

\*) Basics of Pytorch ( Deep learning framework )



based on Torch



provides two high level features :

- i) Tensor computation ( like Numpy ) with strong GPU acceleration ( CUDA )
- ii) DNN built on a tape - based autodiff system ( or autograd )

→ What are tensors ? Generalisation of matrices & vectors .

( not entirely correct )



matrices → data structures , & tensors provide a relation on how these nos. relate to each other & transform each other



Deep learning also has these data structures that have a relation on the basis of how they are transformed on product, summation, activation, etc



matrix like properties along with structural

Properties. In programs, they can be thought of as large dimensional matrices .

①

```
import torch → Pytorch  
import numpy as np  
import matplotlib.pyplot as plt
```

} library inclusions

②  $x = \text{torch.ones}(3, 2)$  → allocates a tensor of  $3 \times 2$  dimension with value = 1

$x = \text{torch.zeros}(3, 2)$

$x = \text{torch.rand}(3, 2)$  → random value b/w 0 & 1

$x = \text{torch.empty}(3, 2)$  → doesn't initialize values, works with garbage values, and requires to be have values associated.

$y = \text{torch.zeros\_like}(x)$  takes  $x$ , and gives us a vector of the shape of  $x$ , but with 0 values.

$x = \text{torch.linspace}(0, 1, \underbrace{\text{steps}}_{\text{initial + final values}} = 5)$  for a tensor, with total 5 steps

$\text{tensor}([0.0, 0.2500, 0.5000, 0.7500, 1.0000])$

$x = \text{torch.tensor}([[1, 2], [3, 4], [5, 6]])$  → explicitly defining tensors

→ Slice Tensors : how we can access parts of a tensor

↓

similar to the way , it is handled in numpy

>> print(x.size()) → torch.Size([3, 2])

>> print(x[:, 1]) → torch.Tensor([2, 4, 6])

↓

pick up all rows, and for each row pick the 1st index.

>> print(x[0, :]) → torch.Tensor([1, 2])

>> y = x[1, 1] → accessing a particular element

↓

tensor(4)

>> y.item() >> ④ → to fetch the item out of a single element tensor

→ Reshaping tensors

x → (3, 2)

rowmajor order

$$x = \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix}$$

>> y = x.view(2, 3)

$$y = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}$$

>> y = x.view(⑥, -1)

↓

we define all dims except the last, and using -1, it figures the required dimension on its own.

$$y = \begin{bmatrix} [1] \\ [2] \\ [3] \\ [4] \\ [5] \\ [6] \end{bmatrix}$$

\* dimensions of tensors play a very important role in deep learning and leads to many errors.  
Imp. to keep track of dimensions

and axes

## → Simple Tensor Operations.

```
>> x = torch.ones(3, 2)
```

```
>> y = torch.ones(3, 2)
```

```
>> x + y
```

```
>> x - y
```

```
>> x * y
```

} point wise operations.

```
>> z = y.add(x)
```

```
>> z = y.add_(x)
```

addition inplace → modifies y as well here

## → Numpy <→ Pytorch (Interfacing Numpy & Pytorch)

```
>> x_np = x.numpy() → convert to numpy by
```

calling underlying member  
fn. of Tensor class.

```
>> type(x_np)
```

```
<class 'numpy.ndarray'>
```

```
>> a = np.random.random(5)
```

```
>> a_pt = torch.from_numpy(a)
```

from numpy array to torch.Tensor

\*) The conversion from numpy to Tensor is not a copy but rather a bridge or in terms a shallow copy meaning changes made to either are reflected in the other

- ) The CPU performance of PyTorch (based on Torch written in Lua) is better than Numpy. Numpy is already an optimization over Python lists.

```
[29] %%time
for i in range(100):
    a = np.random.randn(100,100)
    b = np.random.randn(100,100)
    c = np.matmul(a, b)
```

↳ CPU times: user 187 ms, sys: 97.6 ms, total: 285 ms  
 Wall time: 152 ms

```
[30] %%time
for i in range(100):
    a = torch.randn([100, 100])
    b = torch.randn([100, 100])
    c = torch.matmul(a, b)
```

↳ CPU times: user 31.2 ms, sys: 435 µs, total: 31.6 ms  
 Wall time: 37 ms

```
[31] %%time
for i in range(10):
    a = np.random.randn(10000,10000)
    b = np.random.randn(10000,10000)
    c = a + b
```

↳ CPU times: user 1min 40s, sys: 2.5 s, total: 1min 42s  
 Wall time: 1min 42s

```
[32] %%time
for i in range(10):
    a = torch.randn([10000, 10000])
    b = torch.randn([10000, 10000])
    c = a + b
```

↳ CPU times: user 21.3 s, sys: 124 ms, total: 21.4 s  
 Wall time: 21.4 s

→ CUDA support → language extension by NVIDIA to support programming GPUs directly available in C & Fortran

⇒ torch.cuda.device\_count() → "checks for the presence of any GPUs on the system"

```
>> print (torch.cuda.device(0))
```

device : GPU  
host : CPU

↓

returns a reference to an object within torch referring to the cuda device.

```
>> print (torch.cuda.get_device_name(0))
```

```
>> cuda0 = torch.device('cuda:0')
```

↓

Variable used to refer the CUDA device, where 0 refers to the index based location

```
>> a = torch.ones(3, 2, device=cuda0)
```

↓

This originally used to get made on CPU, but once we specify the cuda device, it gets built on the GPU.

↓

Ops performed on these tensors also run on the GPU

↓

They give a speedup of almost 1000x as compared to

CPU,

which was itself opt.

```
[10] %%time
for i in range(10):
    a_cpu = torch.randn([10000, 10000])
    b_cpu = torch.randn([10000, 10000])
    b_cpu.add_(a_cpu)
```

```
⇒ CPU times: user 18.6 s, sys: 272 ms, total: 18.9 s
Wall time: 18.9 s
```

```
▶ %%time
for i in range(10):
    a = torch.randn([10000, 10000], device=cuda0)
    b = torch.randn([10000, 10000], device=cuda0)
    b.add_(a)
```

```
⇒ CPU times: user 887 µs, sys: 18 ms, total: 18.9 ms
Wall time: 22.7 ms
```

- ) For matmul, on CPU we take around 4 mins 11 secs, but for the same dim tensors, we take around 33.3 ms
- ) The speedup as expected is much more for parallelizable operations like matrix multiplication.

## → Autograd

>>  $x = \text{torch. ones}([3, 2], \text{ requires\_grad} = \text{True})$

↓  
list highlighting the  
shape of the req.  
tensors

↓  
informs pyTorch, that  
this is something that  
could be differentiated  
wrt .

↓  
i.e. I can create other tensors,  
that are relations on  $x$ , and hence  
can be derived w.r.t.  $x$

>>  $y = x + 5$

↓  
Tensor on DS level is about storing  
these multi-dimensional matrices  
but at a structural level, it relates  
different tensors with each other

>>  $z = y * y + 1$

↓  
The ability to model relations  
between different high order  
matrix like DS is the basic premise  
behind a tensor.

>>  $t = \text{torch.sum}(z)$

```
[3] y = x + 5
print(y)
```

```
[4] tensor([[6., 6.],
           [6., 6.],
           [6., 6.]], grad_fn=<AddBackward0>)
```

\* ) to perform the backward pass for a tensor wrt the independent tensor , we use tensor.backward()

↓

>> t.backward()

\* ) Now, to get the value of the derivative of the dependent tensor on which backward was called wrt the independent tensor, we use ind-tensor.grad

↓

>> print(x.grad)

$$\begin{aligned}y &= x + 5 \\z &= y^2 + 1 \\t &= \sum_{i=1}^n z_i\end{aligned}$$

\* ) Now if the value we call the backward fn is a single value, then it runs great, but if we run it wrt multiple values ( a vector), we need to provide an argument which happens to be the gradient of the  $\Rightarrow$  scalar value wrt the current vector.

↓

>> x = torch.ones([3, 2], requires\_grad = True)

>> y = x + 5

>> g = 1 / (1 + torch.exp(-y))  $\therefore [3, 2]$

>> s = torch.sum(g).

>> s.backward()

>> x.grad

$$\begin{aligned}\frac{dt}{dx_i} &= \frac{dt}{dz_i} \frac{dz_i}{dy_i} \frac{dy_i}{dx_i} \\ \left[ \frac{dz}{dx} \right] &= \left[ 2y \frac{dy}{dx} \right] \\ \frac{dz}{dx} &= 2y \\ &= 12\end{aligned}$$

pointwise  
autograd

→ Now to call `r.backward()` in the same effect as `s.backward()`, we would need to provide the gradient of `s` wrt `r` as an argument in backward, with the same size as `s`.

`>> a = torch.ones([3, 2])`

`>> r.backward(a) →` computes derivative of `r`

`>> x.grad()` wrt `x`, but pointwise multiplies with `a` i.e. the derivative of `s` wrt `x`.

$$\frac{\partial s}{\partial x} = \frac{\partial s}{\partial r} \cdot \underbrace{\left( \frac{\partial r}{\partial x} \right)}_{\text{given by } r.backward(a)}$$

→ Loss fn. with autodiff

Training set  
↑

`>> x = torch.randn([20, 1], requires_grad=True)`  
`>> y = 3 * x - 2`

↓  
ground truth value

↓  
independent tensor

`>> w = torch.tensor([1.], requires_grad=True)`  
`>> b = torch.tensor([1.], requires_grad=True)`

`>> y_hat = w^T x + b`

`>> loss = torch.sum((y_hat - y)^2 / 2)`

`>> loss.backward()`

This is the basic workflow around autograd

`>> w.grad, b.grad`

→ Do it in a loop.

$$\text{learning\_rate} = 0.01$$

w = torch.tensor ([1.], requires\_grad=True)

b = torch.tensor ([1.], requires\_grad=True)

for i in range(10):

x = torch.round ([20, 1])

$$y = 3^x - 2$$

$$\hat{y} = w^T x + b$$

$$\text{loss} = \text{torch.sum} ((\hat{y} - y)^{** 2})$$

loss.backward()

with torch.no\_grad():

$$w = \text{learning\_rate} * w.\text{grad}$$

$$b = \text{learning\_rate} * b.\text{grad}$$

w.grad.zero\_()

b.grad.zero\_()

print (w.item(), b.item())

↓  
There is no  
book keeping  
done for ops  
performed in this  
block of code

sets the

grad=0

for next

iteration

→ A cell can be annotated with `%%time` to get the time taken on execution of the cell.

→ the execution loop on GPU

1. 1. time

$$\text{learning-rate} = 0.001$$

$$N = 1000000$$

$$\text{epochs} = 200$$

$w = \text{torch.randn}([N], \text{requires_grad=True}, \text{device=cuda0})$

$b = \text{torch.ones}([1], \text{requires_grad=True}, \text{device=cuda0})$

for i in range(epochs):

$x = \text{torch.randn}([N], \text{device=cuda0})$

$y = \text{torch.dot}(3 + \text{torch.ones}([N], \text{device=cuda0}),$   
 $x) - 2$

$\hat{y} = \text{torch.dot}(w, x) + b$

$\text{loss} = \text{torch.sum}((\hat{y} - y)^2 / 2)$

$\text{loss.backward}()$

with  $\text{torch.no_grad}()$ :

$w = \text{learning-rate} * w.\text{grad}$

$b = \text{learning-rate} * b.\text{grad}$

$w.\text{grad.zero}()$

$b.\text{grad.zero}()$

\* every tensor part of the computation needs to be constructed on the GPU to utilize the speed up. otherwise

it is no use