

6.0 NEED OF Dataset & DataLoader CLASS

24 August 2025 02:48 AM Avinash Yadav

Q) Why do we need the 'dataset' and 'dataloader' classes?

So there is a problem in our code that we wrote while writing a neural network pipeline using 'in-memory'.

This was code that we wrote :-

A simple model with only one neuron

```
PyTorch - Temp.py
1 class MySimpleNN(nn.Module):
2     def __init__(self, num_features):
3         super().__init__()
4         self.linear = nn.Linear(num_features, 1)
5         self.sigmoid = nn.Sigmoid()
6
7     def forward(self, features):
8         out = self.linear(features)
9         out = self.sigmoid(out)
10        return out
11
12 model = MySimpleNN(X_train_tensor.shape[1])
13 learning_rate = 0.1
14 epochs = 50
15
16 loss_function = nn.BCELoss()
17
18 optimizer = torch.optim.SGD(
19     model.parameters(),
20     lr=learning_rate
21 )
22
23 for epoch in range(epochs):
24     y_pred = model(X_train_tensor)
25     loss = loss_function(y_pred, y_train_tensor.view(-1, 1))
26     optimizer.zero_grad()
27     loss.backward()
28     optimizer.step()
29
30     print(f'Epoch: {epoch + 1}, Loss: {loss.item()}')
31
```

The training pipeline where we are building the model, optimizers, and within each epoch we are doing forward pass, loss calculation, clearing gradient, backward pass and parameter updation.

that means to update own parameters for once we are using the whole dataset at once i.e we are sending whole data to forward pass → on the basis of which we are calculating loss → and based on those loss, we are calculating gradient &

→ finally on the basis of those gradient, we are updating old gradient.

So, generally we do not use the batch gradient descent because of bigger reasons :-

1) the first bigger problem is that, batch gradient descent is memory inefficient

2) second problem is that, batch gradient descent does not give better convergence, because of which we may not be able to reach better optimal values of the parameters.

Therefore what we were doing earlier, rather than loading the entire data, load our data in batches.

i.e. Suppose we have a data of 1000 rows, so rather than taking 1000 rows in memory and performing gradient descent on that, it is better to divide our data in 10 batches of each batch having 100 rows & for each batch perform → forward pass → loss calculation → gradient calculation → apply gradient descent.

Once this is done for 1 batch, then load the next batch of 100 rows and perform the same operations.

Repeat it for all the 10 batches.

⇒ This is what called as

Mini Batch Gradient Descent

Simple Solution to Implement Mini Batch Gradient Descent in our above code :-

Same code as above.

```
PyTorch - Temp.py
1 class MySimpleNN(nn.Module):
2     def __init__(self, num_features):
3         super().__init__()
4         self.linear = nn.Linear(num_features, 1)
5         self.sigmoid = nn.Sigmoid()
6
7     def forward(self, features):
8         out = self.linear(features)
9         out = self.sigmoid(out)
10        return out
11
12 model = MySimpleNN(X_train_tensor.shape[1])
13 learning_rate = 0.1
14 epochs = 50
15
16 loss_function = nn.BCELoss()
17
18 optimizer = torch.optim.SGD(
19     model.parameters(),
20     lr=learning_rate
21 )
22
23 batch_size = 32
24 n_samples = len(X_train_tensor)
25
26 for epoch in range(epochs):
27     # Simply Loop over the dataset in chunks of `batch_size`
28     for start_idx in range(0, n_samples, batch_size):
29         end_idx = start_idx + batch_size
30         X_batch = X_train_tensor[start_idx:end_idx]
31         y_batch = y_train_tensor[start_idx:end_idx]
32
33         # Forward Pass
34         y_pred = model(X_batch)
35         loss = loss_function(y_pred, y_batch.view(-1, 1))
36
37         # Update Step
38         optimizer.zero_grad()
39         loss.backward()
40         optimizer.step()
41
42     print(f'Epoch: {epoch + 1}, Loss: {loss.item()}')
```

Now for each epoch, it will take each batch of size 32 rows and perform forward pass → loss calculation → backward pass → parameter updation.

i.e. epoch (1)

Batch (1) → 0-32 rows
Batch (2) → 32-63 rows
Batch (10) → 96-100 rows

forward pass → loss calculation → backward pass → parameter updation

forward pass → loss calculation → backward pass → parameter updation

forward pass → loss calculation → backward pass → parameter updation

epoch (2)

⋮

epoch (10)

same process

same process

So now, this much simple it is to apply mini batch gradient descent using manual code although this code will work but there are some problems in this code.

So the problems are :-

1) No standard interface for data

that means, now we are very easily able to make the batch with the help of X_train_tensor and y_train_tensor

Here we suffered that we have the complete dataset with us. But there are some scenarios where bringing the dataset in itself a big task.

But in our discussed code, we didn't mention where this logic will get implemented

image classification
→ Dog vs Cat
FOLDER-1
1000 images
Data
FOLDER-2
2000 images
Cat

now from this data we have to pick the images & make our X_train_tensor and y_train_tensor
so once we have made them, then only we can create our batches.

2) No Easy way to apply transformations.

There can be scenarios, that we wanted to make batches of 32 images and before sending each batch for training we want to apply some sort of transformation lets say RGB → BGR. But this is not mentioned in our discussed approach.

3) Shuffling and Sampling

1000 Dogs

1000 Dogs

1000 Cats

Shuffled Dogs & Cats Dataset

performing sampling to extract random 32 images to form batches.

But how will this shuffling and sampling will be achieved, is not discussed in our code.

4) Batch Management & Parallelization

suppose we have to extract batches of 32 images, and we have to extract multiple batches parallelly. But the implementation of this parallelization is not discussed in our code.

So in order to solve all these problems, we get 2 classes from PyTorch namely 'Dataset Class' and 'DataLoader Class'.