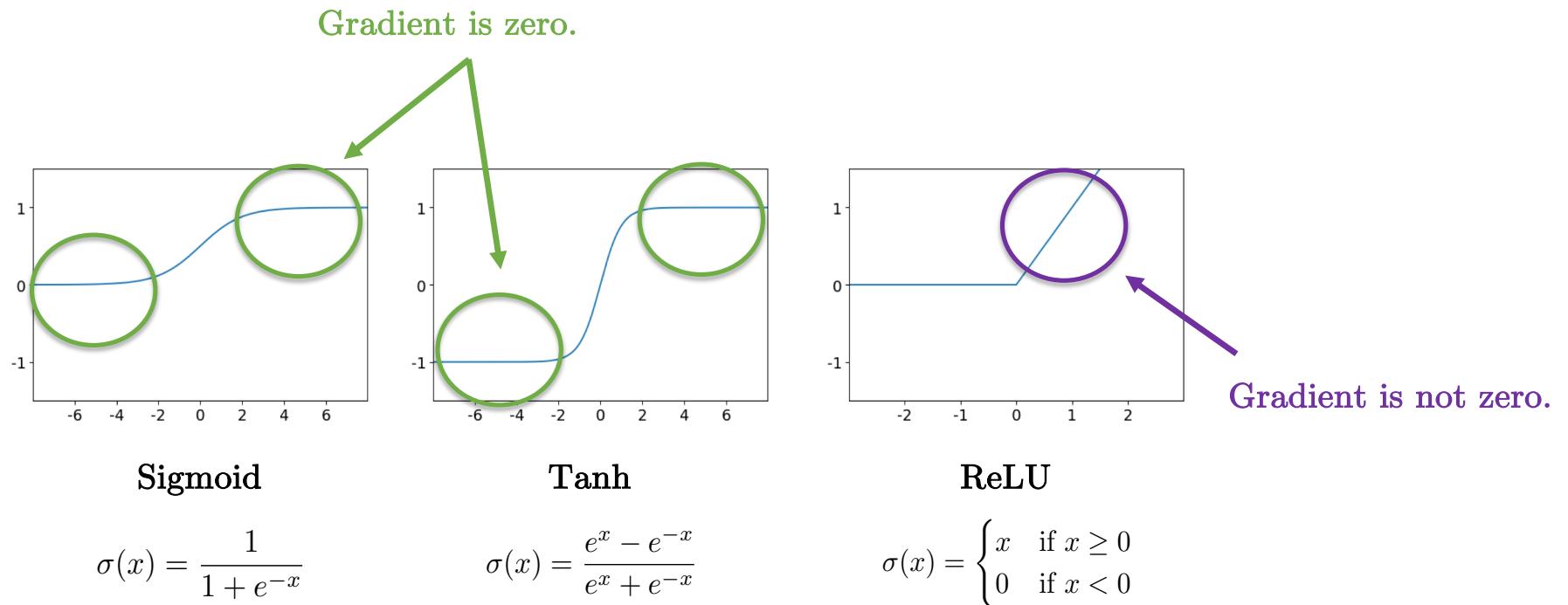
$$\mathbf{s}^{(i)} = W \sigma(V \sigma(U \mathbf{x}^{(i)}))$$



$$\mathbf{s}^{(i)} = A \mathbf{x}^{(i)}$$

Limitations of sigmoid and tanh

- Issues with Sigmoid and Tanh:
 - They are responsible for the vanishing gradient problem \Rightarrow Weights will not be updated.



Outline

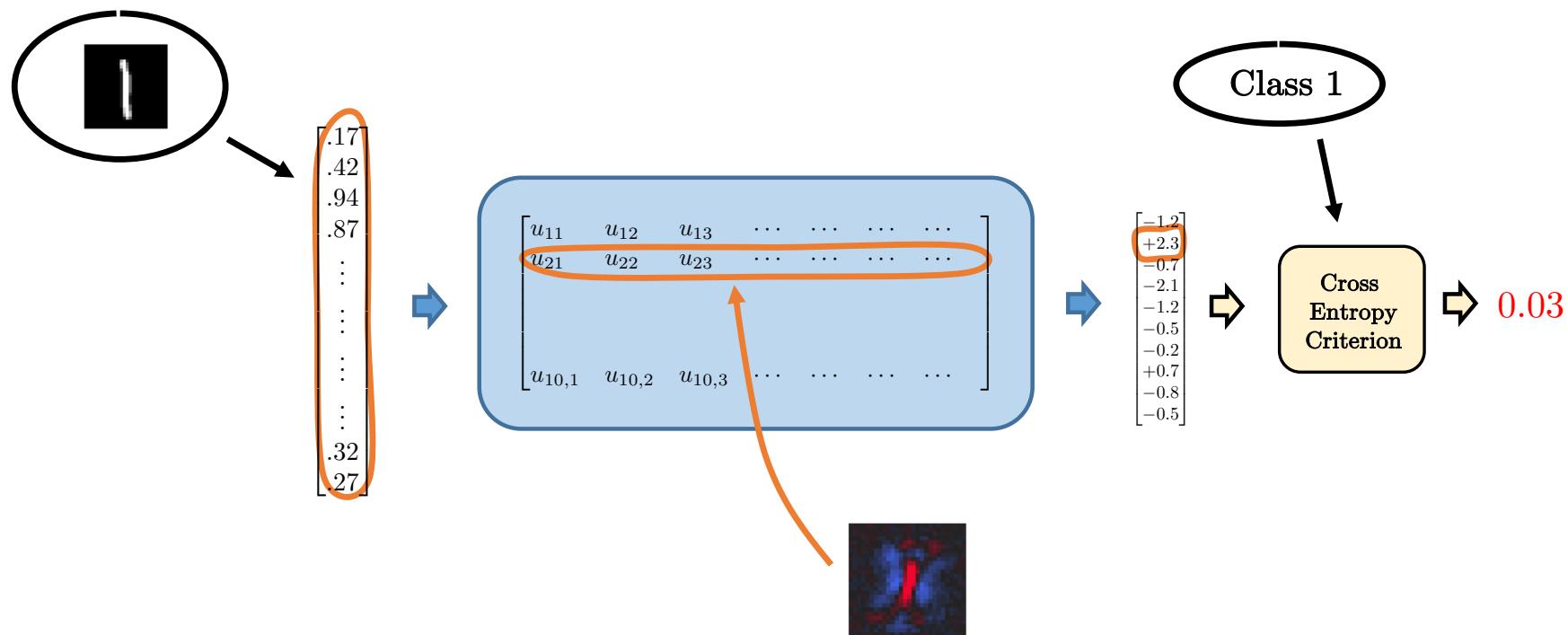
- The three representations of neural networks
- The importance of non-linearities
- **The importance of multiple layers**
- Labs with MNIST and CIFAR

An example

- Let us consider two classification problems to demonstrate the need of multiple layers :
 - An **easy** classification problem
 - A “**hard**” classification problem

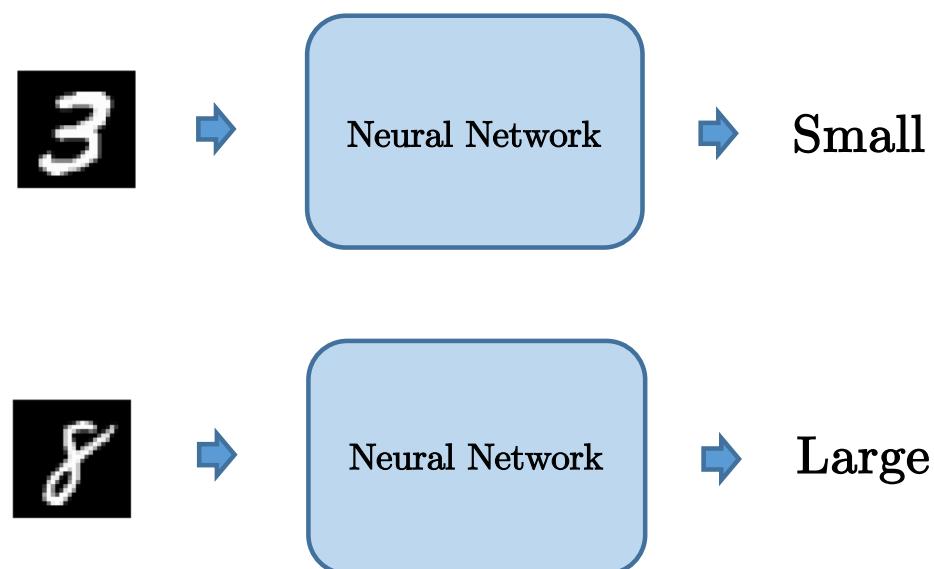
An easy classification problem

- Classify the digit numbers 0,1,2...9 :
 - One-layer net with 10 classes is enough to solve this classification problem.



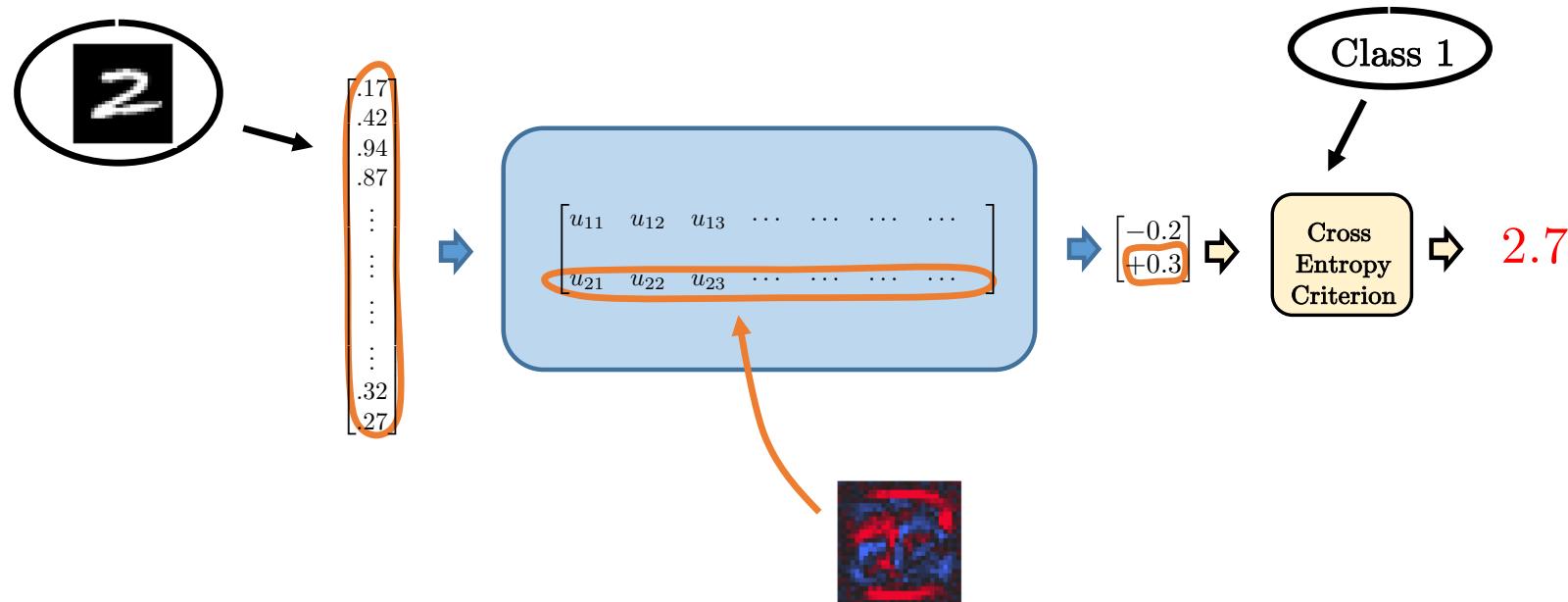
A hard classification problem

- A harder classification problem:
 - Classify small vs large digits into 2 categories
- Class 1: Small (0,1,2,3,4)
Class 2: Large (5,6,7,8,9)



A hard classification problem

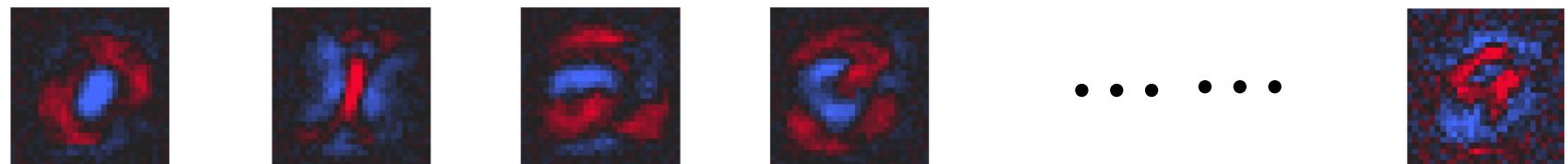
- Classify small vs large digits into 2 categories :
 - One-layer net with 2 classes is **not** enough to solve this classification problem.



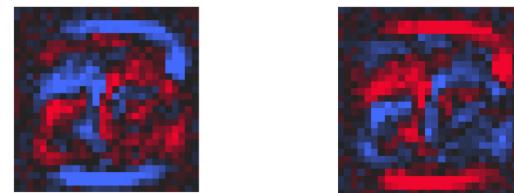
This pattern tries to represent digits
5-9 simultaneously with no success.

Template matching

- It is **easy** to find good templates for **zero, one, two, three, ..., nine** :

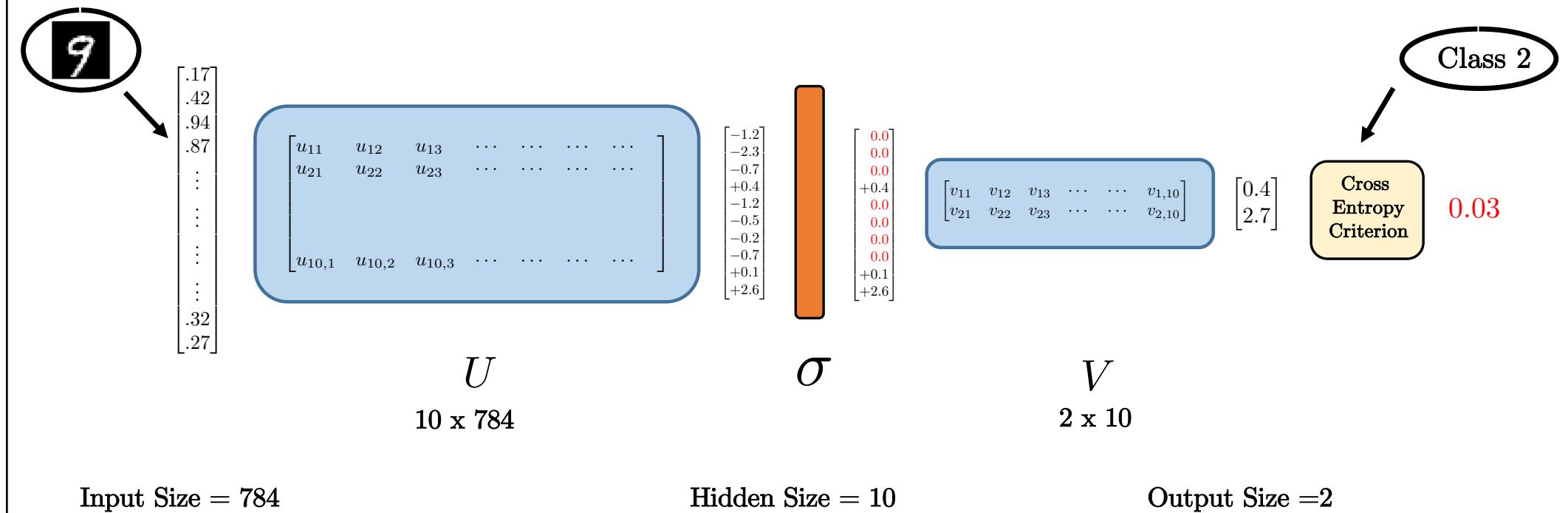


- It is **hard** to find good templates for **small versus large numbers** :



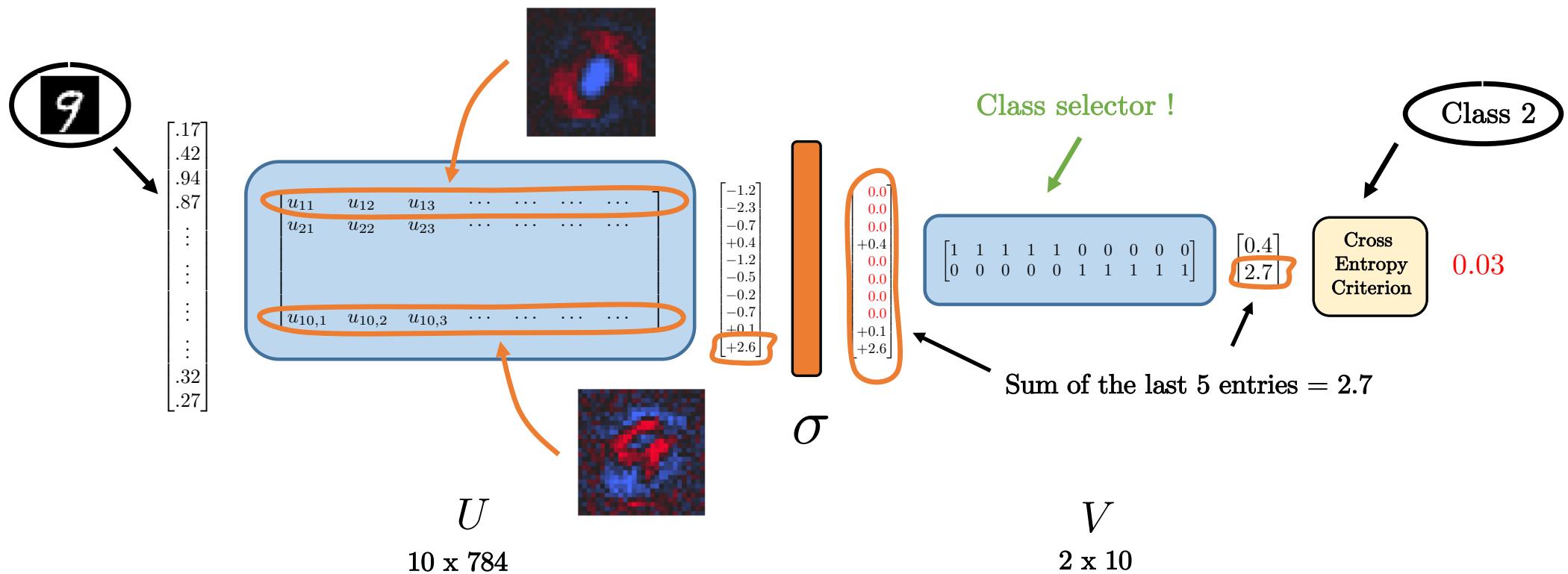
A hard classification problem

- Let's consider a **two-layer network** rather than just one-layer :



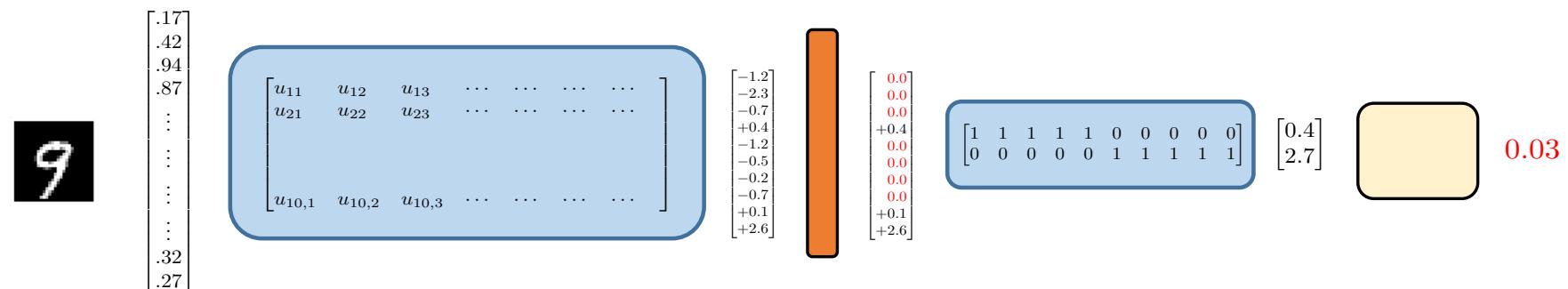
A hard classification problem

- After training, the network parameters U and V will become as follows :



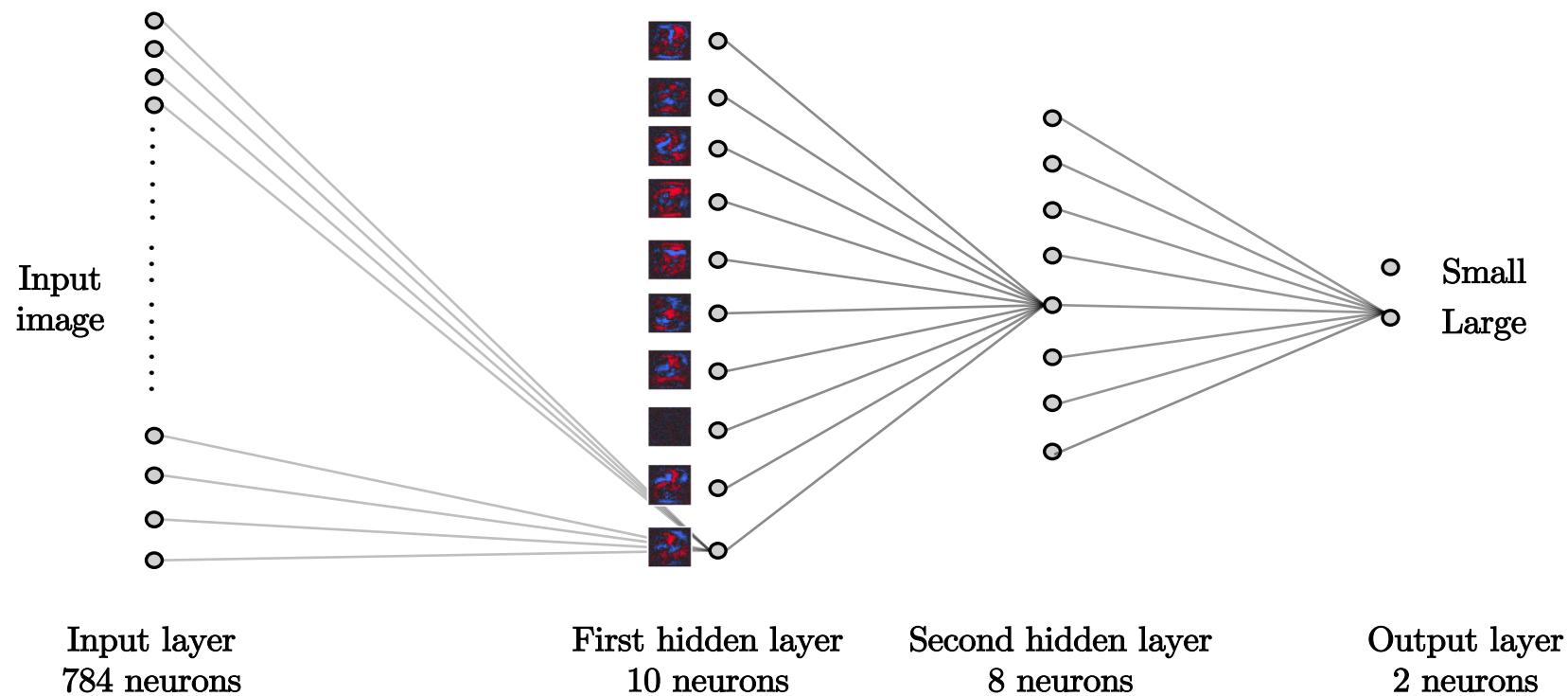
Multiple-step reasoning

- Multiple layers provide multiple-step reasoning :
 - The first layer is used to identify the type of digits
 - The second layer uses the classification made by the first layer, and processes it by pooling together 0-4 into one group, and 5-9 into another group.



Multiple-step reasoning

- What would happen with more layers?
 - The new network would also **succeed**.
 - The network would **learn relevant features** and do multiple stage reasoning.
 - The new features will be **less/no interpretable** (it is fine).

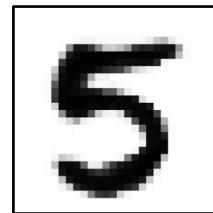


Outline

- The three representations of neural networks
- The importance of non-linearities
- The importance of multiple layers
- **Labs with MNIST and CIFAR**

MNIST dataset

- Number of classes: 10
- Training set: 50,000 labelled images
- Test set: 10,000 labelled images



28 x 28 grayscale image

0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9 9 9 9

Lab 01

- MNIST classification with MLP

Jupyter mnist_multilayer_demo Last Checkpoint: a few seconds ago (autosaved) Logout

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3

In [1]:

```
import torch
import torch.nn as nn
import torch.nn.functional as F
import torch.optim as optim
from random import randint
import time
import utils
```

In [2]:

```
from utils import check_mnist_dataset_exists
data_path=check_mnist_dataset_exists()

train_data=torch.load(data_path+'mnist/train_data.pt')
train_label=torch.load(data_path+'mnist/train_label.pt')
test_data=torch.load(data_path+'mnist/test_data.pt')
test_label=torch.load(data_path+'mnist/test_label.pt')
```

In [3]:

```
class two_layer_net(nn.Module):

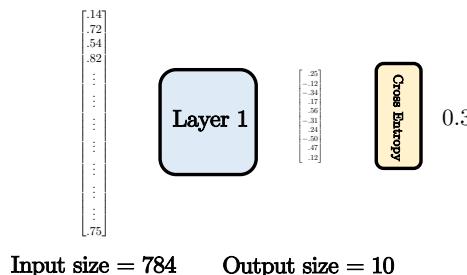
    def __init__(self, input_size, hidden_size, output_size):
        super(two_layer_net , self).__init__()

        self.layer1 = nn.Linear( input_size , hidden_size , bias=False )
        self.layer2 = nn.Linear( hidden_size , output_size , bias=False )
```

Lab 01

- Results

- One-layer network:



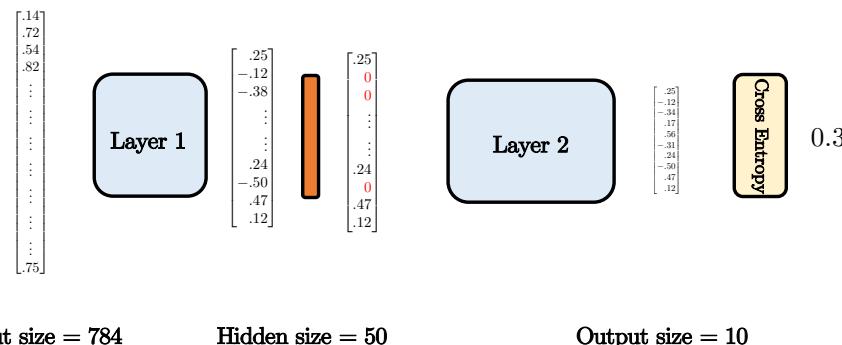
Number of parameters:

$$784 \times 10 = 7,840$$

Error on train set: 8%

Error on test set: 8%

- Two-layer network:



Number of parameters:

$$784 \times 50 + 50 \times 10 = 39,700$$

Error on train set: 0%

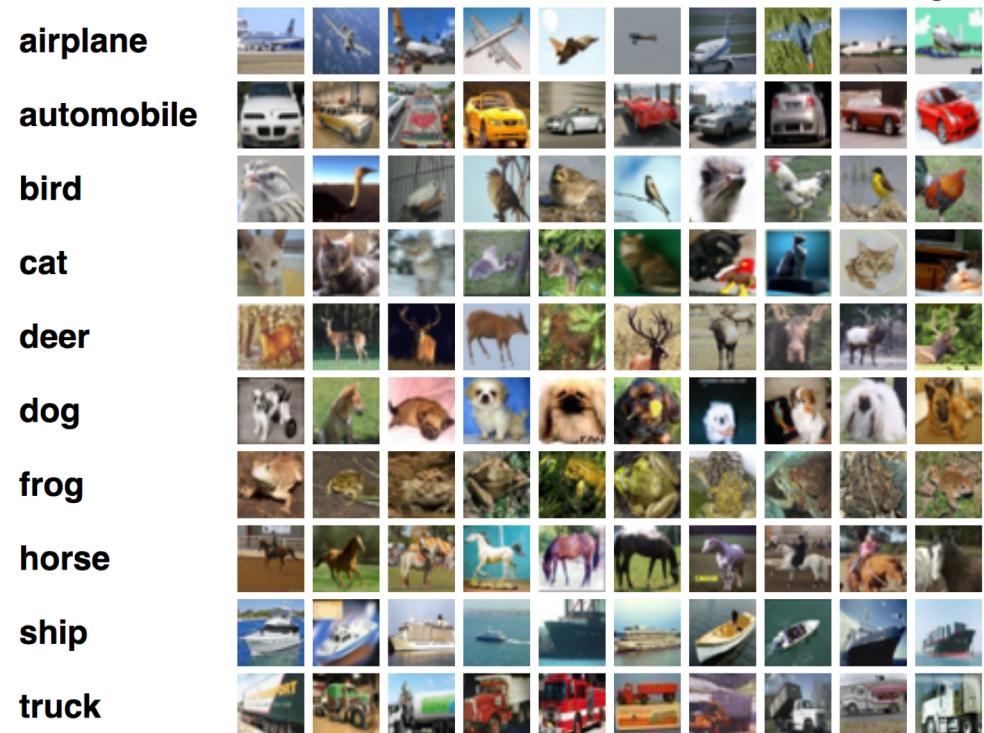
Error on test set: 2.5%

CIFAR dataset

- Number of classes: 10
- Training set: 50,000 labelled images
- Test set: 10,000 labelled images



32 x 32 RGB image



Lab 02

- CIFAR classification with MLP

The image shows two Jupyter notebook interfaces side-by-side.

Left Notebook (GPU Demo):

- Title:** Lab 02 : GPU demo
- In [1]:**

```
import torch
import torch.nn as nn
import torch.nn.functional as F
```
- Text:** Set the device variable to be either cuda or cpu
- In [9]:**

```
device = torch.device("cuda")
#device = torch.device("cpu")
print(device)
```

 Output: cpu
- Text:** Make a random tensor (by default the tensor is created on the CPU)
- In [10]:**

```
x=torch.rand(3,3)
print(x)
```

 Output:

```
tensor([[0.2126, 0.8232, 0.9466],
       [0.0370, 0.5176, 0.7032],
       [0.4544, 0.4457, 0.8515]])
```
- Text:** Send the tensor to the device (which is going to be the GPU if the device was set to be the GPU -- otherwise it will remain on the CPU)

Right Notebook (CIFAR multi-layer exercise):

- Title:** Lab 02 : CIFAR multi-layer -- exercise
- In [1]:**

```
import torch
import torch.nn as nn
import torch.nn.functional as F
import torch.optim as optim
from random import randint
import utils
import time
```
- Text:** With or without GPU?
- Text:** It is recommended to run this code on GPU:
 - Time for 1 epoch on CPU : 5 sec
 - Time for 1 epoch on GPU : 0.6 sec w/ GeForce GTX 1080 Ti
- In []:**

```
#device= torch.device("cuda")
device= torch.device("cpu")
print(device)
```
- Text:** Download the CIFAR dataset
- In [2]:**

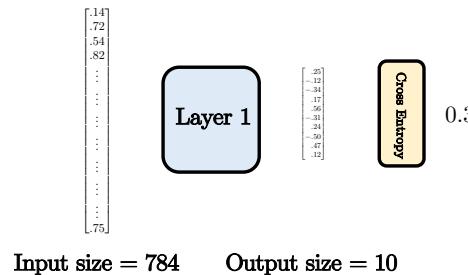
```
from utils import check_cifar_dataset_exists
data_path=check_cifar_dataset_exists()

train_data= torch.load(data_path+'cifar/train_data.pt')
train_label= torch.load(data_path+'cifar/train_label.pt')
test_data= torch.load(data_path+'cifar/test_data.pt')
test_label= torch.load(data_path+'cifar/test_label.pt')
```

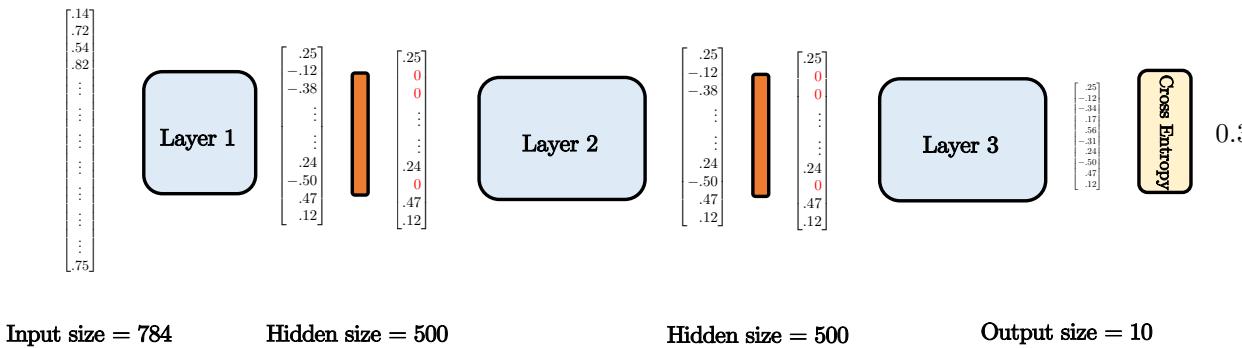
Lab 02

- Results

- One-layer network:



- Three-layer network:





Questions?