

Pytorch

Natural Language Processing

Lecture 7



THE GEORGE
WASHINGTON
UNIVERSITY

WASHINGTON, DC

Pytorch

What is PyTorch?

- Open source machine learning library.
- Developed by Facebook's AI Research lab.
- It leverages the power of GPUs.
- Automatic computation of gradients.
- Makes it easier to test and develop new ideas.

- It is pythonic - concise, close to Python conventions
- Strong GPU support
- Autograd - automatic differentiation
- Many algorithms and components are already implemented
- Similar to NumPy

- Easy to **understand** the code
- Has as many type of layers as Torch (Pool, CONV 1,2,3D, LSTM, MLP)
- Lot's of loss functions
- Very similar to **numpy** library.
- Faster compare to other frameworks.
- Allow to build networks which structure is dependent on the computation itself.

Pytorch Levels of Abstraction

- **Tensor**: Like numpy array, but runs on GPU
- **Variable**: Node in a computational graph; stores data and gradient
- **Module**: A neural network layer; may store state or learnable weights

- Create a tensor
- Math Operation
- <https://pytorch.org/docs/stable/index.html>

- NumPy's main object is the homogeneous multidimensional array.
- It is a table of elements (usually numbers), all of the same type, indexed by a tuple of positive integers.
- In NumPy dimensions are called axes. The number of axes is rank.
- For example, the coordinates of a point in 3D space $[1, 2, 1]$ is an array of rank 1, because it has one axis. That axis has a length of 3.

Initializing a Tensor

- Directly from data: Tensors can be created from data. The data type is automatically inferred.

```
1 | import torch
2 | data = [[1,2,3,4],[4,5,6,7]]
3 | data_tensor = torch.Tensor(data)
4 | print(data_tensor)
```

- From a Numpy array: Tensors can be created from a numpy array

```
1 | import numpy as np
2 | data_numpy = np.array(data)
3 | data_tensor = torch.Tensor(data_numpy)
4 | print(data_tensor)
```

- From another tensors: The new tensor retains the properties (shape and datatype) of the argument tensor, unless explicitly overridden.

```
1 | x = torch.ones_like(data_tensor)
2 | print(x)
3 | y = torch.randn_like(data_tensor)
4 | print(y)
```

- Tensor attributes describe the shape, datatype, and the device on which they are sorted.

```
1 | import torch
2 | T = torch.randn(3,4)
3 | print(T.dtype)
4 | print(T.shape)
5 | print(T.device)
```

- By default, tensors are created on the CPU. We need to explicitly move tensors to GPU. We can use a method **.to**. Note that copying data from cpu to gpu and viceversa can be expansive in terms of time and memory.

```
1 | if torch.cuda.is_available():
2 |     tensor_gpu = T.to('cuda')
3 | else:
4 |     print('Tensors are not in the GPU.')
```

- Standard numpy-like indexing and slicing.

```
1 | T1 = torch.ones(4,4)
2 | print(T1[:,0])
3 | print(T1[:, -1])
4 | print(T1[...,-1])
```

- Matrix multiplication (element wise and matrix multiplication)

```
1 | T2 = torch.rand(4,1)
2 | T3 = T1 @ T2 ; print(T3)
3 | T4 = T1.matmul(T2); print(T4)
4 |
5 | T5 = torch.randn(4,4)
6 | T6 = T1 * T5; print(T6)
7 | T7 = T1.mul(T5); print(T7)
```

• Tensor Operations

```
1 | import torch
2 |
3 | x = torch.Tensor(2, 3) ; print(x)
4 | y = torch.rand(2, 3)
5 | z2 = torch.add(x, y); print(z2)
6 | print(torch.is_tensor(z2))
7 | z1 = torch.Tensor(2, 3)
8 | torch.add(x, y, out=z1)
9 |
10 | print(x.size())
11 | print(torch.numel(x))
12 |
13 | k = x.view(6); print(k)
14 | l = x.view(-1, 2); print(l)
15 |
16 | x1 = torch.randn(5, 3).type(torch.FloatTensor); print(x1)
17 | x2 = torch.randn(5, 3).type(torch.LongTensor); print(x2)
18 |
19 | v = torch.arange(9); print(v)
20 | r1 = torch.cat((x, x, x), 0); print(r1)
21 | r2 = torch.stack((v, v)); print(r2)
22 | r3 = torch.chunk(v, 3); print(r3)
```

• Tensor Operations

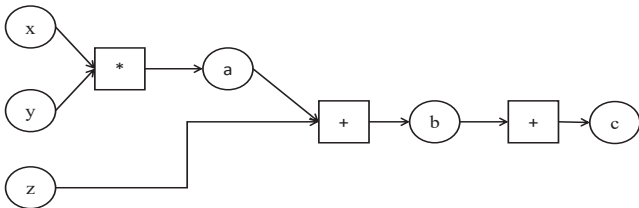
```
1 | import torch
2 |
3 | t = torch.rand(2, 1, 2, 1); print(t)
4 | r = torch.squeeze(t); print(r)
5 | r = torch.squeeze(t, 1); print(r)
6 |
7 |
8 | x = torch.rand([1, 2, 3]); print(x)
9 | r = torch.unsqueeze(x, 0); print(r)
10 | r = torch.unsqueeze(x, 1); print(r)
11 |
12 | v = torch.arange(9).reshape(3,3)
13 | # flatten a Tensor and return elements with given indexes
14 | r = torch.take(v, torch.LongTensor([0, 4, 2]))
15 | r = torch.transpose(v, 0, 1); print(r)
16 |
17 | mat1 = torch.randn(2, 3)
18 | mat2 = torch.randn(3, 4)
19 | r = torch.mm(mat1, mat2)
20 |
21 | v1 = torch.ones(3)
22 | r = torch.diag(v1)
```

- A Pytorch Variable is a node in a computational graph.
- The autograd package provides automatic differentiation for all operations on Tensors.
- “ `autograd.Variable` is the central class of the package. It wraps a Tensor, and supports nearly all of operations defined on it.
- Once you finish your computation you can call `.backward()` and have all the gradients computed automatically. “

Computational Graph - Numpy

- Lets use numpy and code the graph.

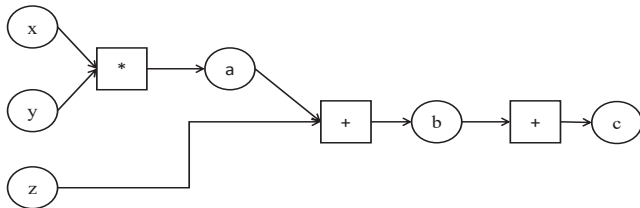
```
1 | import numpy as np
2 | x = np.random.randn(3,4)
3 | y = np.random.randn(3,4)
4 | z = np.random.randn(3,4)
5 | a = x * y
6 | b = a + z
7 | c = np.sum(b)
8 | grad_c = 1.0
9 | grad_b = grad_c * np.ones((3,4))
10 | grad_a = grad_b.copy()
11 | grad_z = grad_b.copy()
12 | grad_x = grad_a * y
13 | grad_y = grad_a * x
```



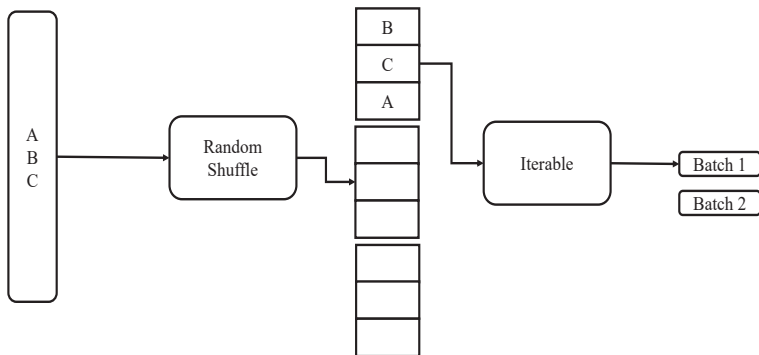
Computational Graph - PyTorch

- Lets use torch Variable and code the graph.

```
1 | import torch
2 | from torch.autograd import Variable
3 | x = Variable(torch.randn(3,4), requires_grad = True)
4 | y = Variable(torch.randn(3,4), requires_grad = True)
5 | z = Variable(torch.randn(3,4), requires_grad = True)
6 | a = x * y
7 | b = a + z
8 | c = torch.sum(b)
9 |
10 | c.backward()
11 | print(x.grad.data)
12 | print(y.grad.data)
13 | print(z.grad.data)
```



- Shuffle and iterate the data.



- We define neural network by subclassing `nn.Module` and initialize the neural network layers in `__init__`. Every `nn.Module` subclass implements the operations on input data in forward methos.

```
1 | from torch import nn
2 | class model(nn.Module):
3 |     def __init__(self, hidden_dim):
4 |         super(model, self).__init__()
5 |         self.linear1 = nn.Linear(1, hidden_dim)
6 |         self.act1 = torch.sigmoid
7 |         self.linear2 = nn.Linear(hidden_dim, 1)
8 |     def forward(self, x):
9 |         return self.linear2(self.act1(self.linear1(x)))
```

Loss Function, Optimizers, Training Loop

- We pass our models output logits to loss function

```
1 | import torch
2 | criterion = nn.CrossEntropyLoss()
3 | optimizer = torch.optim.SGD(model.parameters(), lr = 0.001)
```

- Training Loop

```
1 | def train_loop (dataloader, model, loss, optimizer):
2 |     size = len(dataloader.dataset)
3 |     for batch, (X,y) in enumerate(dataloader):
4 |         pred = model(X)
5 |         loss = criterion(pred,y)
6 |         optimizer.zero_grad()
7 |         loss.backward()
8 |         optimizer.step()
```

Sample Code

- To begin, load the required libraries. The first package you'll import is the torch library.

```
1 | import torch
2 | from torch import nn
3 | from torchtext.data.utils import get_tokenizer
4 | from torchtext.vocab import build_vocab_from_iterator
5 | from torch.utils.data import DataLoader
6 | from torchtext.datasets import AG_NEWS
7 | import time
```

- For example, the AG_NEWS dataset iterators yield the raw data as a tuple of label and text.

```
1 | train_iter = list(AG_NEWS(split='train'))
2 | test_iter = list(AG_NEWS(split='test'))
3 | print(train_iter[0])
```

Prepare data processing pipelines

- The first step is to build a vocabulary with the raw training dataset.

```
1 | device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
2 | tokenizer = get_tokenizer('basic_english')
3 | def yield_tokens(data_iter):
4 |     for _, text in data_iter:
5 |         yield tokenizer(text)
```

- Here we use built in factory function `build_vocab_from_iterator` which accepts iterator that yield list or iterator of tokens.

```
1 | vocab = build_vocab_from_iterator(yield_tokens(train_iter), specials=["<unk>"])
2 | vocab.set_default_index(vocab["<unk>"])
3 | print(vocab(['here', 'is', 'an', 'example']))
```

- Users can also pass any special symbols to be added to the vocabulary.

Prepare data processing pipelines

- Prepare the text processing pipeline with the tokenizer and vocabulary.

```
1 | text_pipeline = lambda x: vocab(tokenizer(x))  
2 | label_pipeline = lambda x: int(x) - 1
```

- The text pipeline converts a text string into a list of integers based on the lookup table defined in the vocabulary. The label pipeline converts the label into integers.

```
1 | print(text_pipeline('here is the an example'))  
2 | print(label_pipeline('10'))
```

Generate data batch and iterator

- **torch.utils.data.DataLoader** is recommended for PyTorch user.
- It works with a map-style dataset that implements the `getitem()` and `len()` protocols, and represents a map from indices/keys to data samples.
- Before sending to the model, `collate_fn` function works on a batch of samples generated from `DataLoader`.
- The input to `collate_fn` is a batch of data with the batch size in `DataLoader`.
- Pay attention here and make sure that `collate_fn` is declared as a top level def.
- This ensures that the function is available in each worker.

- This is a general architecture of any dataloader.

```
1 | from torch.utils.data import DataLoader
2 | import numpy as np
3 | from torch.utils.data import Dataset
4 | import torch
5 | class Custom_Data_loader(Dataset):
6 |     def __init__(self):
7 |         y = np.loadtxt('data-diabetes.csv', delimiter=',', dtype=np.float32)
8 |         self.len = y.shape[0]
9 |         self.x_data = torch.from_numpy(y[:, 0:-1])
10 |        self.y_data = torch.from_numpy(y[:, [-1]])
11 |
12 |    def __len__(self):
13 |        return self.len
14 |
15 |    def __getitem__(self, index):
16 |        return self.x_data[index], self.y_data[index]
17 |
18 |
19 | dataset = Custom_Data_loader()
20 | train_loader = DataLoader(dataset=dataset, batch_size=32, shuffle=True, num_workers=2)
```

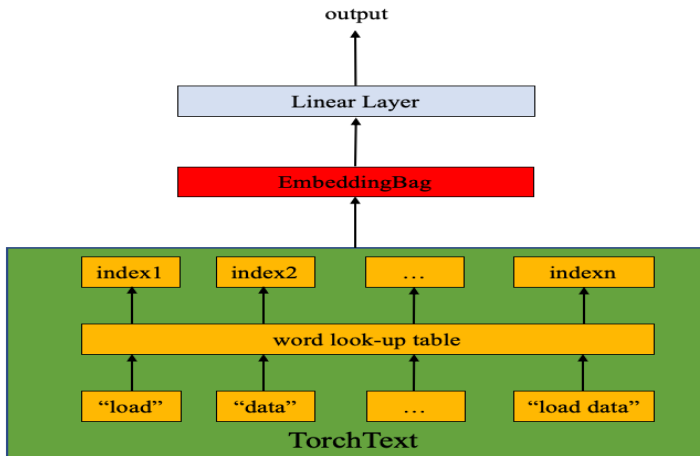
Generate data batch and iterator

- In this example, the text entries in the original data batch input are packed into a list and concatenated as a single tensor for the input of nn.EmbeddingBag.

```
1 | def collate_batch(batch):
2 |     label_list, text_list, offsets = [], [], [0]
3 |     for (_label, _text) in batch:
4 |         label_list.append(label_pipeline(_label))
5 |         processed_text = torch.tensor(text_pipeline(_text), dtype=torch.int64)
6 |         text_list.append(processed_text)
7 |         offsets.append(processed_text.size(0))
8 |     label_list = torch.tensor(label_list, dtype=torch.int64)
9 |     offsets = torch.tensor(offsets[:-1]).cumsum(dim=0)
10 |    text_list = torch.cat(text_list)
11 |    return label_list.to(device), text_list.to(device), offsets.to(device)
```

- The offset is a tensor of delimiters to represent the beginning index of the individual sequence in the text tensor. Label is a tensor saving the labels of individual text entries.

- Model Architecture



- The model is composed of the `nn.EmbeddingBag` layer plus a linear layer for the classification purpose.

```
1 | class TextClassificationModel(nn.Module):  
2 |     def __init__(self, vocab_size, embed_dim, num_class):  
3 |         super(TextClassificationModel, self).__init__()  
4 |         self.embedding = nn.EmbeddingBag(vocab_size, embed_dim)  
5 |         self.fc = nn.Linear(embed_dim, num_class)  
6 |         self.init_weights()  
7 |     def init_weights(self):  
8 |         initrange = 0.5  
9 |         self.embedding.weight.data.uniform_(-initrange, initrange)  
10 |        self.fc.weight.data.uniform_(-initrange, initrange)  
11 |        self.fc.bias.data.zero_()  
12 |    def forward(self, text, offsets):  
13 |        embedded = self.embedding(text, offsets)  
14 |        return self.fc(embedded)
```

- `nn.EmbeddingBag` with the default mode of “mean” computes the mean value of a “bag” of embeddings. Although the text entries here have different lengths, `nn.EmbeddingBag` module requires no padding here since the text lengths are saved in offsets.

- We build a model with the embedding dimension of 64. The vocab size is equal to the length of the vocabulary instance. The number of classes is equal to the number of labels,

```
1 | num_class = len(set([label for (label, text) in train_iter]))  
2 | vocab_size = len(vocab)  
3 | emsize = 64  
4 | model = TextClassificationModel(vocab_size, emsize, num_class).to(device)
```

Define functions to train the model

• Model training

```
1 | def train(dataloader):
2 |     model.train()
3 |     total_acc, total_count = 0, 0
4 |     log_interval = 500
5 |     start_time = time.time()
6 |
7 |     for idx, (label, text, offsets) in enumerate(dataloader):
8 |         optimizer.zero_grad()
9 |         predicted_label = model(text, offsets)
10 |         loss = criterion(predicted_label, label)
11 |         loss.backward()
12 |         torch.nn.utils.clip_grad_norm_(model.parameters(), 0.1)
13 |         optimizer.step()
14 |         total_acc += (predicted_label.argmax(1) == label).sum().item()
15 |         total_count += label.size(0)
16 |         if idx % log_interval == 0 and idx > 0:
17 |             elapsed = time.time() - start_time
18 |             print(' | epoch {:3d} | {:5d}/{:5d} batches '
19 |                 ' | accuracy {:.3f}'.format(epoch, idx, len(dataloader),
20 |                                             total_acc/total_count))
21 |             total_acc, total_count = 0, 0
22 |             start_time = time.time()
```

Define Function to Train the Model

● Model training

```
1 | def train(dataloader):
2 |     model.train()
3 |     total_acc, total_count = 0, 0
4 |     log_interval = 500
5 |     start_time = time.time()
6 |
7 |     for idx, (label, text, offsets) in enumerate(dataloader):
8 |         optimizer.zero_grad()
9 |         predicted_label = model(text, offsets)
10 |         loss = criterion(predicted_label, label)
11 |         loss.backward()
12 |         torch.nn.utils.clip_grad_norm_(model.parameters(), 0.1)
13 |         optimizer.step()
14 |         total_acc += (predicted_label.argmax(1) == label).sum().item()
15 |         total_count += label.size(0)
16 |         if idx % log_interval == 0 and idx > 0:
17 |             elapsed = time.time() - start_time
18 |             print(' | epoch {:3d} | {:5d}/{:5d} batches '
19 |                   ' | accuracy {:.3f}'.format(epoch, idx, len(dataloader),
20 |                                               total_acc/total_count))
21 |             total_acc, total_count = 0, 0
22 |             start_time = time.time()
```

Define Function to Evaluate the Model

● Model evaluation

```
1 | def evaluate(dataloader):
2 |     model.eval()
3 |     total_acc, total_count = 0, 0
4 |
5 |     with torch.no_grad():
6 |         for idx, (label, text, offsets) in enumerate(dataloader):
7 |             predicted_label = model(text, offsets)
8 |             loss = criterion(predicted_label, label)
9 |             total_acc += (predicted_label.argmax(1) == label).sum().item()
10 |            total_count += label.size(0)
11 |    return total_acc/total_count
```

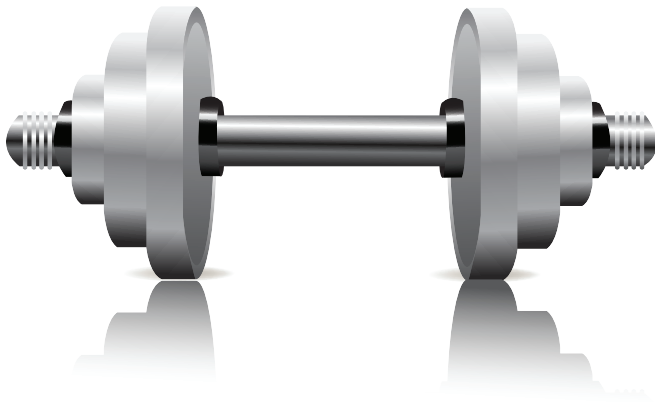
● Hyperparameters

```
1 | EPOCHS = 10
2 | LR = 0.001
3 | BATCH_SIZE = 64
4 | criterion = torch.nn.CrossEntropyLoss()
5 | optimizer = torch.optim.Adam(model.parameters(), lr=LR)
```


- Train the model with dataloader

```
1 | train_dataloader = DataLoader(train_iter, batch_size=BATCH_SIZE,  
2 |                               shuffle=True, collate_fn=collate_batch)  
3 | test_dataloader = DataLoader(test_iter, batch_size=BATCH_SIZE,  
4 |                               shuffle=True, collate_fn=collate_batch)  
5 |  
6 |  
7 | for epoch in range(1, EPOCHS + 1):  
8 |     epoch_start_time = time.time()  
9 |     train(train_dataloader)  
10 |     accu_val = evaluate(test_dataloader)  
11 |     print('-' * 59)  
12 |     print('| end of epoch {:3d} | time: {:.2f}s | '  
13 |           'valid accuracy {:.3f} '.format(epoch,  
14 |                                           time.time() - epoch_start_time, accu_val))  
15 |     print('-' * 59)
```

- Exercise 1 to 1
- [Class-Ex-Lecture7.py](#)



Word2vec

- Instead of capturing co occurrence counts directly.
- Predict surrounding words of every word.
- Predict surrounding words in a window of length m of every word.
- Objective function: Maximize the log probability of any context word given the current center word:

$$J(\theta) = \frac{1}{T} \sum_{t=1}^T \sum_{-m \leq j \leq m} \log p(w_{t+j} | w_t)$$

- Where θ represents all variables we optimize.

- There are two architectures used by Word2vec.

1- Continuous Bag of words (CBOW)

2- Skip gram

- In CBOW, the current word is predicted using the window of surrounding context windows.
- Skip-Gram performs opposite of CBOW which implies that it predicts the given sequence or context from the word.
- Word2vec provides an option to choose between CBOW (continuous Bag of words) and skip-gram. Such parameters are provided during training of the model.

Windowing and Target Word

- Consider a local window of a target word

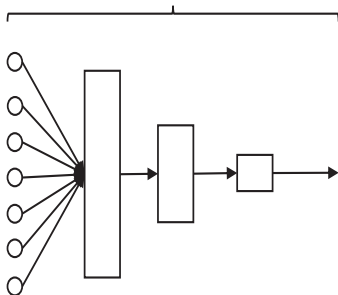
Window 2 : Where there is a **will** there is a way

Window 3 : Where there is a **will** there is a way

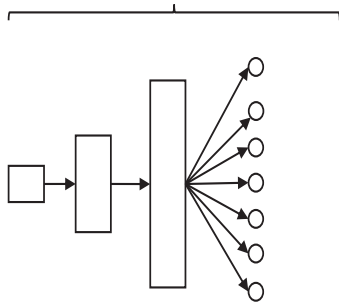
Window 4 : Where there is a **will** there is a way

Word to Vector (word2vec)

Continuous Bag of Word



Skip Gram Model



- Exercise 2 to 2
- [Class-Ex-Lecture7.py](#)

