

CS5242 : Neural Networks and Deep Learning

Lecture 3: Vanilla Neural Networks Inference and Learning

Semester 1 2021/22

Xavier Bresson

<https://twitter.com/xbresson>

Department of Computer Science
National University of Singapore (NUS)



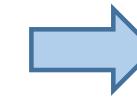
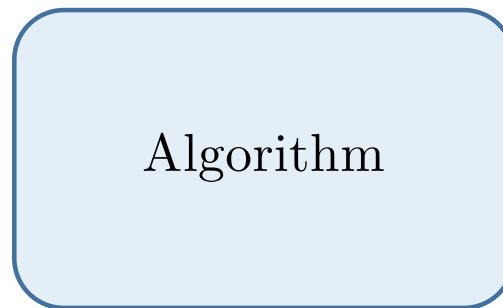
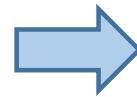
Outline

- Image classification task
- Neural networks for image classification
- Forward pass for inference
- Understanding linear layers
- Understanding the backward pass
- Matrix formulation
- Mini-batch learning

Outline

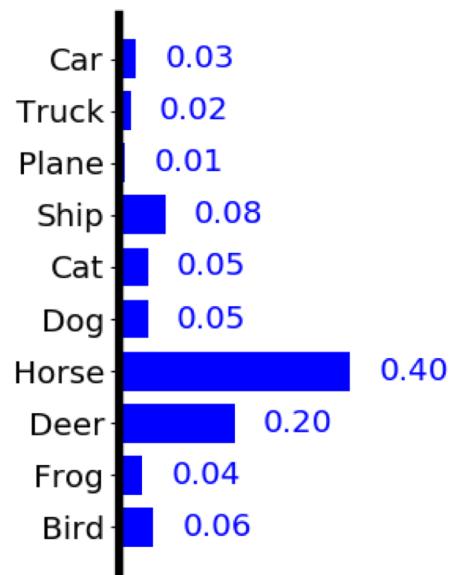
- **Image classification task**
- Neural networks for image classification
- Forward pass for inference
- Understanding linear layers
- Understanding the backward pass
- Matrix formulation
- Mini-batch learning

Classification task



Input

$3 \times 32 \times 32$ tensor



Output

Probabilities over
classes

Outline

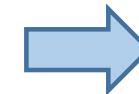
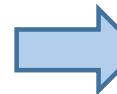
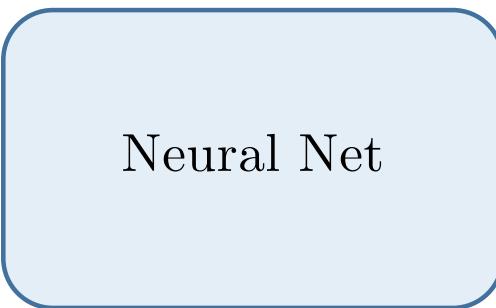
- Image classification task
- Neural networks for image classification
- Forward pass for inference
- Understanding linear layers
- Understanding the backward pass
- Matrix formulation
- Mini-batch learning

Neural network as classification algorithm

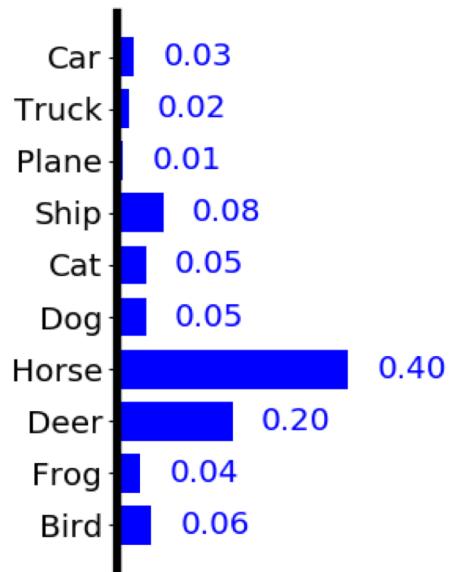


Input

$3 \times 32 \times 32$ tensor



The classification task is solved by the **forward pass**, which consists in passing the input image through a series of linear and non-linear operations to obtain the **solution**.



Output

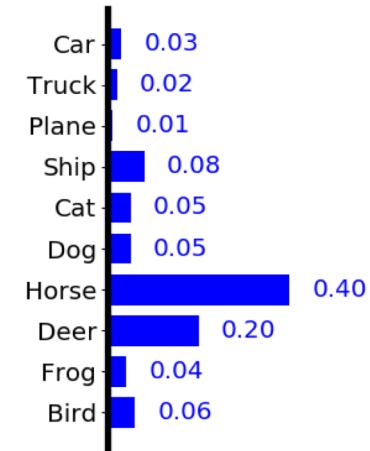
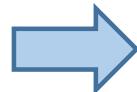
Probabilities over
classes

Outline

- Image classification task
- Neural networks for image classification
- **Forward pass for inference**
- Understanding linear layers
- Understanding the backward pass
- Matrix formulation
- Mini-batch learning

Forward pass

- Inference/prediction:
 - How to take 3072 numbers and convert them into 10 positive numbers that sum to one (probability function)?



Input

$3 \times 32 \times 32 = 3072$ numbers

Output

10 positive numbers that
sum to one

Simplest possible transformations

- Inference is done by a linear layer (LL) + a (non-linear) softmax layer.
 - This is the simplest way to take 3072 numbers and to convert them into 10 positive numbers that sum to one (probability function).



$\begin{bmatrix} .14 \\ .72 \\ .54 \\ .82 \\ \vdots \\ .75 \end{bmatrix}$

Linear
Transformation

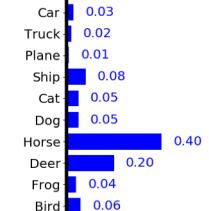


$\begin{bmatrix} -1.6 \\ -0.9 \\ -1.1 \\ -2.5 \\ 0.1 \\ 0.5 \\ 1.7 \\ 1.2 \\ -0.3 \\ 0.2 \end{bmatrix}$



Softmax
Function

$\begin{bmatrix} .01 \\ .03 \\ .02 \\ .01 \\ .08 \\ .11 \\ .38 \\ .23 \\ .05 \\ .08 \end{bmatrix}$



3072 numbers

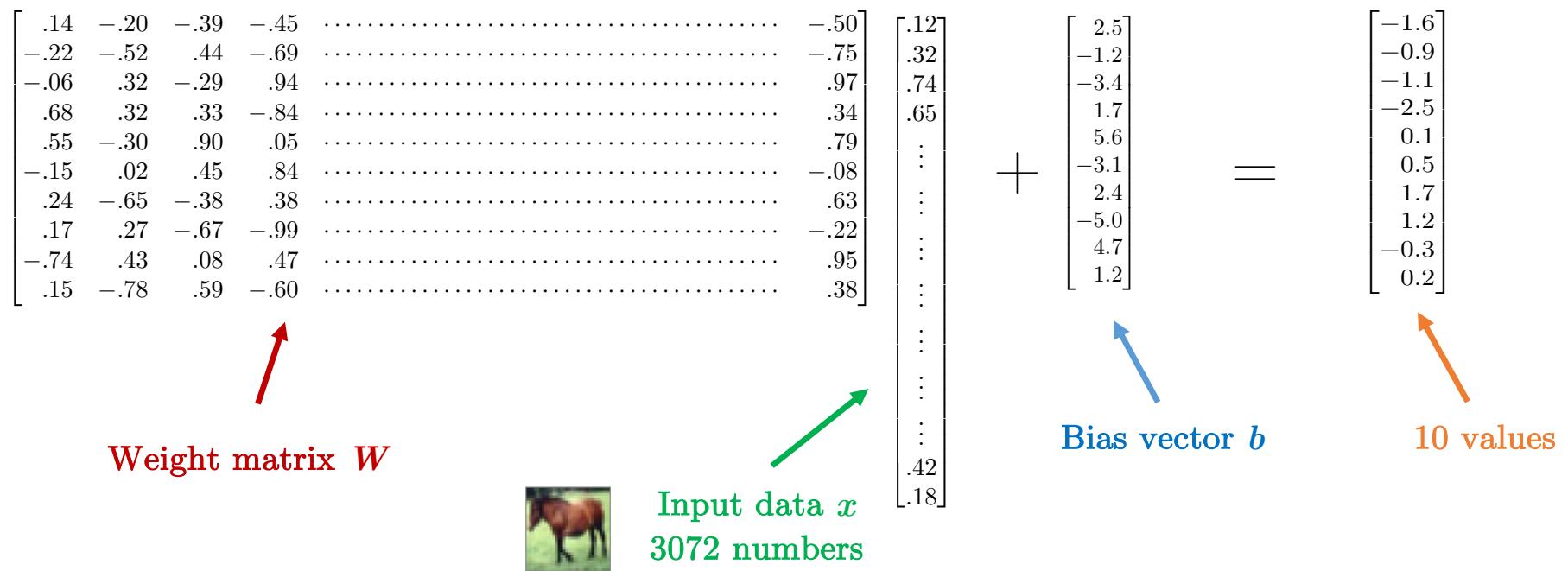
10 numbers

10 positive numbers
that sum to one

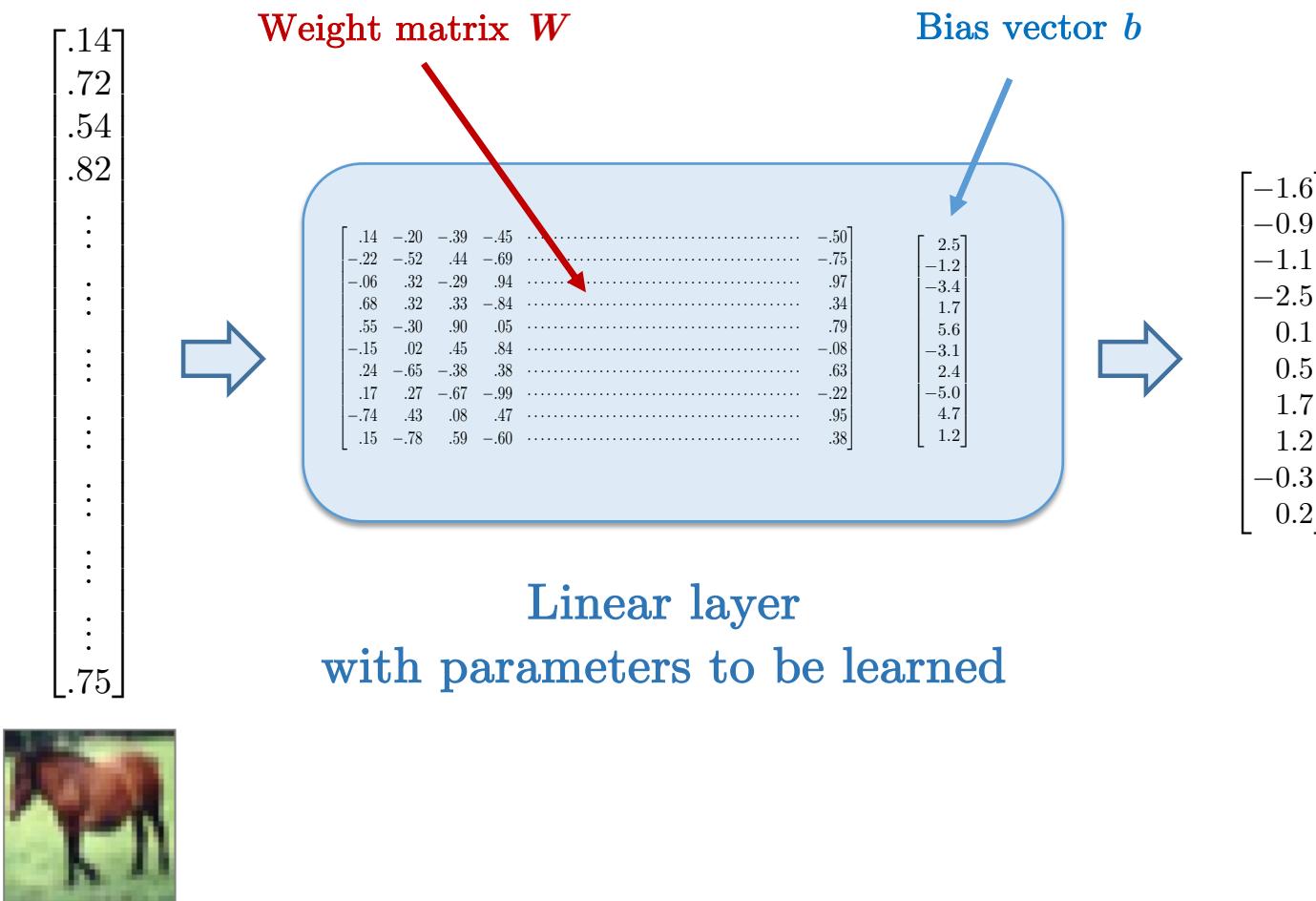
Linear layer

- Linear transformation:
 - Operation that **changes the data dimension** from 3072 to 10.
 - Linear layers have both a **weight matrix** and a **bias vector**:

$$Wx + b$$



Linear layer in PyTorch



Lab 01

- Linear layer module

The image shows two Jupyter notebook interfaces side-by-side, illustrating different approaches to creating and inspecting linear modules.

Left Notebook (demo):

- Title:** Lab 01 : Linear module -- demo
- In [1]:** `import torch
import torch.nn as nn`
- Text:** Make a *Linear Module* that takes input of size 5 and return output of size 3
- In [7]:** `mod = nn.Linear(5,3,bias=True)
print(mod)`
Output: `Linear(in_features=5, out_features=3, bias=True)
torch.Size([3, 5])
torch.Size([3])`
- Text:** Let's make a random tensor of size 5:
- In [4]:** `x=torch.rand(5)
print(x)
print(x.size())`
Output: `tensor(0.5148, 0.6442, 0.5563, 0.4040, 0.9193)
torch.Size([5])`
- Text:** Feed it to the module:

Right Notebook (exercise):

- Title:** Lab 01 : Linear module -- exercise
- In [1]:** `import torch
import torch.nn as nn`
- Text:** Make a linear module WITHOUT bias that takes input of size 2 and return output of size 3.
- In []:** `mod= # complete here`
Output: `print(mod)`
- Text:** Print the internal parameters of the module. Try to print the bias and double check that it return you "None"
- In []:** `print(# complete here)
print(# complete here)`
- Text:** Make a vector $x = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ and feed it to your linear module. What is the output?
- In []:** `x= # complete here
print(x)
y= # complete here
print(y)`

Softmax layer

- Softmax layer converts the 10 scores into a probability distribution over the 10 classes:

$$\begin{bmatrix} -1.6 \\ -0.9 \\ -1.1 \\ -2.5 \\ 0.1 \\ 0.5 \\ 1.7 \\ 1.2 \\ -0.3 \\ 0.2 \end{bmatrix}$$

`exp()`

$$\begin{bmatrix} 0.2 \\ 0.4 \\ 0.3 \\ 0.1 \\ 1.1 \\ 1.6 \\ 5.5 \\ 3.3 \\ 0.7 \\ 1.2 \end{bmatrix}$$

sum = 14.53

Scores

$$\xrightarrow{\text{divide by sum()}}$$

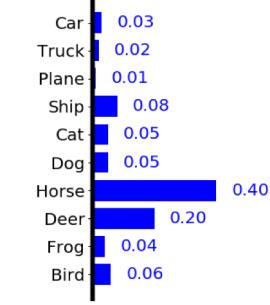
divide by sum()

$$\begin{bmatrix} .01 \\ .03 \\ .02 \\ .01 \\ .08 \\ .11 \\ .38 \\ .23 \\ .05 \\ .08 \end{bmatrix}$$

sum = 1

Positive numbers

Probabilities
over classes



Softmax layer

- These two operations (exponentiate then unit normalization) can be done into one single step by the **Softmax** function:

$$\text{Softmax} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} = \begin{pmatrix} \frac{e^{x_1}}{e^{x_1} + e^{x_2} + e^{x_3} + e^{x_4}} \\ \frac{e^{x_2}}{e^{x_1} + e^{x_2} + e^{x_3} + e^{x_4}} \\ \frac{e^{x_3}}{e^{x_1} + e^{x_2} + e^{x_3} + e^{x_4}} \\ \frac{e^{x_4}}{e^{x_1} + e^{x_2} + e^{x_3} + e^{x_4}} \end{pmatrix}$$

- It is called “Softmax” because **it detects the Max** :

Examples:

$$\text{Softmax} \begin{pmatrix} -2.0 \\ -0.5 \\ 2.0 \\ 1.5 \end{pmatrix} = \begin{pmatrix} .01 \\ .04 \\ .59 \\ .36 \end{pmatrix}$$
$$\text{Softmax} \begin{pmatrix} -2.0 \\ -0.5 \\ 7.0 \\ 1.5 \end{pmatrix} = \begin{pmatrix} .00 \\ .00 \\ .99 \\ .01 \end{pmatrix}$$

Lab 02

- Softmax layer in PyTorch

The screenshot shows a Jupyter Notebook interface with the title "Lab 02 : Softmax function -- demo". The notebook has four cells:

- In [1]:**

```
import torch
import torch.nn.functional as F
```
- Make a vector**
In [3]:

```
x=torch.Tensor([-2 , -0.5 , 2.0 , 1.5 ])
print(x)
print(x.size())
```

tensor([-2.0000, -0.5000, 2.0000, 1.5000])
torch.Size([4])
- Feed it to the softmax**
In [6]:

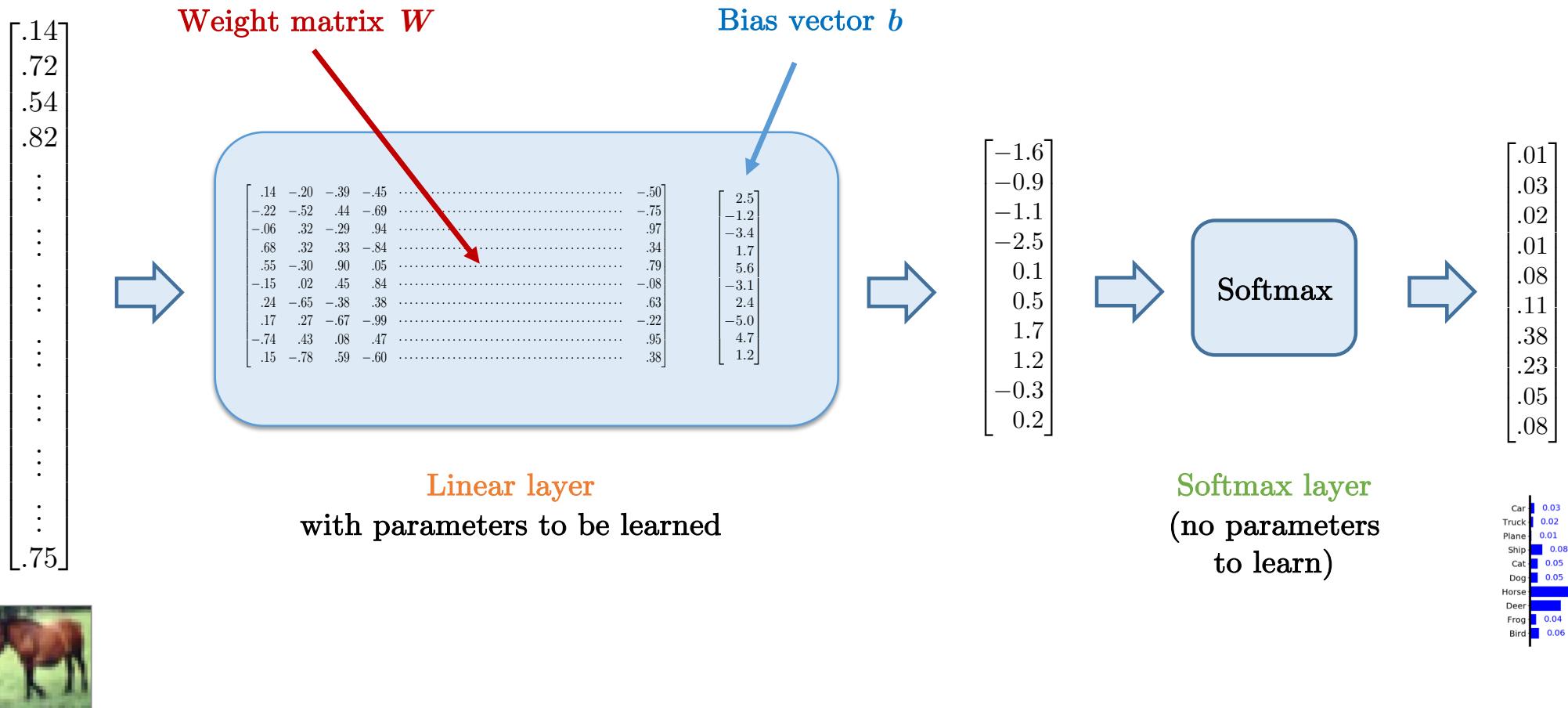
```
p=F.softmax(x , dim=0)
print(p)
print(p.size())
```

tensor([0.0107, 0.0481, 0.5858, 0.3553])
torch.Size([4])
- Check that it sums to one**
In [4]:

```
print( p.sum() )
```

tensor(1.)

One-layer neural network



Lab 03

- Design a vanilla neural network in PyTorch :

The image shows two Jupyter notebook interfaces side-by-side, illustrating the creation of neural networks in PyTorch.

Left Notebook (demo):

- Title:** Lab 03 : Vanilla neural networks -- demo
- Section:** Creating a two-layer network
- Code In [1]:**

```
import torch
import torch.nn as nn
import torch.nn.functional as F
```
- Code In [2]:**

```
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=True)
        self.layer2 = nn.Linear(hidden_size, output_size, bias=True)

    def forward(self, x):
        x = self.layer1(x)
        x = F.relu(x)
        x = self.layer2(x)
        p = F.softmax(x, dim=0)

        return p
```
- Text:** Create an instance that takes input of size 2, then transform it into something of size 5, then into something of size 3

Right Notebook (exercise):

- Title:** Lab 03 : Vanilla neural networks -- exercise
- Section:** Creating a one-layer network
- Code In [1]:**

```
import torch
import torch.nn as nn
import torch.nn.functional as F
```
- Text:** Make a class for a one layer network. Let's call the layer "mylayer". And let's give it a bias.
- Code In [2]:**

```
class one_layer_net(nn.Module):
    def __init__(self, input_size, output_size):
        super(one_layer_net, self).__init__()

        # complete here

    def forward(self, x):
        # complete here
        # complete here

        return p
```
- Text:** Create an instance of a one layer net that take input of size 2 and return output of size 2
- Code In [3]:**

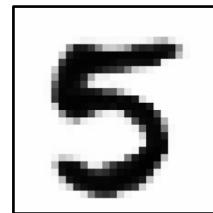
```
net = # complete here
print(net)
```

Outline

- Image classification task
- Neural networks for image classification
- Forward pass for inference
- **Understanding linear layers**
- Understanding the backward pass
- Matrix formulation
- Mini-batch learning

MNIST dataset

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

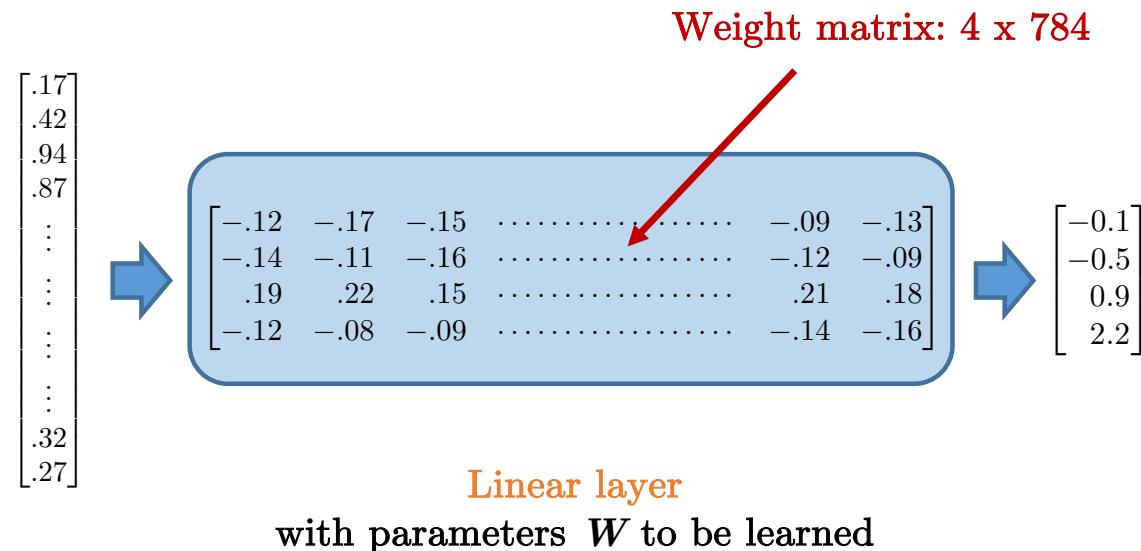


28 x 28 grayscale image

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

MNIST dataset

- For simplicity, we will pretend :
 - We have only 4 classes: Zero, One, Two, Three.
 - We have no bias b .



What linear layers do?

- A linear layer is a template matching algorithm !
- Each row of the weight matrix W is a template (encoding some data property):

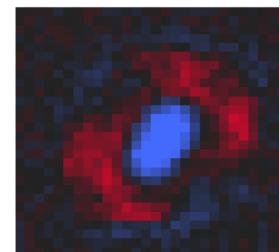
-.12	-.17	-.15	-.09	-.13
-.14	-.11	-.16	-.12	-.09
.19	.22	.1521	.18
-.12	-.08	-.09	-.14	-.16

4 x 784 weight matrix W

Row zero
reshaped as
(28 x 28)

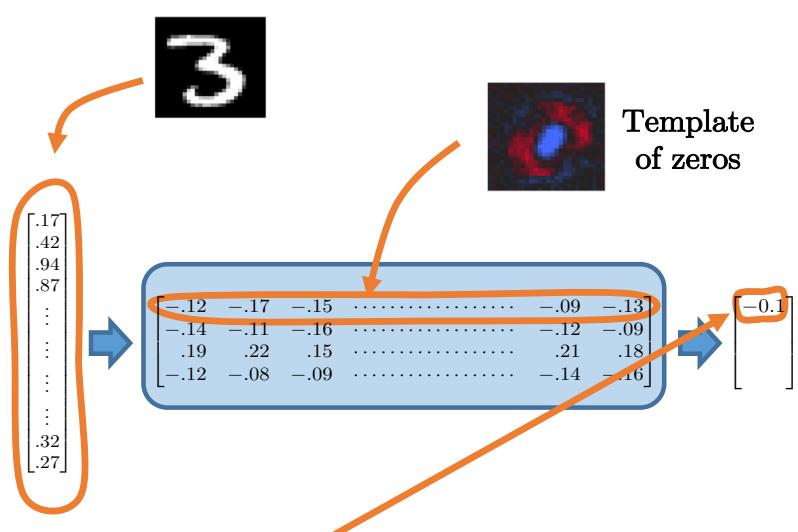
$$\begin{bmatrix} .14 & -.20 & -.39 & -.45 \\ -.94 & -.05 & .57 & .42 \\ .45 & -.46 & .97 & .08 \\ -.72 & .77 & -.81 & -.56 \\ -.34 & -.25 & -.39 & -.50 \end{bmatrix}$$

||

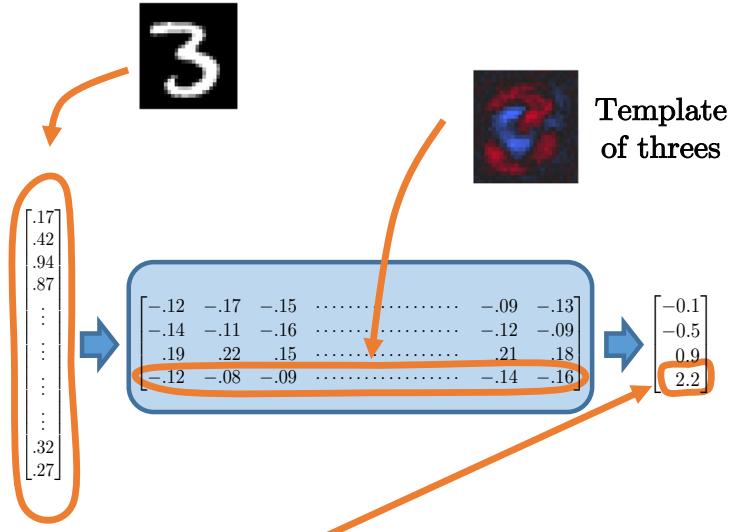


Template matching

- Template matching performs vector-vector multiplication between data and class template.
 - It produces a matching score between data and class template.



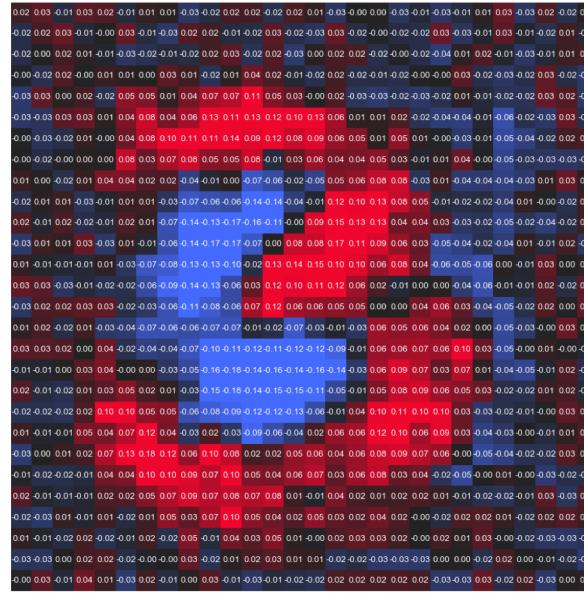
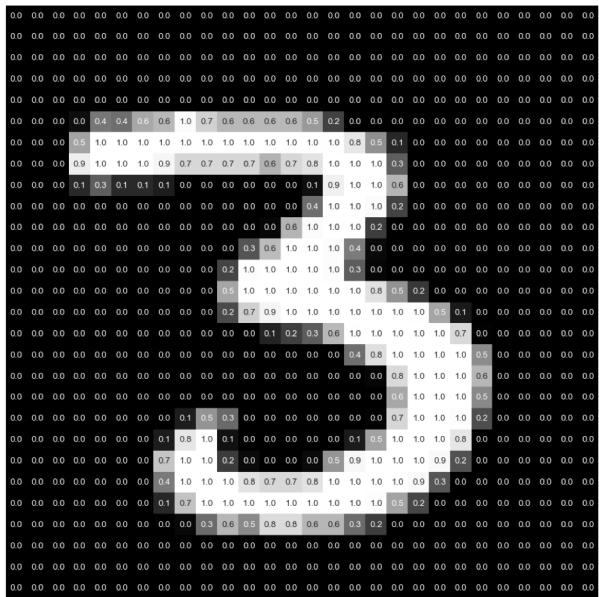
Matching score between the input data and the template of **zeros**



Matching score between the input data and the template of **threes**

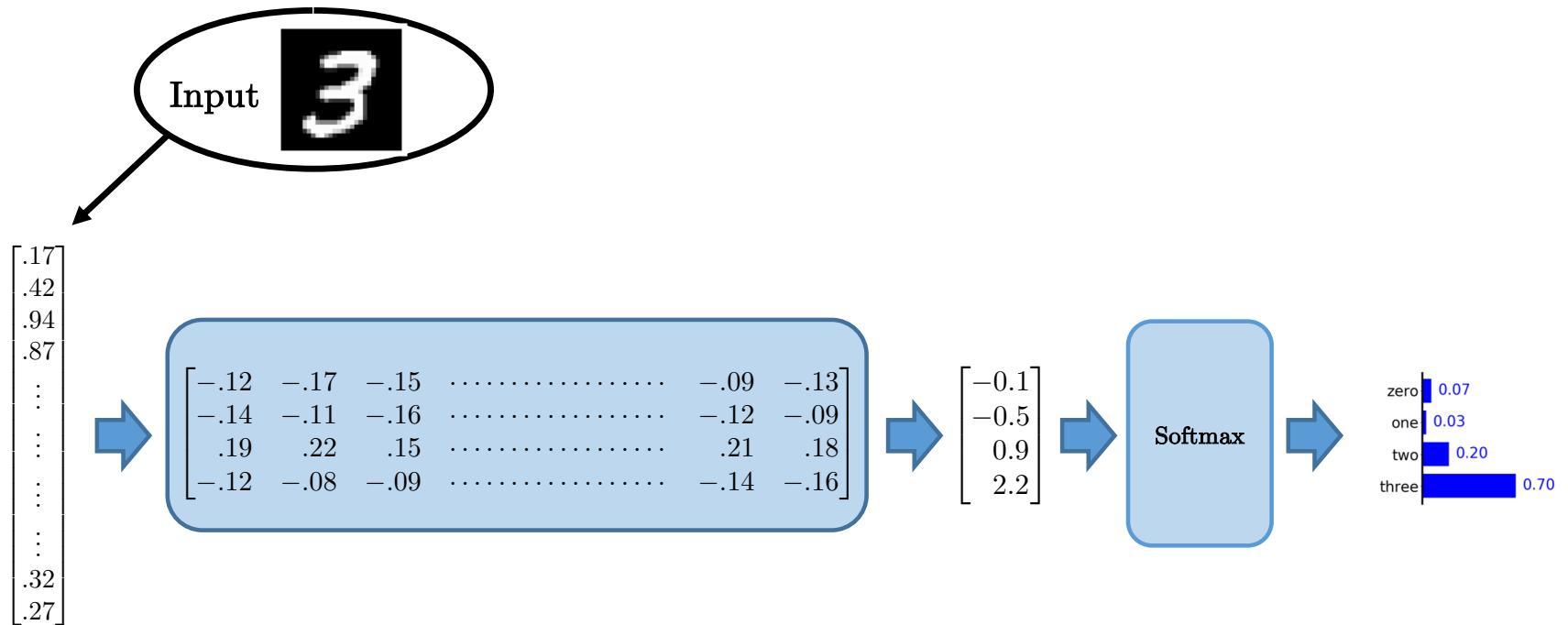
Template matching

- Matching score = inner product of input image and class template



From score to probability

- Scores are converted into probabilities with the softmax layer :



Template matching

- Template matching property of DL architectures is a **first principle**.
- However, **deep learning is more than template matching !**
- Most research focus on **finding the right architecture** to solve complex tasks (that cannot be solved by simple template matching only).

Outline

- Image classification task
- Neural networks for image classification
- Forward pass for inference
- Understanding linear layers
- **Understanding the backward pass**
- Matrix formulation
- Mini-batch learning

How to learn the network parameters?

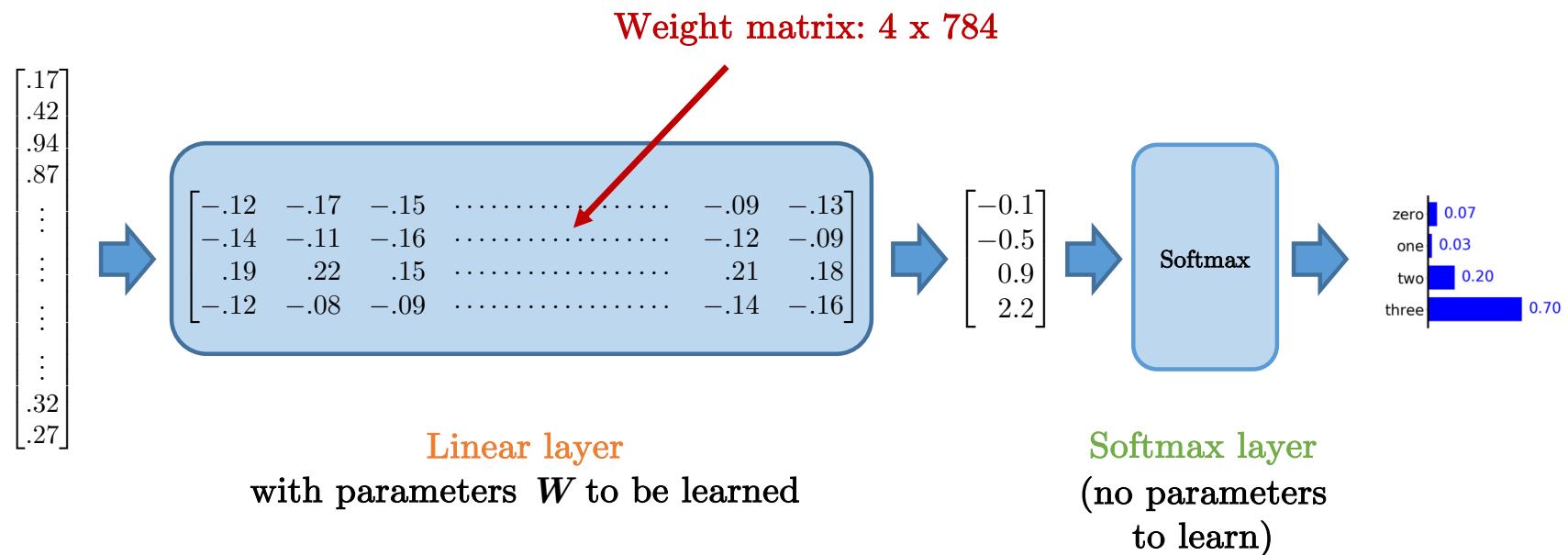
- So far, we have done **inference** : We have used the weight parameters that **best** classify images.
- We need to **learn these parameters** (here the weight and bias of the linear layer).
- Network parameters are learned by **backpropagation**:
 - A backward pass that backpropagates the error of prediction to adjust the weights on the network.
 - **The update rule of the weights is :**

template = template + lr × (error) × (input picture)

$$W = W + lr \text{ (error vector} \otimes \text{input picture})$$

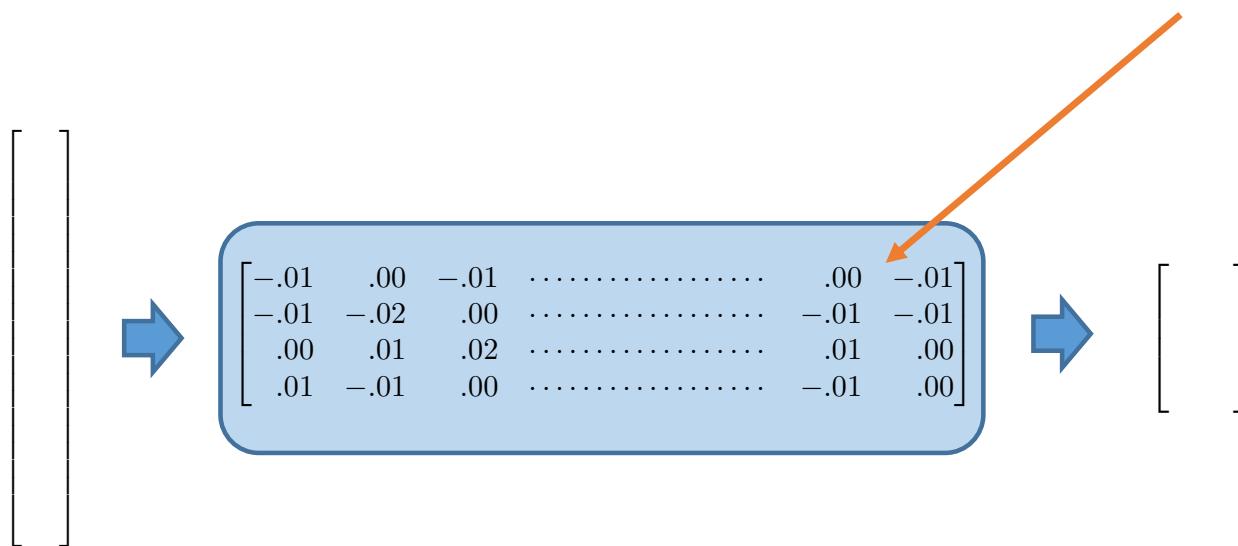
Template learning

- For simplicity, we will still pretend :
 - We have only 4 classes: Zero, One, Two, Three.
 - We have no bias b .



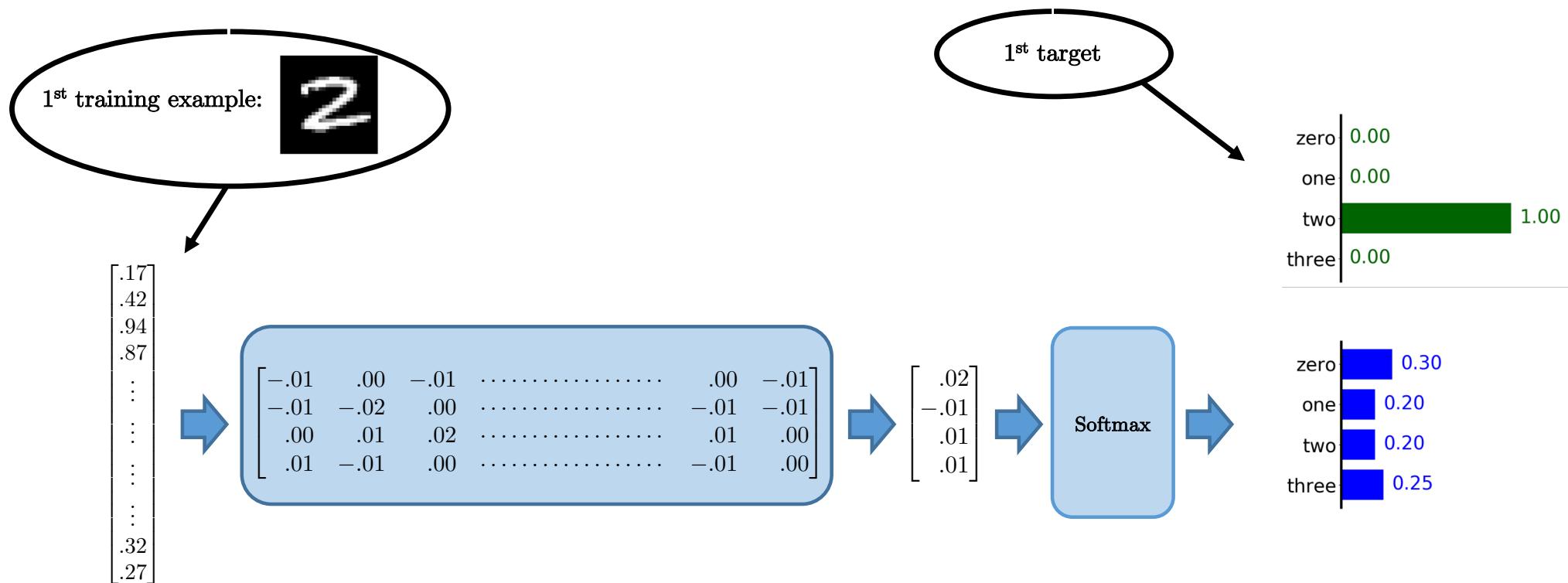
Template learning

- How does the network learn good templates?
- We will use an **iterative step** that will update progressively the templates to minimize the error of prediction.
- **Initialization:**
 - At the beginning of the training the weight matrix is initialized with **small random numbers**.



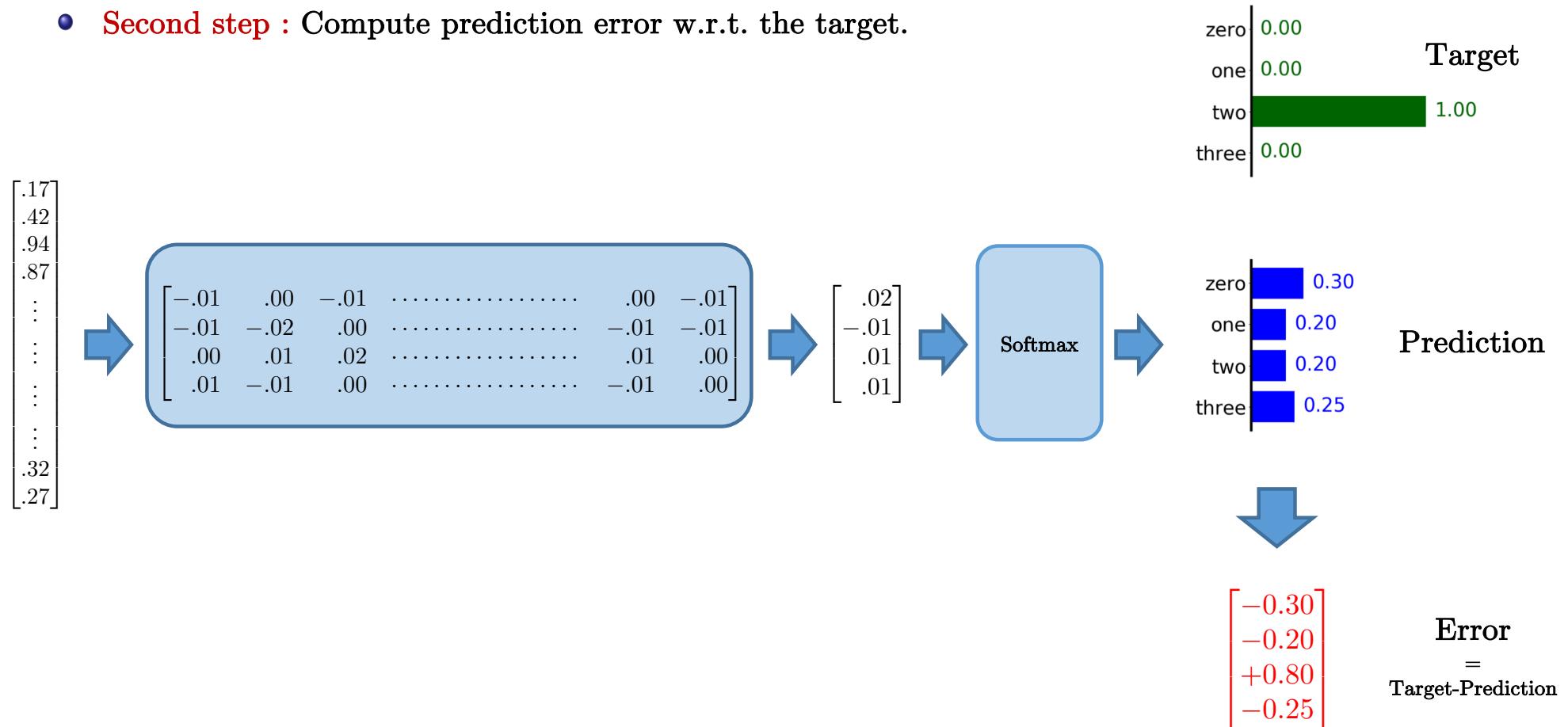
Forward pass

- First step : Compute the class probabilities for the 1st training data.



Prediction error

- Second step : Compute prediction error w.r.t. the target.



Template update

- Third step : Update templates

$$\text{template} = \text{template} + lr \times (\text{error}) \times (\text{input picture})$$

- Each time the network is presented with a picture, let say:



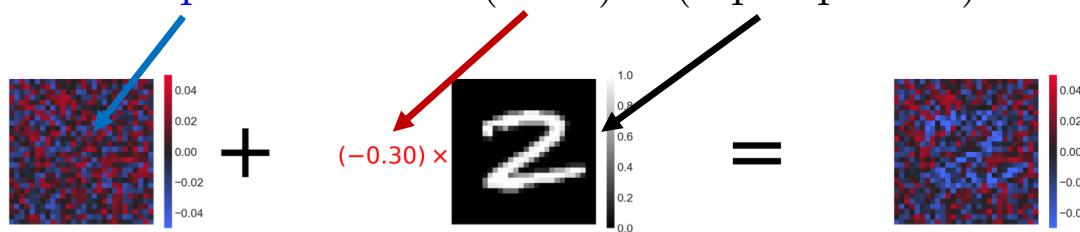
- It adds this picture to the correct class template (here class 2).
- And removes it from the other class templates (here classes 0, 1, 3).

Template update

- Applying the template update:

$$\text{template} = \text{template} + lr \times (\text{error}) \times (\text{input picture})$$

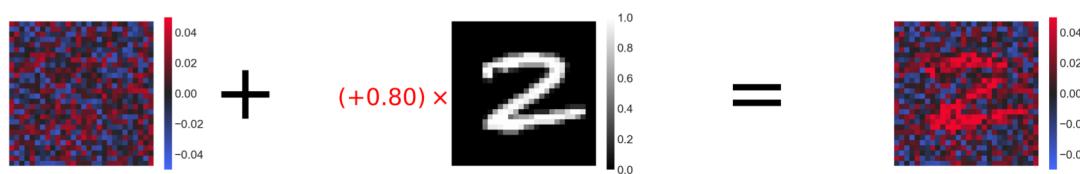
Update template 0:



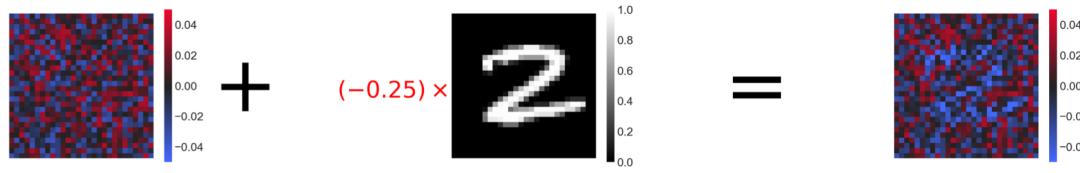
Update template 1:



Update template 2:



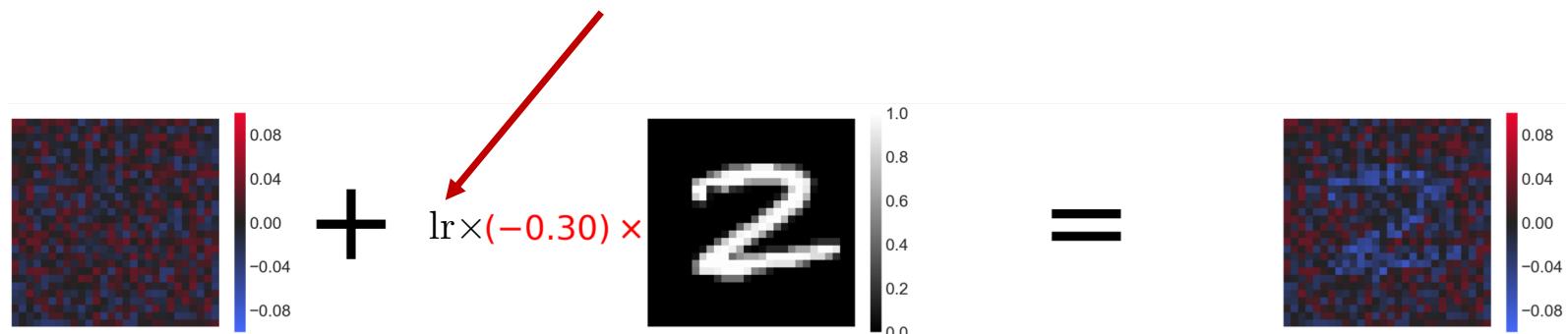
Update template 3:



Learning rate

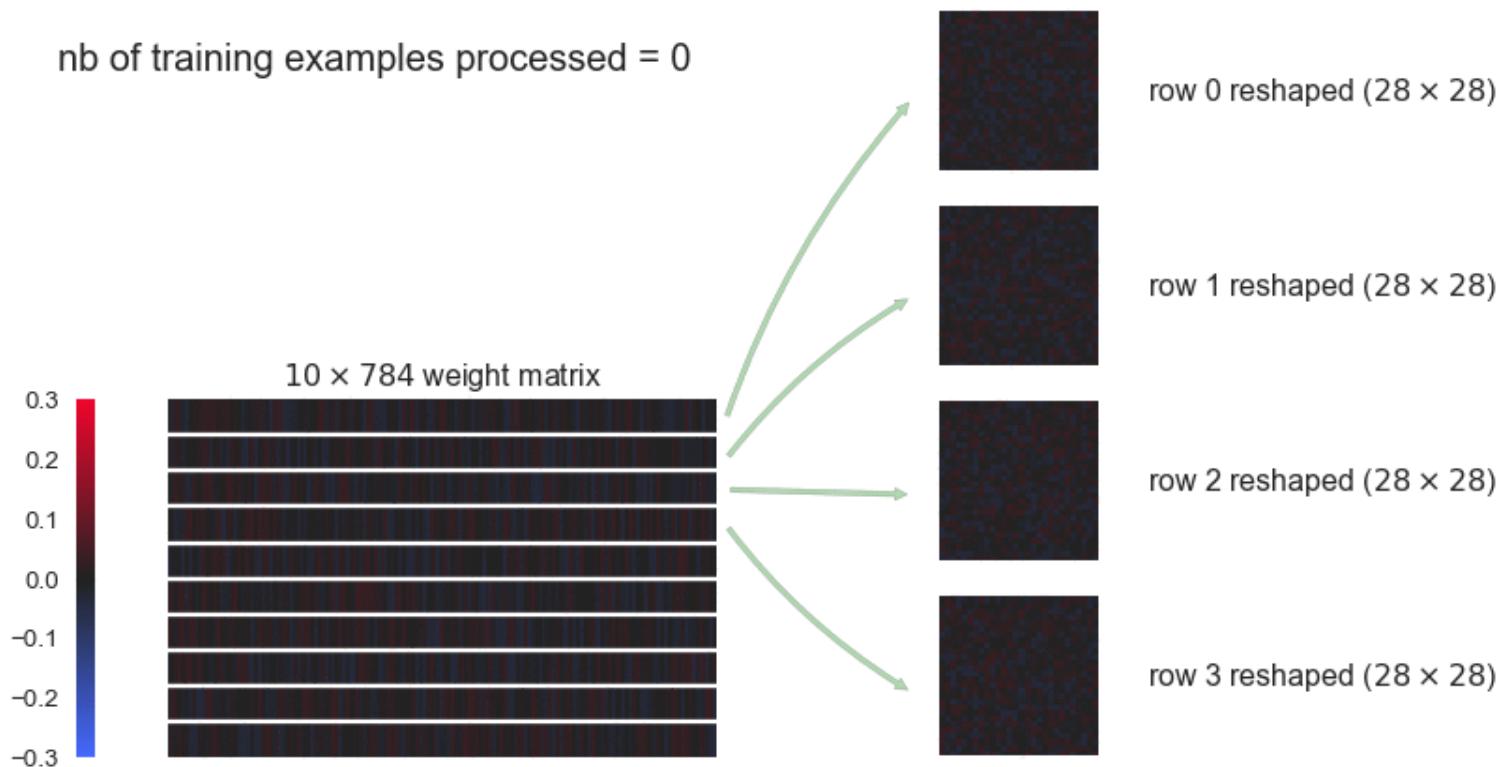
- The learning rate (lr) is a very important hyper-parameter.
- The learning rate controls the learning speed:
 - Too small : Learning is too long.
 - Too large: Learning is too fast and it does not work.
- Typical choice to update the templates is

$$\text{lr} \approx 0.01$$



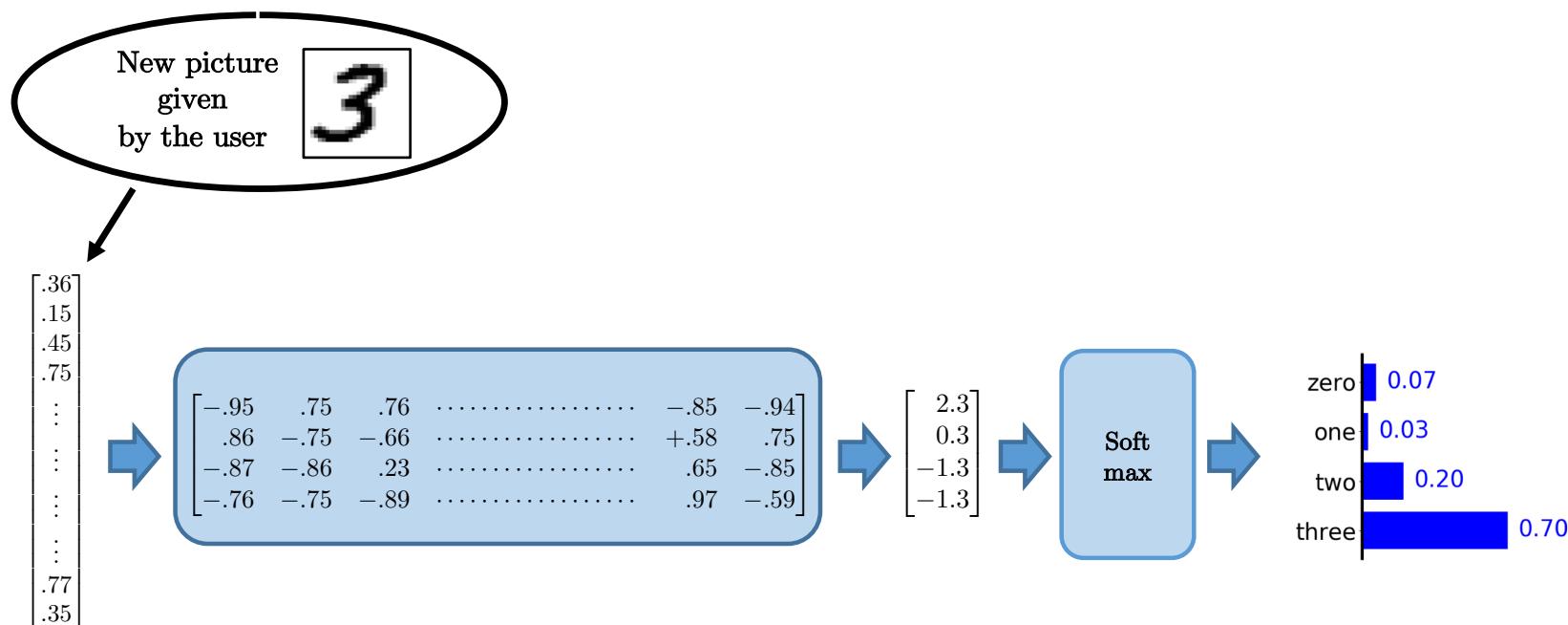
Iterate

- Iterate this process with all training data and multiple times:
 - After many iterations, the templates look good :



Network ready for inference

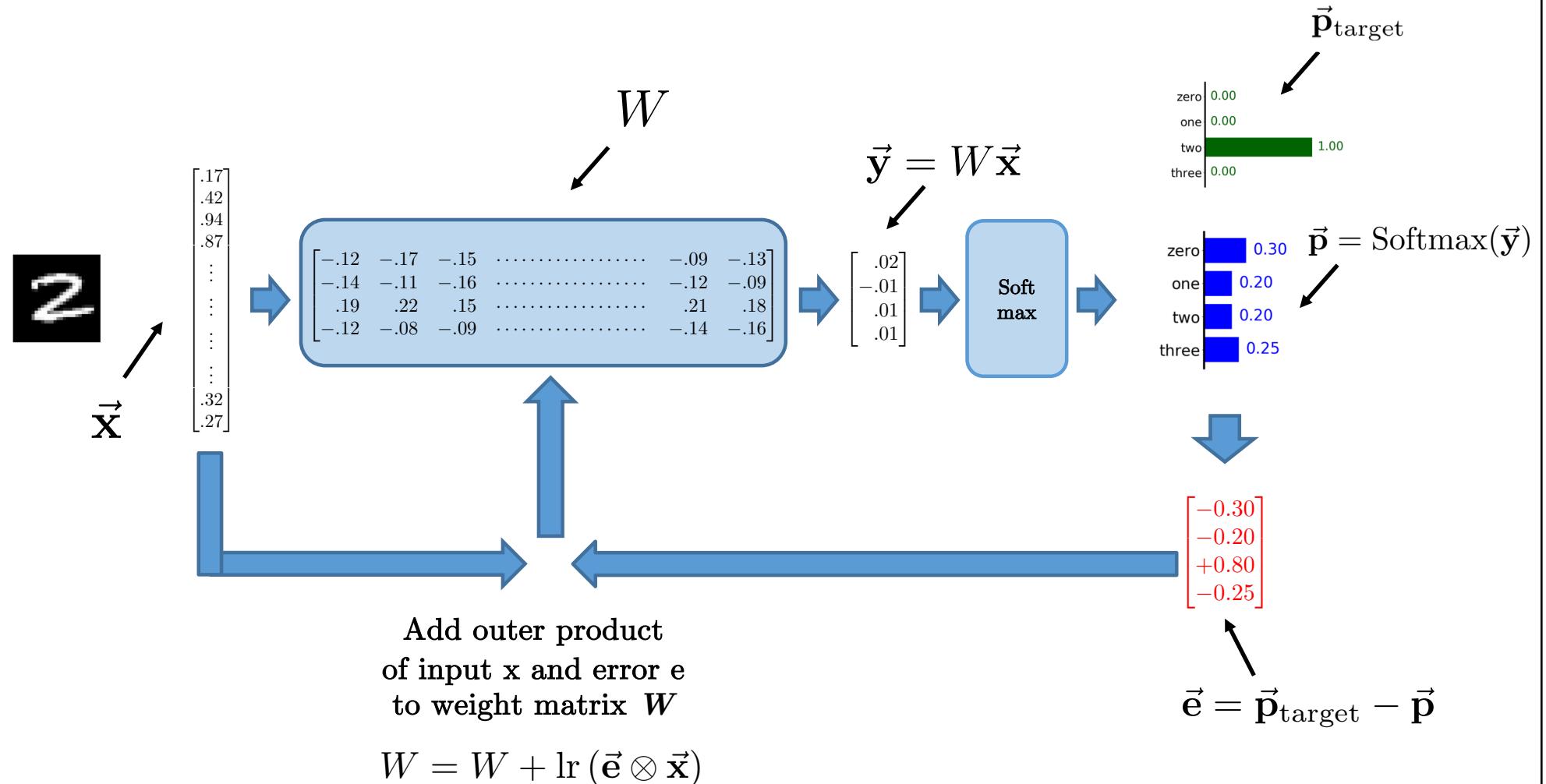
- After many iterations, the templates do not update anymore, the training is finished.
- The network is ready for inference for new data.



Outline

- Image classification task
- Neural networks for image classification
- Forward pass for inference
- Understanding linear layers
- Understanding the backward pass
- **Matrix formulation**
- Mini-batch learning

Matrix notation



Matrix formulation of outer product

$$\vec{e} \otimes (\text{input picture}) = \begin{bmatrix} -0.30 \\ -0.20 \\ +0.80 \\ -0.25 \end{bmatrix} [0.17 \quad 0.42 \quad 0.94 \quad \cdots \quad \cdots \quad 0.32 \quad 0.27]$$
$$= \begin{bmatrix} -0.30 \times \text{input picture} \\ -0.20 \times \text{input picture} \\ +0.80 \times \text{input picture} \\ -0.25 \times \text{input picture} \end{bmatrix}$$

Forward and backward pass

- Forward pass:

$$\vec{y} = W \vec{x} \quad \text{Compute the scores}$$

$$\vec{p} = \text{Softmax}(\vec{y}) \quad \text{Compute the probabilities}$$

- Backward pass:

$$\vec{e} = \vec{p}_{\text{target}} - \vec{p} \quad \text{Compute the error}$$

$$W = W + \text{lr} (\vec{e} \otimes \vec{x}) \quad \text{Update the templates}$$

Lab 04

- Train a one layer net

The image shows two Jupyter notebook interfaces side-by-side, illustrating the steps to train a one-layer neural network on different datasets.

Left Notebook (train_vanilla_nn_demo):

- Title:** Lab 04 : Train vanilla neural network -- demo
- Section:** Training a one-layer net on MNIST
- In [1]:** Python code to import PyTorch modules and define a random integer generator.
- In [2]:** Python code to generate a random integer between 5 and 10.
- In [3]:** Python code to check if the MNIST dataset exists and download it if not.
- In [4]:** Python code to load the MNIST training data and print its size.

Right Notebook (train_vanilla_nn_exercise):

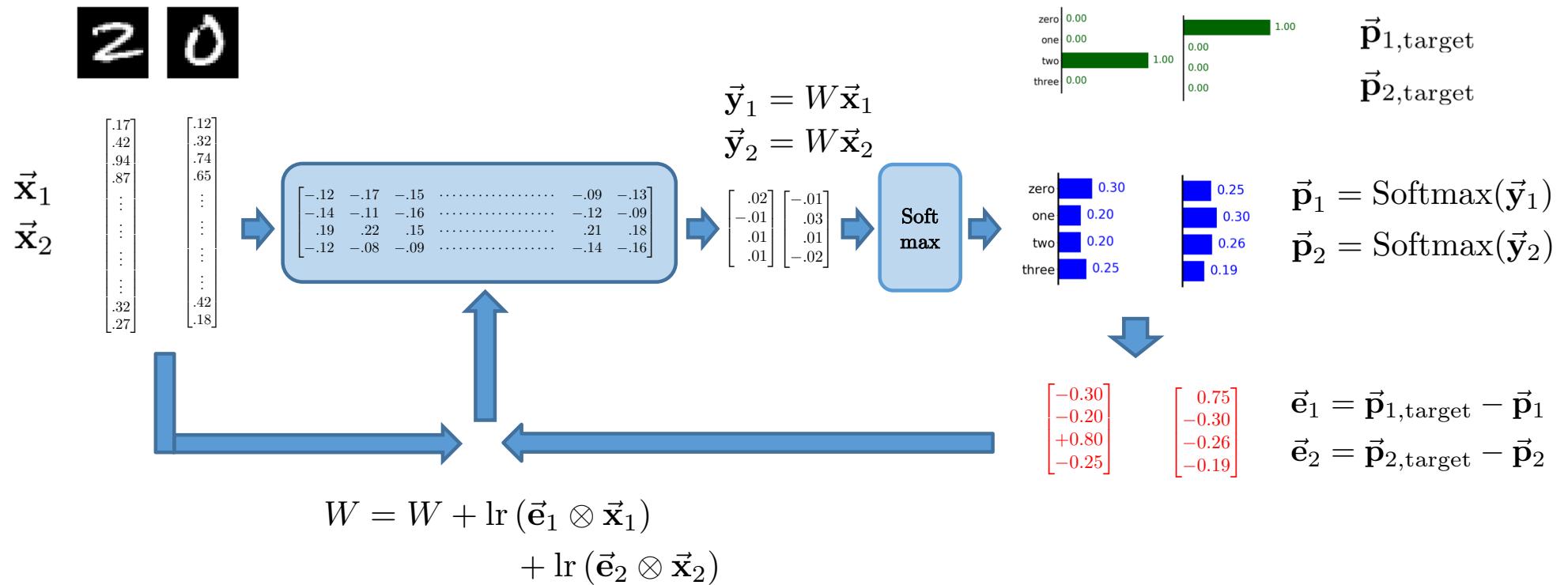
- Title:** Lab 04 : Train vanilla neural network -- exercise
- Section:** Training a one-layer net on FASHION-MNIST
- In [1]:** Python code to import PyTorch modules and define a random integer generator.
- In [2]:** Python code to check if the FASHION-MNIST dataset exists and download it if not.
- In [3]:** Python code to load the FASHION-MNIST training data and print its size.
- In [4]:** Python code to load the FASHION-MNIST test data and print its size.

Outline

- Image classification task
- Neural networks for image classification
- Forward pass for inference
- Understanding linear layers
- Understanding the backward pass
- Matrix formulation
- **Mini-batch learning**

Mini-batch

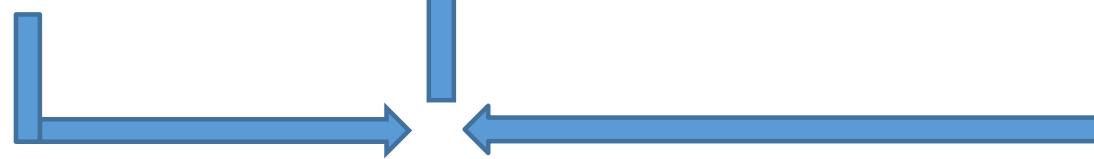
- Rather than processing one image data at a time, let us process **multiple images** at a time :



Mini-batch matrix formulation

$$\begin{matrix} \text{z} & \text{o} \\ \vec{x}_1 & \vec{x}_2 \\ \vdots & \vdots \\ X & \end{matrix}$$

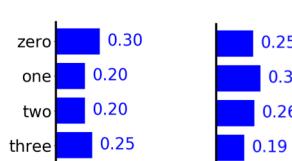
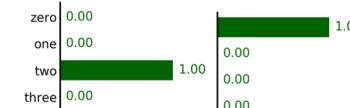
$$\begin{bmatrix} .17 \\ .42 \\ .94 \\ .87 \\ \vdots \\ \vdots \\ .32 \\ .27 \end{bmatrix} \quad \begin{bmatrix} .12 \\ .32 \\ .74 \\ .65 \\ \vdots \\ \vdots \\ .42 \\ .18 \end{bmatrix}$$



$$W = W + lr (\vec{e}_1 \otimes \vec{x}_1 + \vec{e}_2 \otimes \vec{x}_2)$$

$$W = W + lr (EX^T)$$

Sum of outer product can
be written as a mat-mat
multiplication.



$$\vec{p}_{1,\text{target}}$$

$$\vec{p}_{2,\text{target}}$$



$$P_{\text{target}}$$

$$\vec{p}_1 = \text{Softmax}(\vec{y}_1)$$

$$\vec{p}_2 = \text{Softmax}(\vec{y}_2)$$



$$P = \text{Softmax}(Y)$$

$$\vec{e}_1 = \vec{p}_{1,\text{target}} - \vec{p}_1$$

$$\vec{e}_2 = \vec{p}_{2,\text{target}} - \vec{p}_2$$



$$E = P_{\text{target}} - P$$

Matrix formulation of outer product

$$\vec{e}_1 \otimes \vec{x}_1 + \vec{e}_2 \otimes \vec{x}_2 = \begin{bmatrix} -0.30 \\ -0.20 \\ +0.80 \\ -0.25 \end{bmatrix} \begin{bmatrix} .17 & .42 & .94 & .87 & \cdots & \cdots & \cdots & \cdots & .32 & .27 \end{bmatrix} + \begin{bmatrix} +0.75 \\ -0.30 \\ -0.26 \\ -0.19 \end{bmatrix} \begin{bmatrix} .12 & .32 & .74 & .65 & \cdots & \cdots & \cdots & \cdots & .42 & .18 \end{bmatrix}$$



$$EX^T = \begin{bmatrix} -0.30 & +0.75 \\ -0.20 & -0.30 \\ +0.80 & -0.26 \\ -0.25 & -0.19 \end{bmatrix} \begin{bmatrix} .17 & .42 & .94 & .87 & \cdots & \cdots & \cdots & \cdots & .32 & .27 \\ .12 & .32 & .74 & .65 & \cdots & \cdots & \cdots & \cdots & .42 & .18 \end{bmatrix}$$

Forward and backward pass for mini-batch

- Forward pass:

$$Y = WX \quad \text{Compute the scores}$$

$$P = \text{Softmax}(Y) \quad \text{Compute the probabilities}$$

- Backward pass:

$$E = P_{\text{target}} - P \quad \text{Compute the error}$$

$$W = W + \text{lr} (EX^T) \quad \text{Update the templates}$$

Serial vs Batch

- What is the difference between serial and batch learning?
 - Serial learning :



is first processed by

$$\begin{bmatrix} -.01 & .00 & -.01 & \dots & .00 & -.01 \\ -.01 & -.02 & .00 & \dots & -.01 & -.01 \\ .00 & .01 & .02 & \dots & .01 & .00 \\ .01 & -.01 & .00 & \dots & -.01 & .00 \end{bmatrix}$$

The errors is computed and then the weights are updated.

Then



is processed by

$$\begin{bmatrix} -.06 & .03 & .03 & \dots & .00 & -.01 \\ .08 & .03 & -.04 & \dots & -.01 & -.01 \\ -.04 & -.09 & .05 & \dots & .01 & .15 \\ .12 & -.05 & .03 & \dots & -.01 & .18 \end{bmatrix}$$

The errors is computed and then the weight are updated.

- Serial mode : The weights have been updated twice.
 - When the second image is processed, it benefits of better freshly updated weights.
As a consequence the network trains fast.

Serial vs Batch

- Batch learning :



are both processed by



$$\begin{bmatrix} -.01 & .00 & -.01 & \dots & .00 & -.01 \\ -.01 & -.02 & .00 & \dots & -.01 & -.01 \\ .00 & .01 & .02 & \dots & .01 & .00 \\ .01 & -.01 & .00 & \dots & -.01 & .00 \end{bmatrix}$$

The errors are computed and then the weight are updated.

Batch mode : The weights have been updated only once.

Serial vs Batch

- **No batch (single data):** So processing one image after the other, and updated the weights each time, seems clearly better.

Batch size = 1

- **Mini-batch:** On the other hand, processing images by batches of e.g. 200, and updating the weights only every 200 images, is also an attractive option for **GPUs**.

Batch size = 200

- **Full batch:** What about doing batch of 60,000 (all training set) and processing everything in parallel? Not good:
 - Slow update: updating the weights only every 60,000 images.
 - Overfitting issue

Batch size = 60,000

Best choice!

Lab 05

- Train a net with minibatch

jupyter minibatch_training_demo Last Checkpoint: 3 minutes ago (autosaved)

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 utils
```

Warm up

Make a random matrix with 5 rows

In [2]:

```
A=torch.rand(5,3)
print(A)

tensor([[0.2976, 0.5865, 0.8216],
       [0.7939, 0.6613, 0.5833],
       [0.7961, 0.5487, 0.4077],
       [0.2582, 0.0407, 0.4789],
       [0.0092, 0.7814, 0.7776]])
```

Choose at random two indices in {0,1,2,3,4}

In [3]:

```
indices=torch.LongTensor(2).random_(0,5) # generate an integer randomly from 0 to 5-1.
print(indices)

tensor([3, 2])
```

jupyter minibatch_training_exercise Last Checkpoint: 3 minutes ago (autosaved)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3

In []:

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

Lab 05 : Train with mini-batches -- exercise

Download the data and print the sizes

In []:

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

In []:

```
train_data=torch.load(data_path+'fashion-mnist/train_data.pt')
print(train_data.size())
```

In []:

```
train_label=torch.load(data_path+'fashion-mnist/train_label.pt')
print(train_label.size())
```

In []:

```
test_data=torch.load(data_path+'fashion-mnist/test_data.pt')
print(test_data.size())
```

Make a one layer net class



Questions?