



PyTorch

Tutorial II

Course CS 419

Instructor Prof Abir De

Aditya Pande

[22m2108 @iitb.ac.in]

Ninad Gandhi

[22m2151 @iitb.ac.in]

</ 0. Recap />

Previously Discussed

1. MNIST Dataset
2. Neural Network Architecture
 - a. Activations
 - b. Mathematical Form
 - c. Probability and Logits

Previously Discussed: Activations

Sigmoid

$$\sigma(x) = \frac{1}{1 + e^{-(x)}}$$

Tanh

$$\sigma(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

ReLU

$$\sigma(x) = \max(x, 0)$$

LeakyReLU

$$\sigma(x) = \max(0.1x, x)$$

[Visualization](#)

Previously Discussed: Mathematical Form

$$y_i = f(x_i)$$

$$y_i = W_1^T \cdot x_i + b_1$$

$$h_1 = \sigma_1(W_1^T \cdot x_i + b_1)$$

$$h_2 = \sigma_2(W_2^T \cdot (h_1) + b_2)$$

$$\vdots$$

$$h_l = \sigma_l(W_l^T \cdot (h_{l-1}) + b_l)$$

$$logits = W_{out}^T \cdot (h_l) + b_{out}$$

Next in line...

1. Tensors
2. Datasets and Dataloaders
3. Batching
4. Model Architecture
5. Forward Pass
6. Optimizer
7. Loss Objective
8. Device: GPU v CPU
9. Train Trajectory
10. Early Stopping
11. Reproducibility

</ 1. Tensors />

Tensor: What is a tensor?

A torch tensor is a multi-dimensional array used for data and computations in PyTorch

```
import torch

# Create a tensor
x = torch.tensor([1.0, 2.0, 3.0])

# Print the tensor
print("Tensor x:", x)

# Perform basic operations
y = x + 2.0          # Addition
z = x * 3.0          # Multiplication

# Print results
print("x + 2:", y)    # Output: tensor([3.0, 4.0, 5.0])
print("x * 3:", z)    # Output: tensor([3.0, 6.0, 9.0])
```


</ TENSOR : go in BATCHES />

Always use BATCHES !!!!!!!!!

- Groups of data samples processed together
- Batch processing speeds up training and improves efficiency by enabling parallel computation and reducing memory usage

```
any_data_tensor .shape == [batch_size, .,  
.]
```

</ **TENSOR : Attributes** : .shape />

How to access attribute ?

An attribute in is a characteristic or of that tensor

An attribute can be accessed with **.**



```
import torch

# Create a larger tensor
tensor = torch.randn(10, 20, 30) # A tensor
with shape (10, 20, 30)

# Print the shape of the tensor
print("Shape of the tensor:", tensor.shape)
```

</ DataLoader />

- handles batching, shuffling, and loading data efficiently for training and evaluation
- It works with a Dataset to streamline data handling during model training

```
import torch
from torch.utils.data import DataLoader, TensorDataset

# Generate random data
num_samples = 1000
num_features = 20
# Random features
X = torch.randn(num_samples, num_features)
# Random 10-ary labels
y = torch.randint(0, 10, (num_samples,))

# Create a TensorDataset from the random data
dataset = TensorDataset(X, y)

# Create DataLoader
batch_size = 64
data_loader = DataLoader(dataset, batch_size=batch_size,
                          shuffle=True, num_workers=2)
```

Tensor: Methods

- `transpose(T)`
- `squeeze`
- `unsqueeze`
- `size`
- `log`, `log10`, `log2`, ...
- `mean`, `median`, `mode`, `std`
- `nelement`, `element_size`
- `dtype`: `fp16`, `fp32`, `fp64`
- `relu`, `sigmoid`, `tanh`, `softmax`
- `Topk`, `sort`, `sum`

Tensor: Methods

- transpose(T)
- squeeze
- unsqueeze
- size
- log, log10, log2, ...
- mean, median, mode, std
- nelement, element_size
- dtype: fp16, fp32, fp64
- relu, sigmoid, tanh, softmax
- Topk, sort, sum
- flatten, reshape
- view
- max, min
- logsumexp
- ndim
- norm
- Permute
- multinomial, random
- requires_grad

Tensor: Methods (more)

-	abs	-	argsort	-	cholesky_solve	-	diagflat	-	float_power_	-	index_add_
-	abs_	-	argwhere	-	chunk	-	diagonal	-	floor	-	index_copy
-	absolute	-	as_strided	-	clamp	-	diagonal_scatter	-	floor_	-	index_copy_
-	absolute_	-	as_strided_	-	clamp_	-	diff	-	floor_divide	-	index_fill
-	acos	-	as_strided_scatter	-	clamp_max	-	digamma	-	floor_divide_	-	index_fill_
-	acos_	-	as_subclass	-	clamp_max_	-	digamma_	-	fmax	-	index_put
-	acosh	-	asin	-	clamp_min	-	dim	-	fmin	-	index_put_
-	acosh_	-	asin_	-	clamp_min_	-	dim_order	-	fmod	-	index_reduce
-	add	-	asinh	-	clip	-	dist	-	fmod_	-	index_reduce_
-	add_	-	asinh_	-	clip_	-	div	-	frac	-	index_select
-	addbmm	-	atan	-	clone	-	div_	-	frac_	-	indices
-	addbmm_	-	atan2	-	coalesce	-	divide	-	frexp	-	inner
-	addcddiv	-	atan2_	-	coalesce_	-	divide_	-	gather	-	int
-	addcddiv_	-	atan_	-	conj	-	dot	-	gcd	-	int_repr
-	addcmul	-	atanh	-	conj_physical	-	double	-	gcd_	-	inverse
-	addcmul_	-	atanh_	-	contiguous	-	dsplit	-	ge	-	ipu
-	addmm	-	backward	-	copy_	-	dtype	-	ge_	-	is_coalesced
-	addmm_	-	baddbmm	-	copy_	-	eig	-	geometric_	-	is_complex
-	addmv	-	baddbmm_	-	copysign	-	element_size	-	geqr	-	s_conj
-	addmv_	-	bernoulli	-	copysign_	-	eq	-	get	-	is_contiguous
-	addr	-	bernoulli_	-	corrcoef	-	eq_	-	get_device	-	is_cpu
-	addr_	-	bfloat16	-	cos	-	equal	-	grad	-	is_cuda
-	adjoint	-	bincount	-	cos_	-	erf	-	grad_fn	-	is_distributed
-	align_as	-	bitwise_and	-	cosh	-	erf_	-	greater	-	is_floating_point
-	align_to	-	bitwise_and_	-	cosh_	-	erfc	-	greater_	-	is_inference
-	all	-	bitwise_left_shift	-	count_nonzero	-	erfc_	-	greater_equal	-	is_ipu
-	allclose	-	bitwise_left_shift_	-	cov	-	erfinv	-	greater_equal_	-	is_leaf
-	amax	-	bitwise_not	-	cpu	-	erlinv	-	gt	-	is_maia
-	amin	-	bitwise_not_	-	cross	-	exp	-	gt_	-	is_meta
-	aminmax	-	bitwise_or	-	crow_indices	-	exp2	-	half	-	is_mkldnn
-	angle	-	bitwise_or_	-	cuda	-	exp2_	-	hardshrink	-	is_mps
-	any	-	bitwise_right_shift	-	cummax	-	exp_	-	has_names	-	is_mta
-	apply_	-	bitwise_right_shift_	-	cummin	-	expand	-	heaviside	-	is_neg
-	arccos	-	bitwise_xor	-	cumprod	-	expand_as	-	heaviside_	-	is_nested
-	arccos_	-	bitwise_xor_	-	cumprod_	-	expm1	-	histc	-	is_nonzero
-	arccosh	-	bmm	-	cumsun	-	expm1_	-	histogram	-	is_pinned
-	arccosh_	-	bool	-	cumsum_	-	exponential_	-	hsplit	-	is_quantized
-	arcsin	-	broadcast_to	-	data	-	fill	-	hypot	-	is_same_size
-	arcsin_	-	byte	-	deg2rad	-	fill_diagonal_	-	hypot_	-	is_set_to
-	arcsinh	-	cauchy_	-	deg2rad_	-	fix	-	i0	-	is_shared
-	arcsinh_	-	ccol_indices	-	dense_dim	-	fix_	-	i0_	-	is_signed
-	arctan	-	double	-	dequantize	-	flatten	-	igamma	-	is_sparse
-	arctan2	-	ceil	-	del	-	flip	-	igamma_	-	is_sparse_csr
-	arctan2_	-	ceil_	-	detach	-	flipr	-	igammac	-	is_vulkan
-	arctan_	-	cfloat	-	detach_	-	flipud	-	igammac_	-	is_xla
-	arctanh	-	chalf	-	device	-	float	-	imag	-	is_xpu
-	arctanh_	-	char	-	diag	-	float_power	-	index_add	-	isclose...
-	argmax	-	cholesky	-	diag_embed	-		-		-	
-	argmin	-	cholesky_inverse	-		-		-		-	

</ 2. Dataset and Dataloaders />

Built-in Datasets

Torchvision: (<https://pytorch.org/vision/stable/datasets.html>)

- MNIST
- CIFAR-10
- CIFAR-100
- ImageNet

Torchaudio: (<https://pytorch.org/audio/stable/datasets.html>)

- VoxCeleb1Identification
- CMUARCTIC

</ MNIST : dataset example />

1. Data Loading and Preprocessing

```
transform = transforms.Compose([
    transforms.ToTensor(), # Convert image to tensor
    transforms.Normalize((0.5,), (0.5,)) # Normalize data
])
```

Download and load training data

```
trainset = torchvision.datasets.MNIST(root='./data', train=True, download=True, transform=transform)
trainloader = torch.utils.data.DataLoader(trainset, batch_size=64, shuffle=True, num_workers=2)
```

Download and load test data


```
testset = torchvision.datasets.MNIST(root='./data', train=False, download=True, transform=transform)
testloader = torch.utils.data.DataLoader(testset, batch_size=64, shuffle=False, num_workers=2)
```

Custom Dataset

Object type: torch.utils.data.Dataset

Implement 3 methods:

1. `__init__()`
2. `__len__()`
3. `__getitem__()`



```
class CustomDataset(torch.utils.data.Dataset):  
    def __init__(self, data, labels):  
        self.data = data  
        self.labels = labels  
  
    def __len__(self):  
        return len(self.data)  
  
    def __getitem__(self, idx):  
        return self.data[idx], self.labels[idx]
```

Dataloaders

- The `DataLoader` class in PyTorch handles *batching*, *shuffling*, and *parallel* data loading using multiple workers.

Dataloaders

- The **DataLoader** class in PyTorch handles *batching*, *shuffling*, and *parallel* data loading using multiple workers.
- Dataloader wraps around the Dataset object.

Dataloaders

- The **DataLoader** class in PyTorch handles *batching*, *shuffling*, and *parallel* data loading using multiple workers.
- Dataloader wraps around the Dataset object.
- You must split the dataset into train test and val before using the dataloader.

</ MNIST : dataloader example />

- handles batching, shuffling, and loading data efficiently for training and evaluation
- It works with a Dataset to streamline data handling during model training

```
import torch
from torch.utils.data import DataLoader, TensorDataset

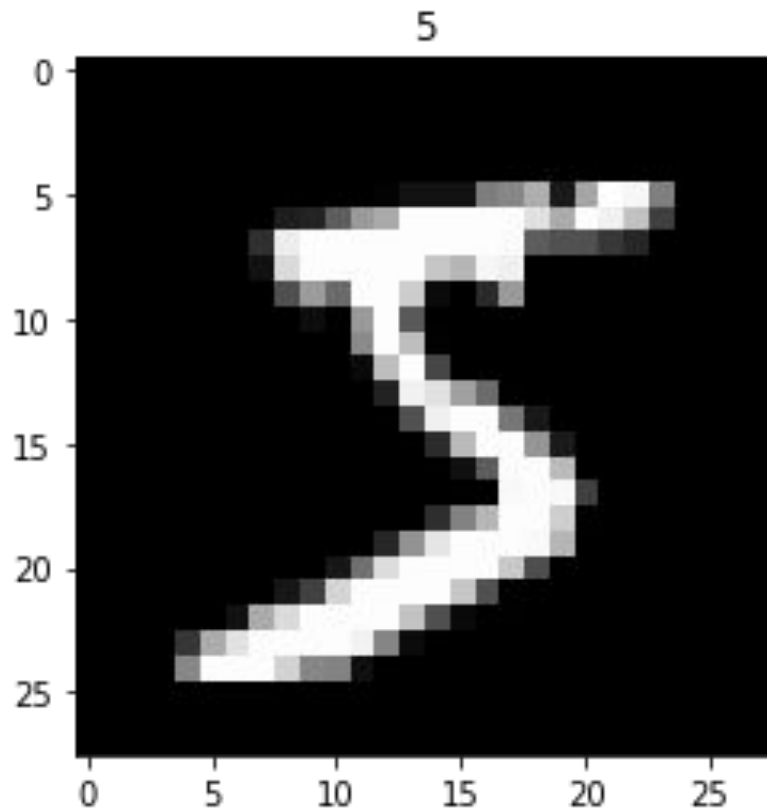
# Generate random data
num_samples = 1000
num_features = 20
# Random features
X = torch.randn(num_samples, num_features)
# Random 10-ary labels
y = torch.randint(0, 10, (num_samples,))

# Create a TensorDataset from the random data
dataset = TensorDataset(X, y)

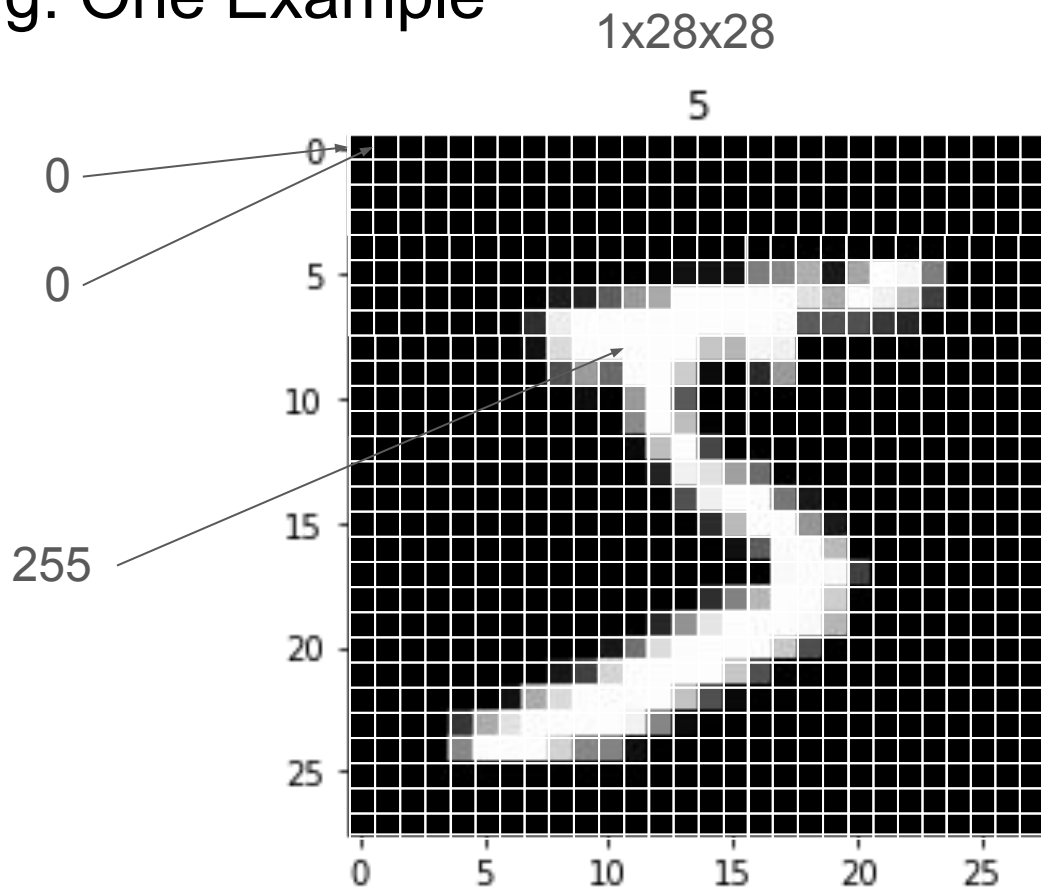
# Create DataLoader
batch_size = 64
data_loader = DataLoader(dataset, batch_size=batch_size,
                          shuffle=True, num_workers=2)
```

</ 3. Batching />

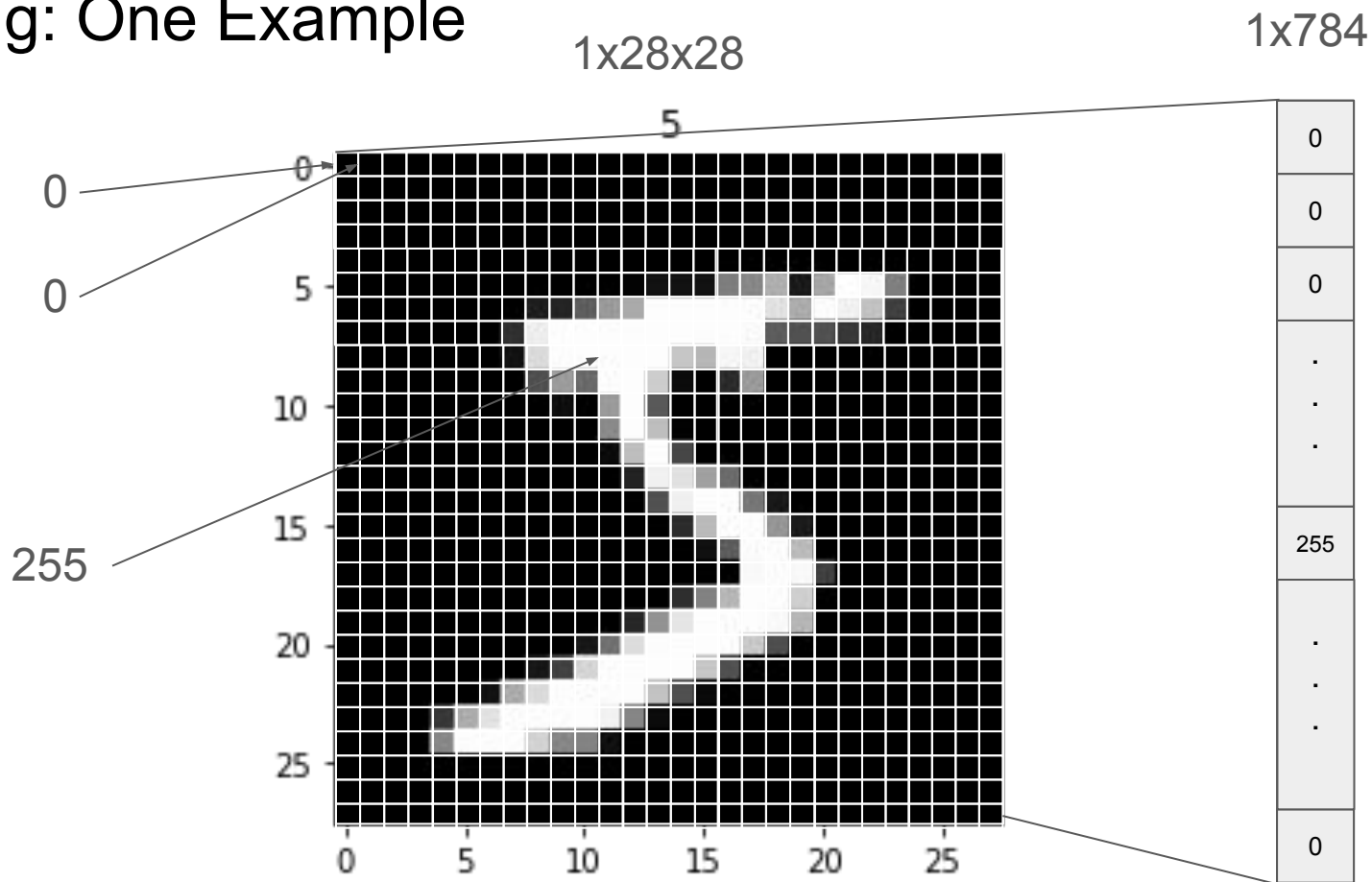
Batching: One Example



Batching: One Example



Batching: One Example



Batching: Sending multiple examples

0	0	0	0	117	0	<div></div>	0	0
0	0	0	0	140	0		0	0
0	120	0	50	20	0		0	0
.
.
.
255	255	200	0	255	255		255	255
.
.
.
0	0	0	0	0	0		0	0

Bx784

</ 4. Model Architecture />

Architecture: Mathematical Form

$$h_1 = \sigma_1(W_1^T \cdot x_i + b_1)$$

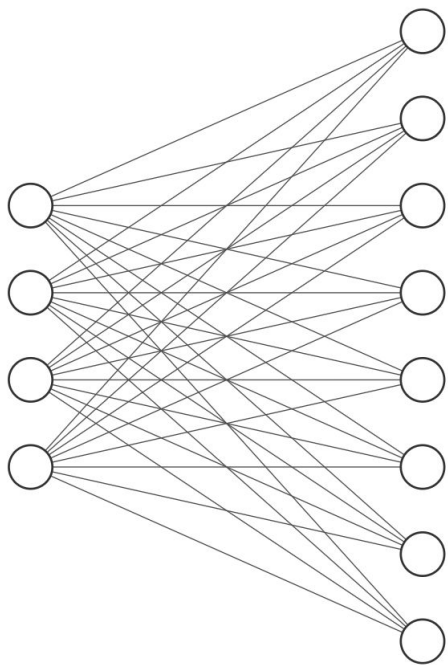
$$h_2 = \sigma_2(W_2^T \cdot (h_1) + b_2)$$

$$\vdots$$

$$h_l = \sigma_l(W_l^T \cdot (h_{l-1}) + b_l)$$

$$logits = W_{out}^T \cdot (h_l) + b_{out}$$

</ Building Block : `nn.Linear(in_dim, out_dim)` />



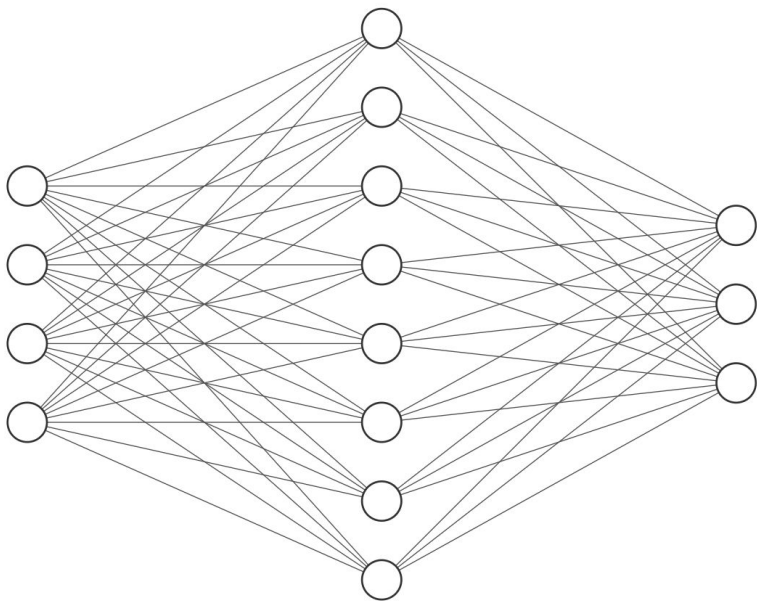
Input Layer $\in \mathbb{R}^4$

Output Layer $\in \mathbb{R}^8$



```
layer = nn.Linear(4, 8)  
out = layer(x)
```

</ Building Block : **Linear & Activations** />



Input Layer $\in \mathbb{R}^4$

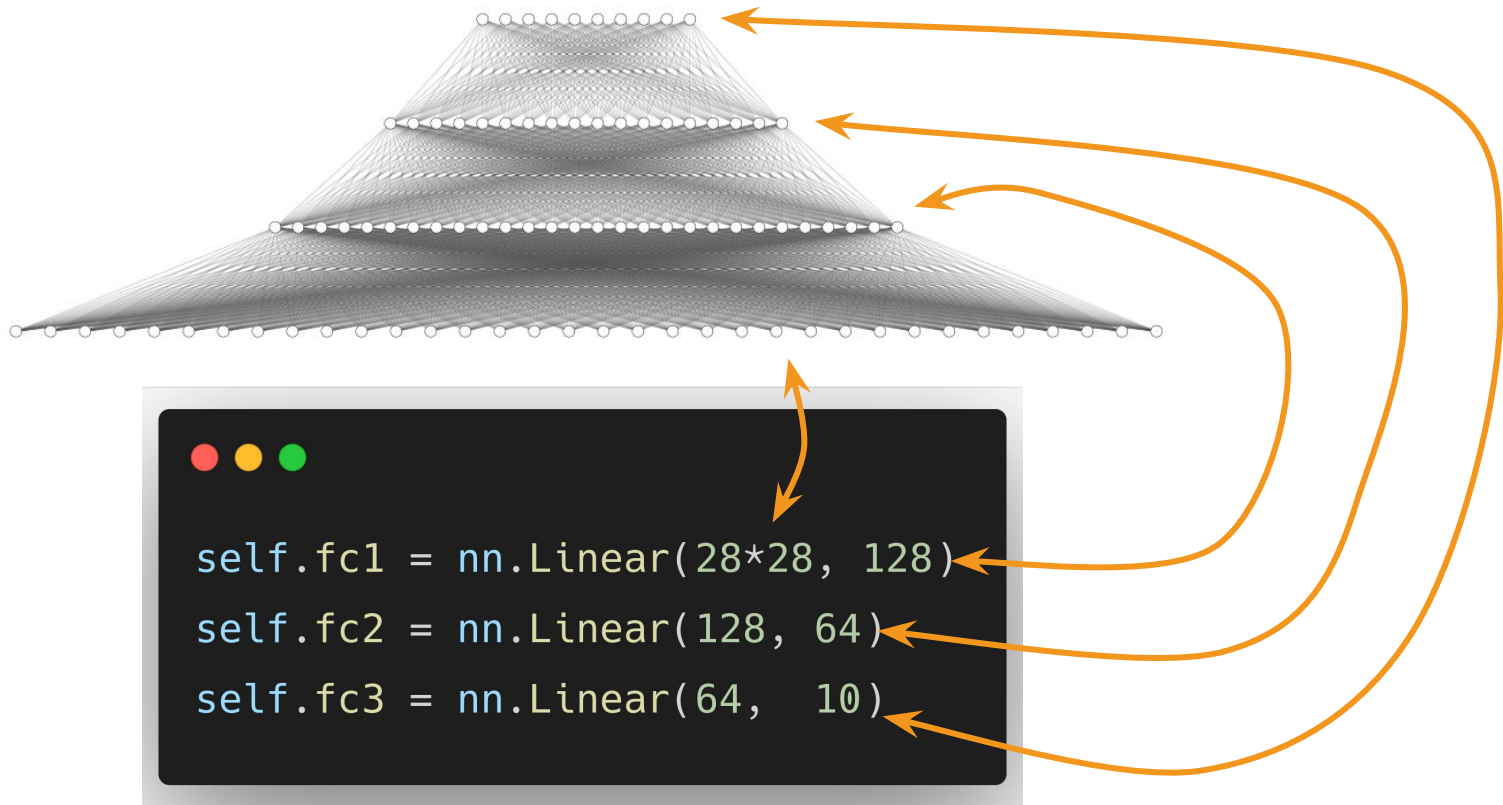
Hidden Layer $\in \mathbb{R}^8$

Output Layer $\in \mathbb{R}^3$



```
layer1 = nn.Linear(4, 8)
layer2 = nn.Linear(8, 3)
out = layer(x)
out = torch.relu(out)
out = torch.relu(out)
```

</ Building Block : **OUR MODEL** />



</ 5. Forward Method />

</ .forward() />

```
class Net(nn.Module):
    def __init__(self):
        super(Net, self).__init__()
        self.fc1 = nn.Linear(28*28, 128)
        self.fc2 = nn.Linear(128, 64)
        self.fc3 = nn.Linear(64, 10)

    def forward(self, x):
        x = x.view(-1, 28*28)

        h1 = torch.relu(self.fc1(x))
        h2 = torch.relu(self.fc2(h1))
        # ...
        h3 = self.fc3(h2)

        return x
```

$$h_1 = \sigma_1(W_1^T \cdot x_i + b_1)$$

$$h_2 = \sigma_2(W_2^T \cdot (h_1) + b_2)$$

⋮

$$h_l = \sigma_l(W_l^T \cdot (h_{l-1}) + b_l)$$

$$\text{logits} = W_{out}^T \cdot (h_l) + b_{out}$$

</ 6. Loss Objective />

</ Loss : Cross Entropy />

For K -class classification using the softmax function, the cross-entropy loss $J(\theta)$ is given by:

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log \left(h_{\theta}(x^{(i)})_k \right)$$

Where:

- m is the number of training examples.
- K is the number of classes.
- $y_k^{(i)}$ is a binary indicator (0 or 1) that represents whether the class label k is the correct classification for the i -th example.
- $h_{\theta}(x^{(i)})_k$ is the predicted probability of the i -th example belonging to class k , given by the softmax function:

$$h_{\theta}(x^{(i)})_k = \frac{\exp(\theta_k^T x^{(i)})}{\sum_{j=1}^K \exp(\theta_j^T x^{(i)})}$$



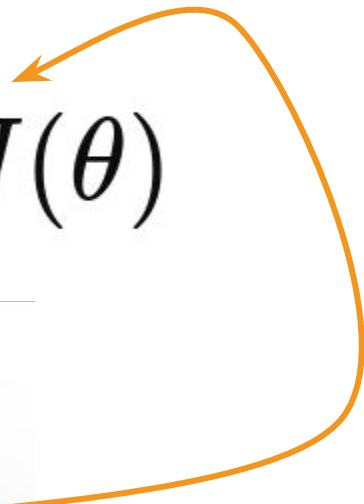
```
criterion = nn.CrossEntropyLoss()
```

</ Loss : Gradient Descent />

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \nabla_{\theta} J(\theta)$$

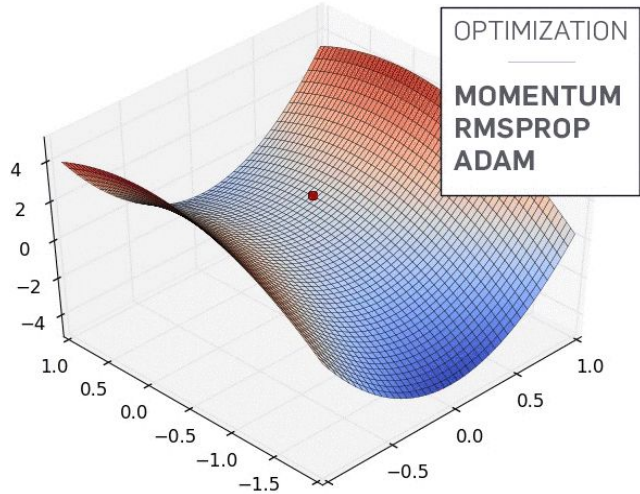


```
loss.backward()
```



</ 7. Optimizer />

Various Optimizers



```
optimizer = optim.Adam(model.parameters(), lr=3e-4)
```

<https://blog.paperspace.com/intro-to-optimization-momentum-rmsprop-adam/>

</ Loss : Gradient Descent />

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \nabla_{\theta} J(\theta)$$



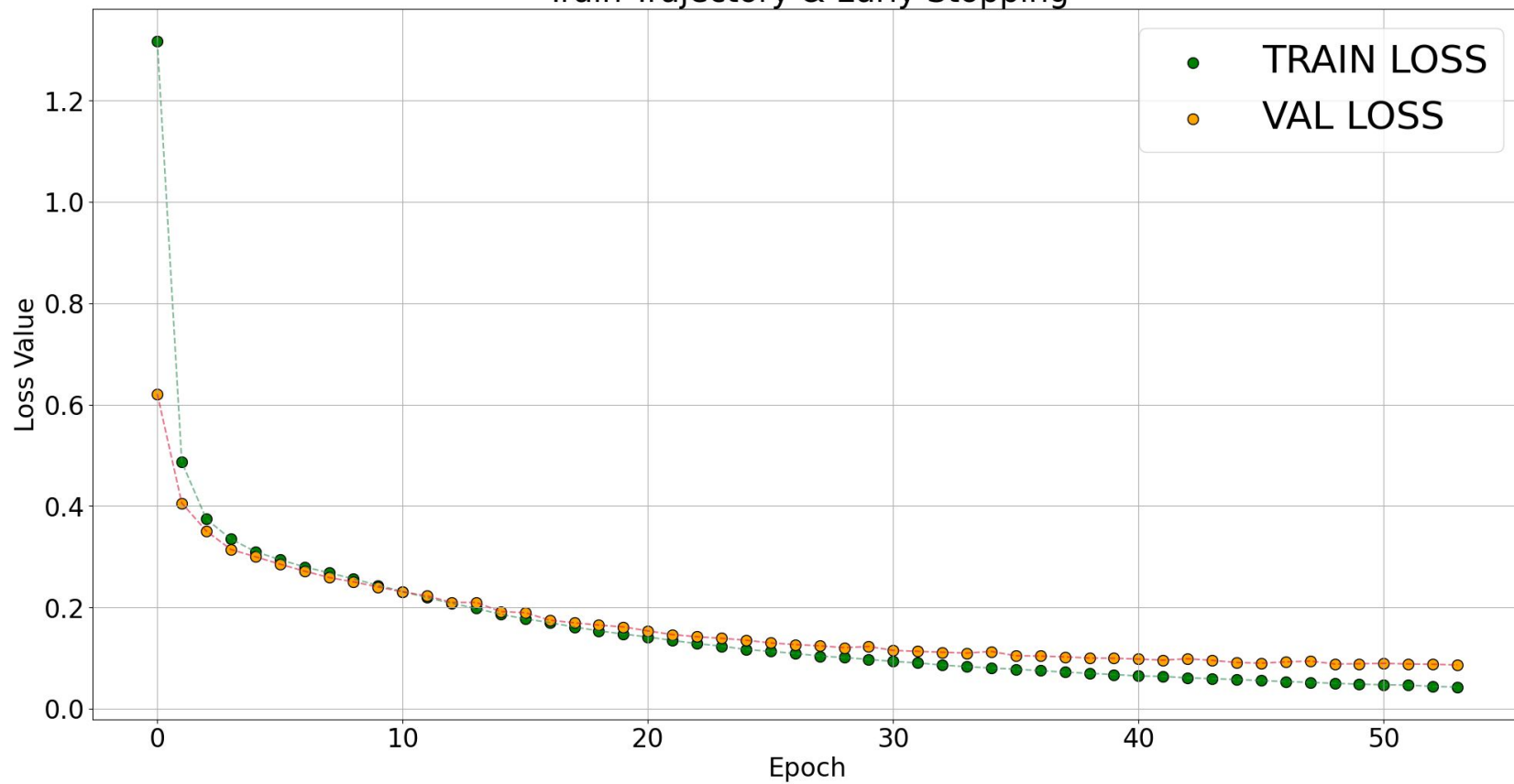
```
loss.backward()
```

```
optimizer.step()
```

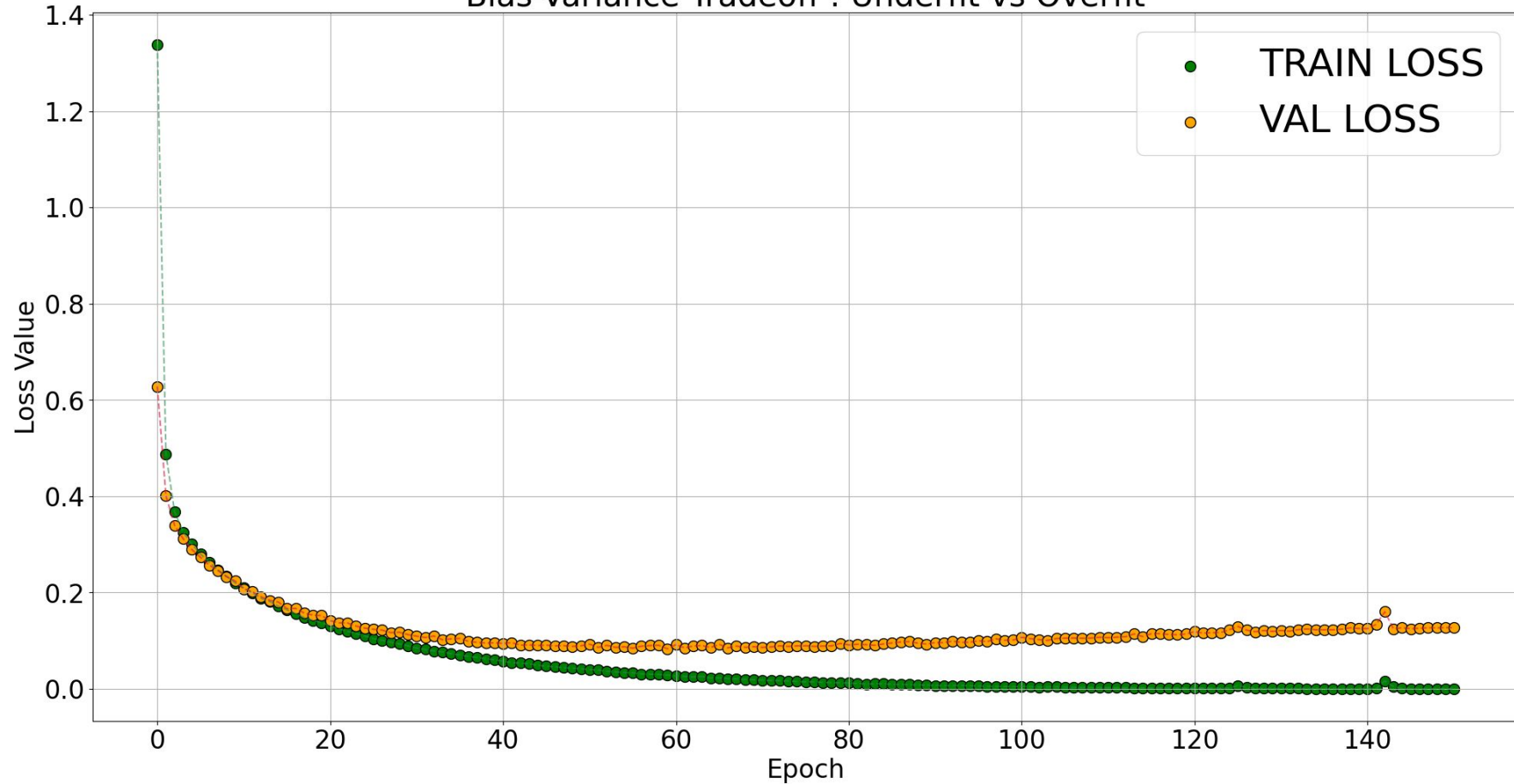

</ 8. Device: GPU v CPU />

</ 9. Train Trajectory />

Train Trajectory & Early Stopping



Bias-Variance Tradeoff : Underfit vs Overfit



</ 10. Reproducibility />

Random Number Generators

- **torch**, **numpy**, and **random** libraries have their own random number generators.

Random Number Generators

- **torch**, **numpy**, and **random** libraries have their own random number generators.
- **Initializing models** generates parameter vectors with random weights.

Random Number Generators

- **torch**, **numpy**, and **random** libraries have their own random number generators.
- **Initializing models** generates parameter vectors with random weights.
- **Shuffling** of data also uses an RNG.

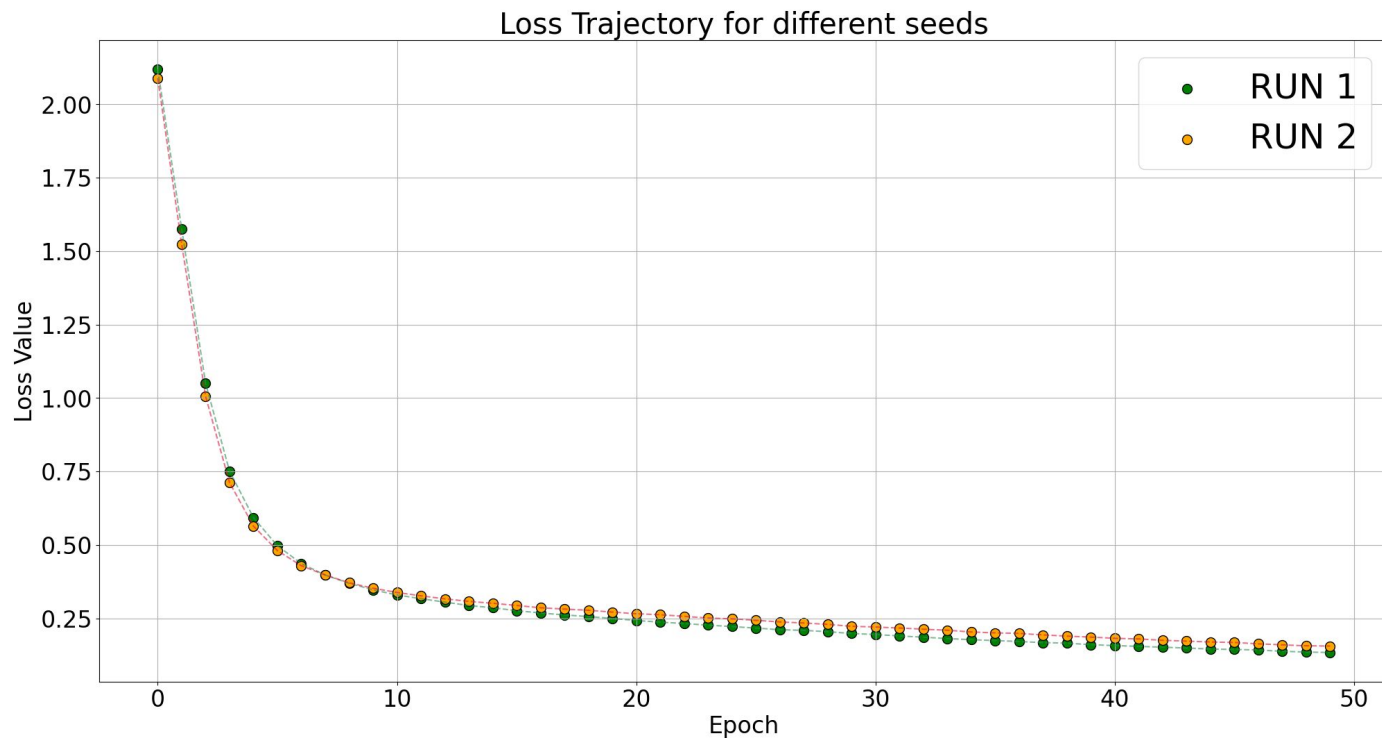
Random Number Generators

- **torch**, **numpy**, and **random** libraries have their own random number generators.
- **Initializing models** generates parameter vectors with random weights.
- **Shuffling** of data also uses an RNG.
- **Dropout** uses a RNG.

Random Number Generators

- **torch**, **numpy**, and **random** libraries have their own random number generators.
 - **Initializing models** generates parameter vectors with random weights.
 - **Shuffling** of data also uses an RNG.
 - **Dropout** uses a RNG.
-
- These RNGs are used in the backend implicitly.
 - **What are the consequences?**

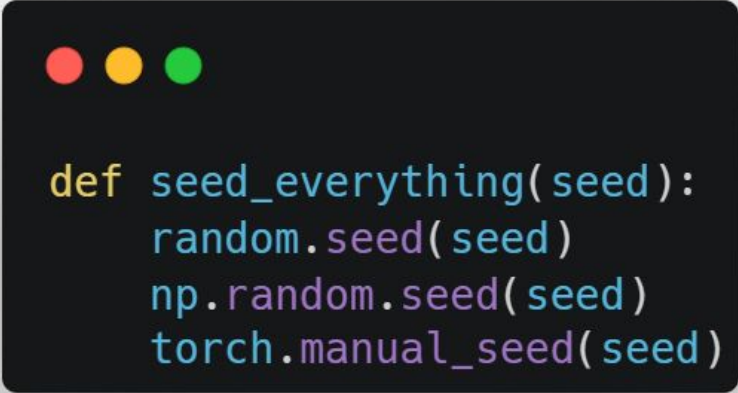
Multiple train runs through the epochs



How to reproduce runs?

Initialize the random seeds for each RNG

- One can write a function as shown.
- Seeding everything is a good practice in machine learning.
- Reproducibility is a big aspect of machine learning research so that people can verify your work.



```
def seed_everything(seed):  
    random.seed(seed)  
    np.random.seed(seed)  
    torch.manual_seed(seed)
```

Non Determinism (Pro stuff)

- It is possible that the training runs are not reproducible even after seeding.
- Non-Deterministic code is because of the hardware and non-deterministic algorithms used by torch in the background.

```
void nll_loss2d_forward_out_cuda_template(
    Tensor& output,
    Tensor& total_weight,
    const Tensor& input,
    const Tensor& target,
    const c10::optional<Tensor>& weight_opt,
    int64_t reduction,
    int64_t ignore_index) {
    // See Note [Writing Nondeterministic Operations]
    // Nondeterministic because of atomicAdd usage in 'sum' or 'mean' reductions.
    if (reduction != at::Reduction::None) {
        at::globalContext().alertNotDeterministic("nll_loss2d_forward_out_cuda_template");
    }
}
```

<https://github.com/pytorch/pytorch/blob/8f1c3c68d3aba5c8898bfb3144988aab6776d549/aten/src/ATen/native/cuda/NLLLoss2d.cu#L236-L240>

Non Determinism (Pro stuff)

- It is possible that the training runs are not reproducible even after seeding.
- Non-Deterministic code is because of the hardware and non-deterministic algorithms used by torch in the background.
- Sometimes asking pytorch to use deterministic algorithms can solve problems. For example, NLLLoss becomes deterministic.



```
torch.use_deterministic_algorithms(True)
```

Non Determinism (Pro stuff)

- It is possible that the training runs are not reproducible even after seeding.
 - Non-Deterministic code is because of the hardware and non-deterministic algorithms used by torch in the background.
 - Sometimes asking pytorch to use deterministic algorithms can solve problems. For example, NLLLoss becomes deterministic.
-
- But you may still have other dependencies that introduce non-determinism. :(

</ 11. Important Points />

These are the things you should do in practice

Always...

These are the things you should do in practice

Always:

1. Use git.

These are the things you should do in practice

Always:

1. Use git.
2. Write modular code.

These are the things you should do in practice

Always:

1. Use git.
2. Write modular code.
3. Document your code.

These are the things you should do in practice

Always:

1. Use git.
2. Write modular code.
3. Document your code.
4. Seed everything before training.

These are the things you should do in practice

Always:

1. Use git.
2. Write modular code.
3. Document your code.
4. Seed everything before training.
5. Reuse the available implementations. (Don't reinvent the wheel)

These are the things you should do in practice

Always:

1. Use git.
2. Write modular code.
3. Document your code.
4. Seed everything before training.
5. Reuse the available implementations. (Don't reinvent the wheel)
6. Use chatGPT to code.

...

</ Thank You />