

6.0 NEED OF Dataset & DataLoader CLASS

24 August 2025 02:48 AM

Avinash Yadav

★ 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 'nn.module'.

This was code that we wrote :-

```
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
14 learning_rate = 0.1
15 epochs = 50
16
17 loss_function = nn.BCELoss()
18
19 optimizer = torch.optim.SGD(
20     model.parameters(),
21     lr=learning_rate
22 )
23
24 for epoch in range(epochs):
25     y_pred = model(X_train_tensor)
26     loss = loss_function(y_pred, y_train_tensor.view(-1, 1))
27     optimizer.zero_grad()
28     loss.backward()
29     optimizer.step()
30
31     print(f'Epoch: {epoch + 1}, Loss: {loss.item()}')
```

The training pipeline where we are building the model, optimizing, and within each epoch we are doing

forward pass, loss calculation, clearing gradient, backward

The biggest flaw in this training pipeline is that for training the model, we are using batch gradient descent

that means, to update our 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 on
clearing gradient, backward
pass and parameter
updation.

loss
→ And based on those
loss, we are calculating
gradient &

so, generally we do not
use the batch gradient
descent because of 2
bigger reasons :-

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

- 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.

perform forward pass → compute error → 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 ⇒

```
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
14 learning_rate = 0.1
15 epochs = 50
16
17 loss_function = nn.BCELoss()
18
19 optimizer = torch.optim.SGD(
20     model.parameters(),
21     lr=learning_rate
22 )
23
24 batch_size = 32
25 n_samples = len(X_train_tensor)
26
27 for epoch in range(epochs):
28     # Simply loop over the dataset in chunks of `batch_size`
29     for start_idx in range(0, n_samples, batch_size):
30         end_idx = start_idx + batch_size
31         X_batch = X_train_tensor[start_idx:end_idx]
```

Same
Code as
that of
above:



defined batch size

$$\frac{320}{32} = 10 \text{ Batches}$$

Starting the training process

Now
we are
using a {

total of
320 rows
90+ divided
into 10 Batches
with each
batch of
size 32 rows

we are using of loops & for loops

for first epoches and for running through batches

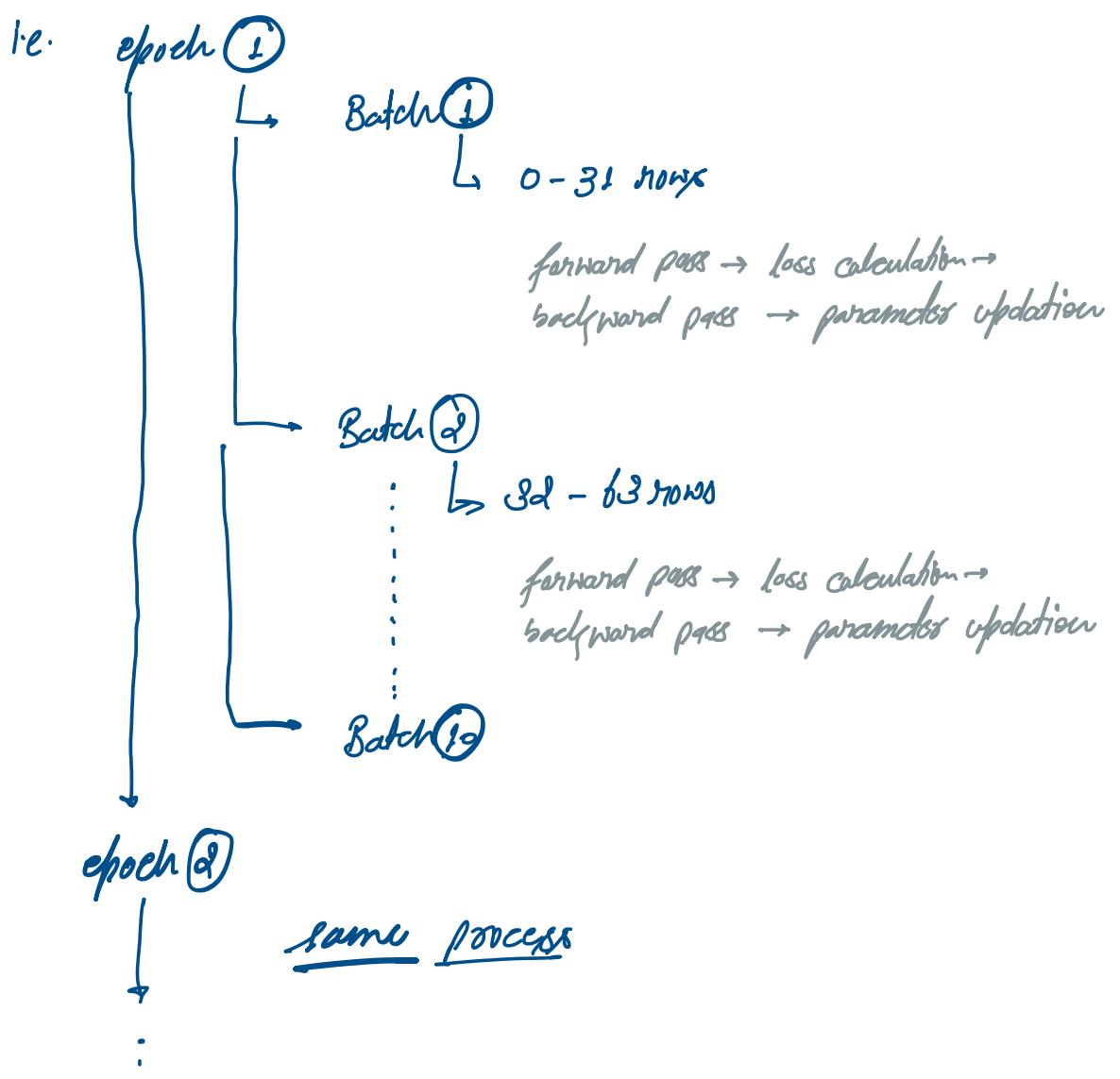
```

30     for start_idx in range(0, n_samples, batch_size):
31         end_idx = start_idx + batch_size
32         X_batch = X_train_tensor[start_idx:end_idx]
33         y_batch = y_train_tensor[start_idx:end_idx]
34
35         # Forward Pass
36         y_pred = model(X_batch)
37         loss = loss_function(y_pred, y_batch.view(-1, 1))
38
39         # Update Step
40         optimizer.zero_grad()
41         loss.backward()
42         optimizer.step()
43
44         print(f'Epoch: {epoch + 1}, Loss: {loss.item()}')

```

size 3d rows

Now for each epoch, it will take each batch of size 3d rows and perform forward pass → loss calculation → backward pass → parameter updation



epoch (50)

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:

1.1 No standard interface for data

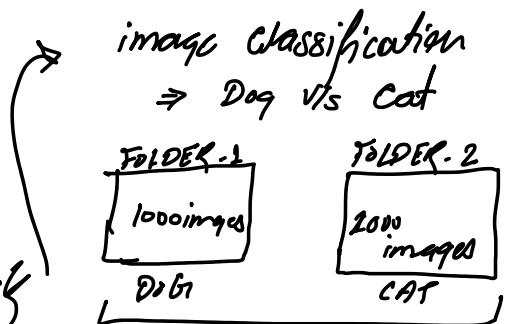
that means, now we are very easily able to make the batch with the help of X_train_tensor and y_train_tensor

```
PyTorch - Temp.py
1 batch_size = 32
2 n_samples = len(X_train_tensor)
3
4 for epoch in range(epochs):
5     # Simply Loop over the dataset in chunks of 'batch_size'
6     for start_idx in range(0, n_samples, batch_size):
7         end_idx = start_idx + batch_size
8         X_batch = X_train_tensor[start_idx:end_idx]
9         y_batch = y_train_tensor[start_idx:end_idx]
10
11         # Forward Pass
12         y_pred = model(X_batch)
13         loss = loss_function(y_pred, y_batch.view(-1, 1))
14
15         # Update Step
16         optimizer.zero_grad()
17         loss.backward()
18         optimizer.step()
19
20     print(f'Epoch: {epoch + 1}, Loss: {loss.item()}')
```

Here we supposed that we have the complete dataset with us.

But there are some scenarios where bringing the dataset is itself a big task

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



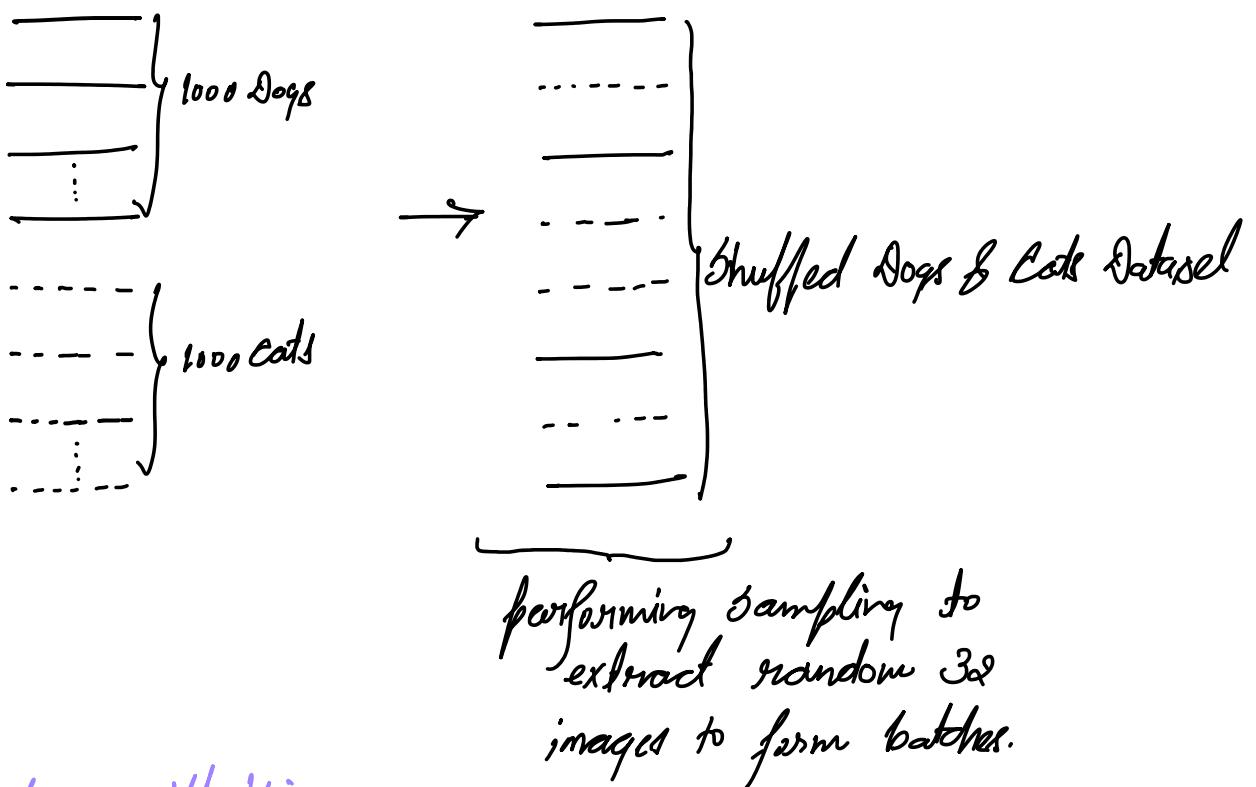
Now from this data we have to pick the images & make our X_train_tensor and y_train_tensor

so once we have made them, then only we can create our batches.

2.8 No Easy way to apply transformations:-

There can be scenarios, that we wanted to make batches of 3d images and before sending each batch for training we want to apply some sort of transformation lets say RGB \rightarrow B&W
But this is not mentioned in our discussed approach.

3.4 Shuffling and Sampling



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

4.4 Batch Management & Parallelization

suppose we have to extract batches of 3d 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'