



Laboratory Lecture for the “**Generative and Deep Learning**” course of the
Master’s Degree in Computer Science (2025–2026)

Riccardo **Massidda**
riccardo.massidda@di.unipi.it

Why PyTorch?

Tensor Manipulation.

Tensor operations on a MATLAB/NumPy-like API.

Accelerator Support.

Seamless execution on CPU, GPU, and TPU devices.

Automatic Differentiation.


Only need to define forward computation → chain rule! 🧠

High-Level API.

Readily available neural networks layers, losses, optimizers, ...

Getting Started

For **this lecture**:

1. Clone the repository [di-unipi/gdl-lab](#) from GitHub,
2. Install PyTorch, either using an environment manager (conda, pipenv, poetry, etc.) or using Docker/**Podman** .

In a hurry? Just open the repository in [Google Colab](#)!

Up-to-date instructions to install **PyTorch** here: [Start Locally | PyTorch](#)

Basics of Tensor Operations and Manipulation

Tensors

Tensors are the main data structure and represent **multidimensional arrays**.

As for NumPy arrays, they support advanced **indexing** and **broadcasting**.

Attributes:

- **dtype**: determine the type of the tensor elements (float{16, 32, 64}, int{8, 16, 32, 64}, uint8). Can be specified during the initialization.
- **device**: memory location, as in CPU or GPU
- **layout**: dense tensors (strided) or sparse (sparse_coo)

Tensor Initialization

- **Existing Array:** `torch.tensor(list)`
- **Constants:** `torch.zeros(*dims)`, `torch.ones(*dims)`
- **Random:** `torch.randn(*dims)`, `torch.rand(*dims)`
- **Range:** `torch.linspace(start, end, steps=100)`
- **NumPy:** `torch.from_numpy(arr)`

Tensor Operations

Some operators are **overloaded**:

- `+`, `-` for addition and subtraction (support broadcasting)
- `*` is the elementwise multiplication (not the matrix product, supports broadcasting)
- `@` for matrix multiplication (`torch.matmul`)

In-place operations are defined with a suffix underscore:

- `add_`, `sub_` are the in-place equivalent for the previous operators, and also support broadcasting.

Check the documentation: <http://pytorch.org/docs/stable/tensors.html>

Broadcasting Rules

PyTorch broadcasting semantics follows NumPy own semantics.

Two tensors are “**broadcastable**” if the following rules hold:

1. Each tensor has **at least one dimension**.
2. When iterating over the dimension sizes, starting at the trailing dimension, the dimension sizes must either be **equal**, one of them is **1**, or one of them **does not exist**.

Broadcasting Rules

```
>>> x=torch.empty(5,7,3)
>>> y=torch.empty(5,7,3)
# same shapes are always broadcastable (i.e. the above rules always hold)

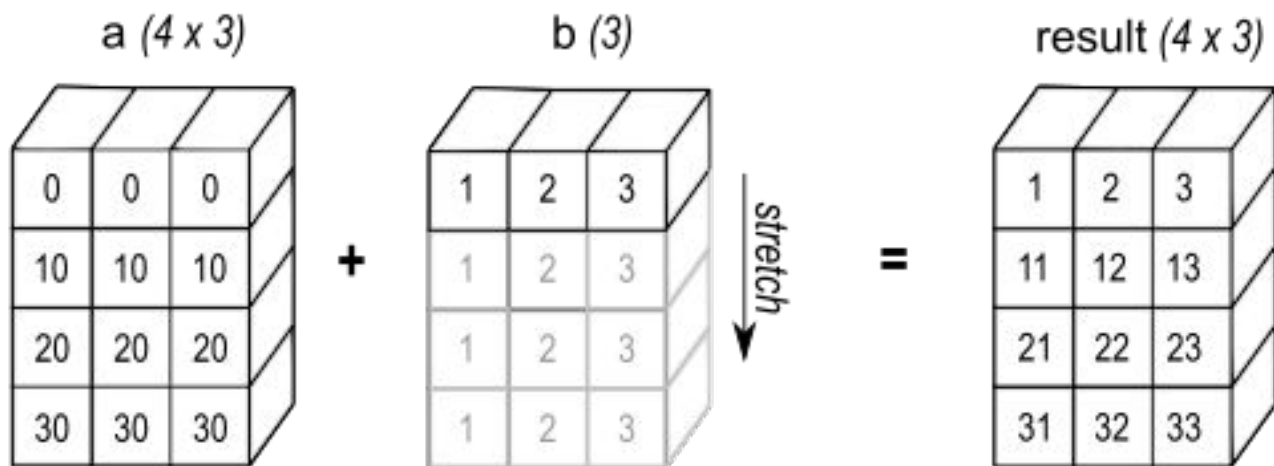
>>> x=torch.empty((0,))
>>> y=torch.empty(2,2)
# x and y are not broadcastable, because x does not have at least 1 dimension

# can line up trailing dimensions
>>> x=torch.empty(5,3,4,1)
>>> y=torch.empty( 3,1,1)
# x and y are broadcastable.
# 1st trailing dimension: both have size 1
# 2nd trailing dimension: y has size 1
# 3rd trailing dimension: x size == y size
# 4th trailing dimension: y dimension doesn't exist

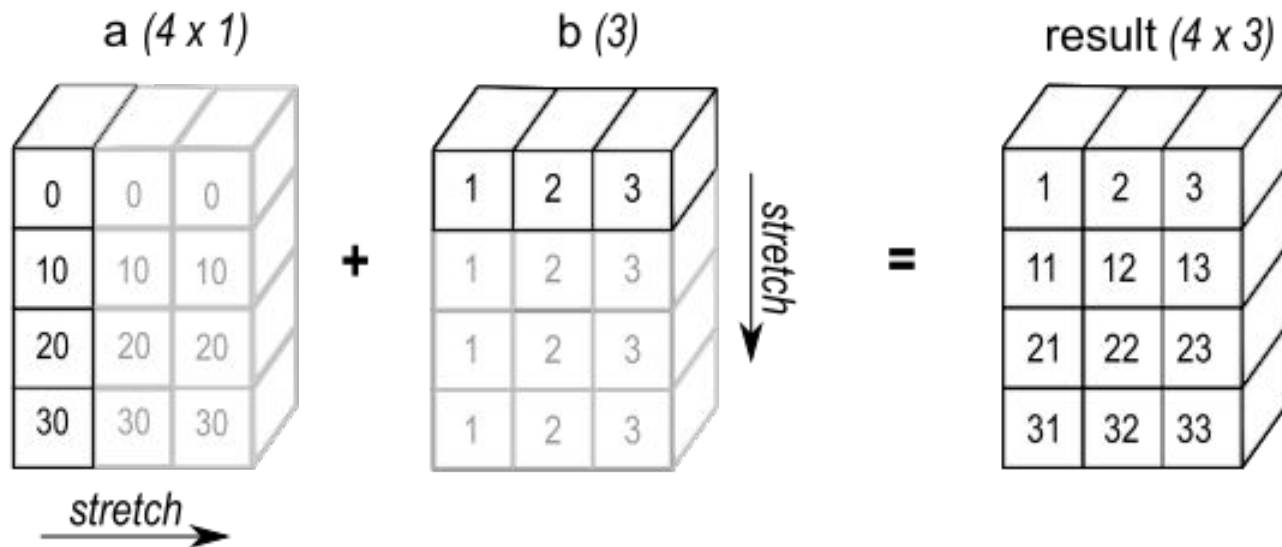
# but:
>>> x=torch.empty(5,2,4,1)
>>> y=torch.empty( 3,1,1)
# x and y are not broadcastable, because in the 3rd trailing dimension 2 != 3
```



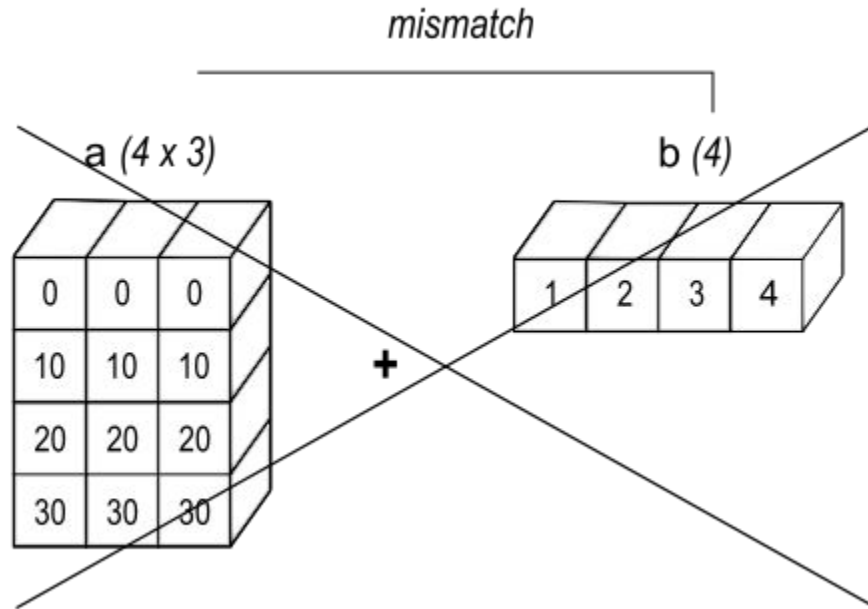
Broadcasting Rules



Broadcasting Rules



Broadcasting Rules



Tensors in GPU

The **submodule** `torch.cuda` provides the API for **GPU management**.

Check availability of the GPU:

```
torch.cuda.is_available()
```

Create or move to GPU:

```
torch.tensor([2., -1.], device="cuda")  
tensor.to("cuda")
```

In all operations, all the tensor must **reside on the same device** and result on the same device.

You can move tensors back to the CPU with the `tensor.cpu()` method.

Tensors in GPU

On a server, you typically have access to **multiple shared GPU** and you must select one:

1. **Manually** selecting with the device argument ('cuda:0', 'cuda:1'...), or
2. Using the **context manager** `torch.cuda.device`

Changing the shell environment variable `CUDA_VISIBLE_DEVICES` to limit the visible GPUs

```
export CUDA_VISIBLE_DEVICES=0
```

Note that the indices of the GPU IDs will always start from 0.

⚠ Remember to de-allocate tensors from the GPU if you're not using it!

Tensors on Apple Silicon

PyTorch supports the **MPS** (Metal Performance Shaders) backend for **Apple Silicon** (M1, ..., M4).

⚠ MPS support is still maturing: not all operations are implemented, and numerical results may differ slightly from CUDA/CPU.

```
torch.backends.mps.is_available()  
  
tensor.to("mps")  
  
torch.tensor(  
    [2., -1.], device="mps")
```

Tensor Indexing

Basic tensor **indexing** is similar to list indexing, but with multiple dimensions.

Boolean arrays can be used to filter elements that satisfy some condition.

If the indices are less than the number of dimensions, the **missing indices** are considered **complete slices**.

```
# first k elements
x = arr[:k]
# all but the first k
x = arr[k:]
# negative indexing
x = arr[-k:]
# mixed indexing
arr[:t_max, b:b+k, :]

# indexing with Boolean condition
def relu(x):
    x[x < 0] = 0
    return x
```

Tensor Reshaping

Reshaping is fundamental to combine tensors.

tensor.**squeeze**() removes all singleton dimensions

tensor.**unsqueeze**(dim) add a singleton dimension at the provided dimension

tensor.**transpose**(dim1, dim2) transposes the two dimensions of the tensor

tensor.**permute**(*dims) re-arranges the dimensions as in *dims

```
x = torch.randn(5,1,5)
```

```
x.squeeze() → [5,5]
```

```
x.unsqueeze(3) → [5,1,5,1]
```

```
x.transpose(1, 2) → [5,5,1]
```

```
x.permute(1,0,2) → [1,5,5]
```

Tensor Reshaping

`tensor.view(*new_shape)` returns a tensor with the same data but a different shape. It requires the tensor to be contiguous in memory.

`tensor.reshape(*new_shape)` does the same, but works even on non-contiguous tensors (may copy data if needed).

⚠ Use `reshape` unless you specifically want the error as a safety check that no copy is happening.

```
x = torch.randn(3, 4)
# OK: x is contiguous
x.view(12)
# OK: same data, new shape
x.view(2, 6)
# transpose → non-contiguous
y = x.t()
# RuntimeError!
y.view(12)
# OK (silently copies)
y.reshape(12)
# OK (explicit copy)
y.contiguous().view(12)
```

Tensor Reduce

Reduction operations collapse the tensor dimensionality.

tensor.**sum**(dim)

tensor.**mean**(dim)

tensor.**prod**(dim)

tensor.**amin**(dim)

tensor.**amax**(dim)

The `keepdim` parameter keeps an empty dimension in place.

```
x = torch.randn(5,1,5)
```

```
x.sum(0) → [1, 5]
```

```
x.mean(1) → [5, 5]
```

```
x.amin(2) → [5,1]
```

Your Turn!

The **Kaiming** uniform initialization scheme provides a standard baseline to train Neural Networks with rectified activation functions.

Write the following functions:

`_relu_kaiming_init_(weights: torch.Tensor)`
that modifies **in-place** the provided tensor,

`_relu_kaiming_init(in_size, out_size)` that
returns a new tensor with shape $(\text{out_size} \times \text{in_size})$

$$\mathcal{U} \left(-\sqrt{\frac{6}{\text{in_size}}}, \sqrt{\frac{6}{\text{in_size}}} \right)$$

Autograd

**Automatic
Differentiation
in PyTorch**

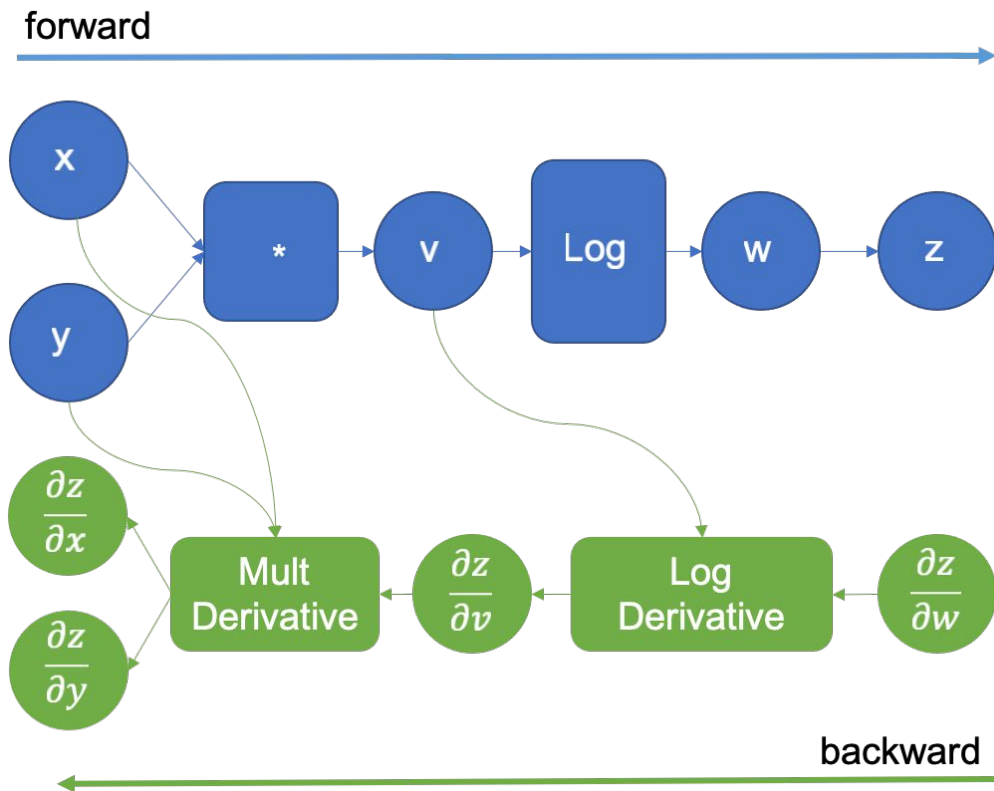
Autograd

The **submodule** `torch.autograd` is responsible for **automatic differentiation**.

Each operation creates a `Function` node in a dynamic computational graph, connected to its `Tensor` arguments.

The gradient is computed on each tensor by calling the `backward()` method.

Autograd



Autograd

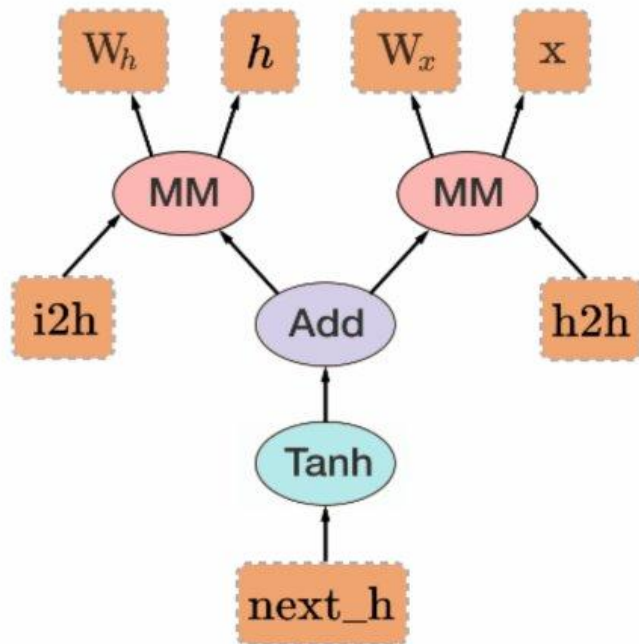
The main **Tensor attributes** related to the graph structure are:

- **data**: Tensor containing the data itself
- **grad**: Tensor containing the gradient (initially set to None)
- **grad_fn**: the function used to compute the gradient

Each **Function** implements two methods:

- **forward**: function application
- **backward**: gradient computation

Autograd



Autograd

requires_grad controls whether autograd tracks operations on a tensor.

For optimizable model parameters

→ **requires_grad=True**

For input data or constant values

→ **requires_grad=False**

detach() returns a new tensor sharing the same storage but outside the graph. Useful to "freeze" an intermediate result (e.g., generating target data from a model, storing predictions for logging).

```
x = torch.randn(3, requires_grad=True)
y = x ** 2

# z shares data with y
# but grad won't flow back through z
z = y.detach()

# False
z.requires_grad == False

# True - same memory!
z.data_ptr() == y.data_ptr()
```

Autograd

`torch.no_grad()` is a context manager that disables gradient tracking for all operations inside. Saves memory and speeds up computation at inference time

While `detach` acts on a single tensor, `no_grad` acts on a block of code.

⚠ In-place modification of a leaf tensor with `requires_grad=True` raises an error as it would silently corrupt the computation graph.

```
# no graph built, no grad stored
with torch.no_grad():
    y_hat = net(x_test)

w = torch.randn(
    3, requires_grad=True)
w.add_(1) # RuntimeError!

# OK: autograd is off
with torch.no_grad():
    w.add_(1)
```

Autograd

torch.inference_mode() is a stricter version of no_grad() available on newer versions of PyTorch (≥ 1.9)

```
with torch.no_grad():
    y = torch.randn(1)
y.requires_grad = True
z = y + 1
z.grad_fn # AddBackward0 object

with torch.inference_mode():
    x = torch.randn(1)
    y = x + 1
# runtime error!!
y.requires_grad = True
```

Autograd

The **requires_grad** attribute is used to specify if the gradient computation should propagate into the Tensor or not, which also stops the **backpropagation**.

For optimizable model parameters \Rightarrow `requires_grad=True`

For input data or constant values \Rightarrow `requires_grad=False`

The method **detach** removes the tensor from the graph, truncating the gradient.

In-place modification is not allowed, as it breaks the automatic differentiation.

At inference time, the **context manager** `torch.no_grad` speeds up computation.

Autograd documentation: <http://pytorch.org/docs/stable/autograd.html>

Your Turn!

Implement a linear regression training loop using only tensors and autograd and gradient descent.

What's `torch.manual_seed(42)`?

```
torch.manual_seed(42)
X = torch.randn(100, 1)
y = 3*X + 2 + torch.randn(100, 1)

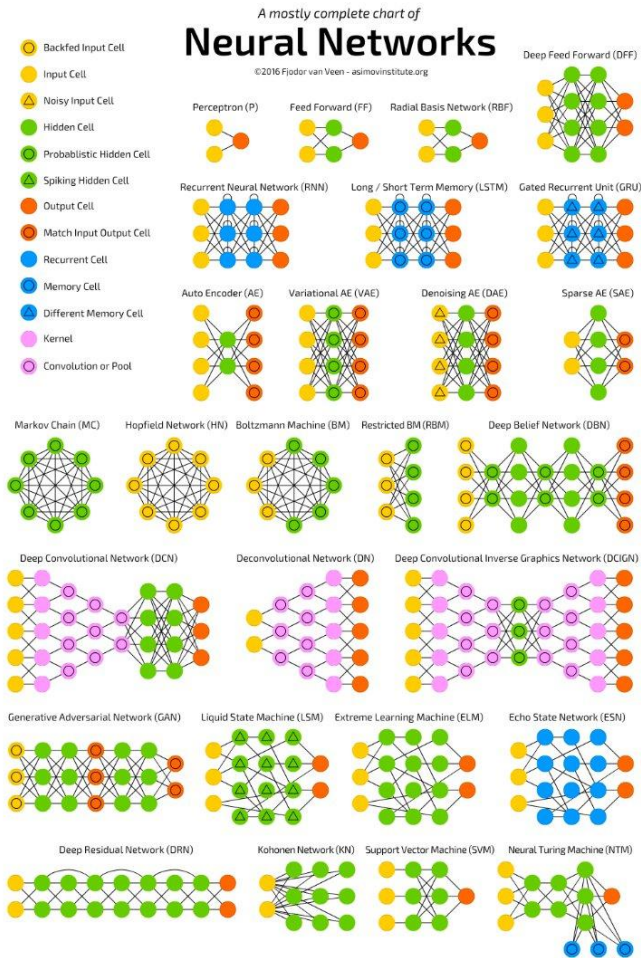
w = ...
b = ...
lr = 0.1

for epoch in range(100):
    ...
```

Building Models and Pipelines

Model Interface

torch.nn contains the basic components to define your neural networks, loss functions, regularization techniques and optimizers.



Module and Parameters

Module is the base class for all the neural network components: Linear, Convolutional, Recurrent Layers...

Each Module contains **Parameter** objects (tensors with a name and `requires_grad=True`). The **parameters()** method returns an iterator over model parameters.

Given a list of modules, PyTorch provides container classes:

- **nn.Sequential(*modules)** chains modules in order; forward is implicit.
- **nn.ModuleList(modules)** when you need custom forward logic (loops, conditionals, skip connections).

⚠ If you use a regular list, the parameters will not be registered!

Forward and Backward

The logic of the module is defined in the **forward()** method, which you can call either as `net(in_tensor)` or `net.forward(in_tensor)`.

The **backward()** step is automatically defined by Autograd, but can be overridden!

It is possible to define forward and backward hooks to debug your model!

Modules can operate in **train** or **eval mode**: `net.train()` or `net.eval()`

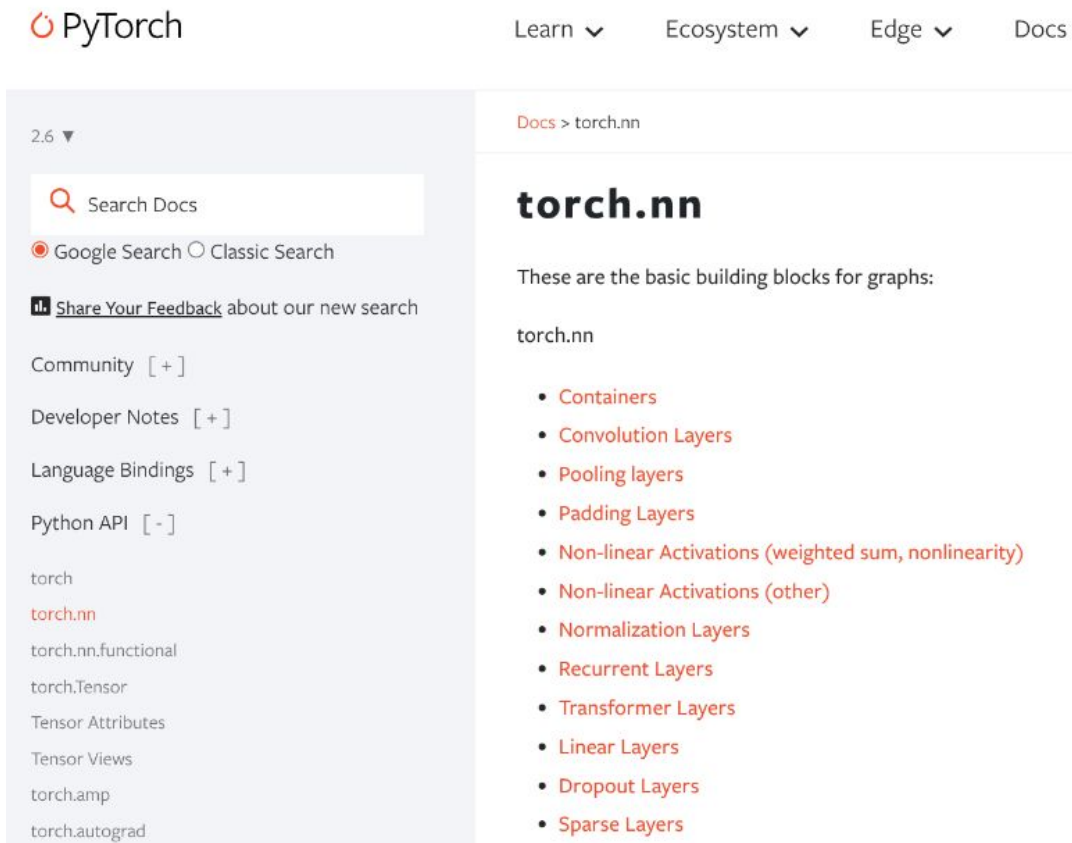
This is useful for layers that define a different behavior during train and test, e.g. Dropout, BatchNormalization...

Existing Modules

There's no need to
reinvent the wheel!

(in most cases, but sometimes you really do: good luck)

PyTorch provides lots
of **common modules**
that can be easily glued
together.



The screenshot shows the PyTorch documentation website. At the top, there is a navigation bar with the PyTorch logo and links for "Learn", "Ecosystem", "Edge", and "Docs". Below the navigation bar, the page title is "torch.nn". The main content area lists various modules under the heading "These are the basic building blocks for graphs:". The listed modules are: Containers, Convolution Layers, Pooling layers, Padding Layers, Non-linear Activations (weighted sum, nonlinearity), Non-linear Activations (other), Normalization Layers, Recurrent Layers, Transformer Layers, Linear Layers, Dropout Layers, and Sparse Layers. On the left side of the page, there is a sidebar with a search bar and a list of navigation links including "Community", "Developer Notes", "Language Bindings", "Python API", and a list of sub-modules like "torch", "torch.nn.functional", "torch.Tensor", etc.

<https://pytorch.org/docs/stable/nn.html>

Functional vs Module API

Apart from modules, PyTorch offers a **functional API** with stateless functions which you can call directly in the forward function.

👍 Rule of thumb: use modules for layers with learnable parameters (Linear, Conv2D) or state (BatchNorm, Dropout). Use functional for pure activations and stateless operations inside forward.

```
# Module style
nn.Sequential(nn.Linear(5, 10),
              nn.ReLU(), nn.Linear(10, 1))
```

```
# Functional style
def forward(self, x):
    x = F.relu(self.fc1(x))
    return self.fc2(x)
```

Datasets and Data Loaders

The module **torch.utils.data** defines classes to handle datasets and load them from data.

DataLoader automatizes mini-batching, shuffling of the dataset, sampling techniques and any pre-processing, and allows parallel loading.



Training Loop

To define a training loop, we need a **loss function** and an **optimizer**.

Always check the **documentation** for the correct **shape** and **input** arguments (does the loss need logits or probabilities? Which dimension should be the last? Is the average for each element or for each sample?)

⚠ Remember to reset gradients using the `zero_grad()` method!

(less talk, more code)

Compilation

torch.compile() JIT-compiles the model's computation graph, fusing operations and reducing Python overhead. In many cases, a single line gives a significant speedup.

No changes to your training loop are needed. The first call is slower (compilation), subsequent calls are faster.

```
model = MyModel()  
model = torch.compile(model)  
  
# Training and inference as usual  
y = model(x)
```

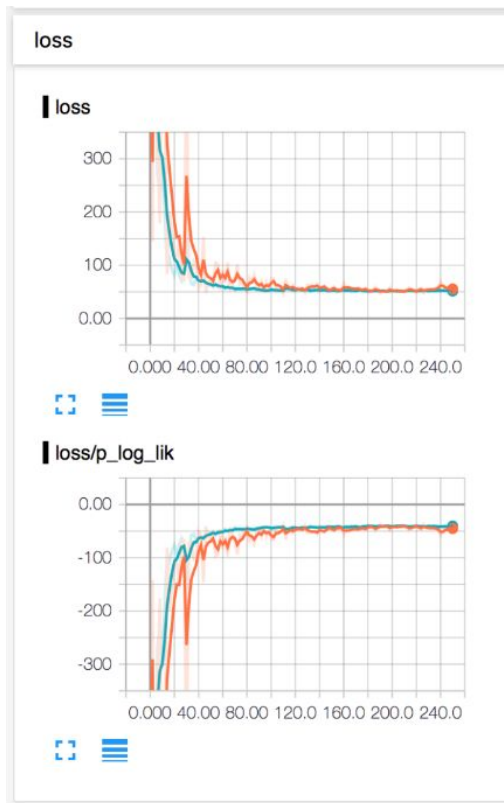
Logging

Several **metrics** can help to understand your model.

Logging them, it's always a good idea!

TensorBoard works great for PyTorch as well.

Otherwise, there are cloud-based commercial products (Weights & Biases, neptune.ai, ...)



Model Serialization

Last, but not least, how do I store my model?

The **state dictionary** stores the value of all model **parameters**.

```
torch.save(the_model.state_dict(), PATH)
```

Then, instantiate the object and reload the state dictionary.

```
net = MyModelClass(*args, **kwargs)
```

```
net.load_state_dict(torch.load(PATH, weights_only=True))
```

PyTorch Ecosystem

To know how things work **under-the-hood** is worth the effort.

... but in practice, most “routine” operations can be abstracted away.

Both **Lightning** and **Transformers by HuggingFace** 🙌 provide APIs for common practices such as training, logging, evaluating and performing inference on Machine Learning models.

Also, tons of libraries in the **PyTorch Ecosystem**: graph neural networks, interpretability, continual learning, federated learning, quantum ML...

Your Turn!

Implement and train a **Convolutional Neural Network** to perform image classification on the **MNIST dataset**.

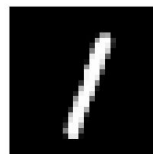


Side Quests:

1. Monitor the performance with a logger,
2. Play around with dropout, batch_norm, etc.
(remember of train vs eval!)
3. Try a Lightning implementation



4 (4)



1 (1)



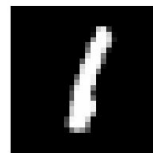
0 (0)



7 (7)



8 (8)



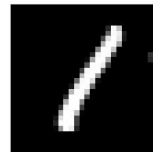
1 (1)



2 (2)



7 (7)



1 (1)