



# PyTorch Tutorial

## 08. Dataset and DataLoader

数据集序列      数据加载 MiniBatch

# Revision: Manual data feed

numpy → Tensor

```
xy = np.loadtxt( 'diabetes.csv.gz' , delimiter= ',' , dtype=np.float32)
x_data = torch.from_numpy(xy[:, :-1])
y_data = torch.from_numpy(xy[:, [-1]])

.....

for epoch in range(100):
    # 1. Forward
    y_pred = model(x_data)
    loss = criterion(y_pred, y_data)
    print(epoch, loss.item())
    # 2. Backward
    optimizer.zero_grad()
    loss.backward()
    # 3. Update
    optimizer.step()
```

梯度下降 = 全部Batch, 计算整个  
随机梯度下降: 1个样本, 随机性可克服驻点  
Mini-Batch (见Lecture 8-3 p20)

Use all of the data

# Terminology: Epoch, Batch-Size, Iterations

Mini-Batch: 组成数据集的循环

```
# Training cycle
for epoch in range(training_epochs):
    # Loop over all batches
    for i in range(total_batch):
```

轮数, 所有样本训练几轮  
Mini-Batch 个数

Definition: **Epoch**

One forward pass and one backward pass of **all the training examples**.



Definition: **Batch-Size**

The **number of training examples** in one forward backward pass.

每轮训练样本的数量大小  
Mini-Batch

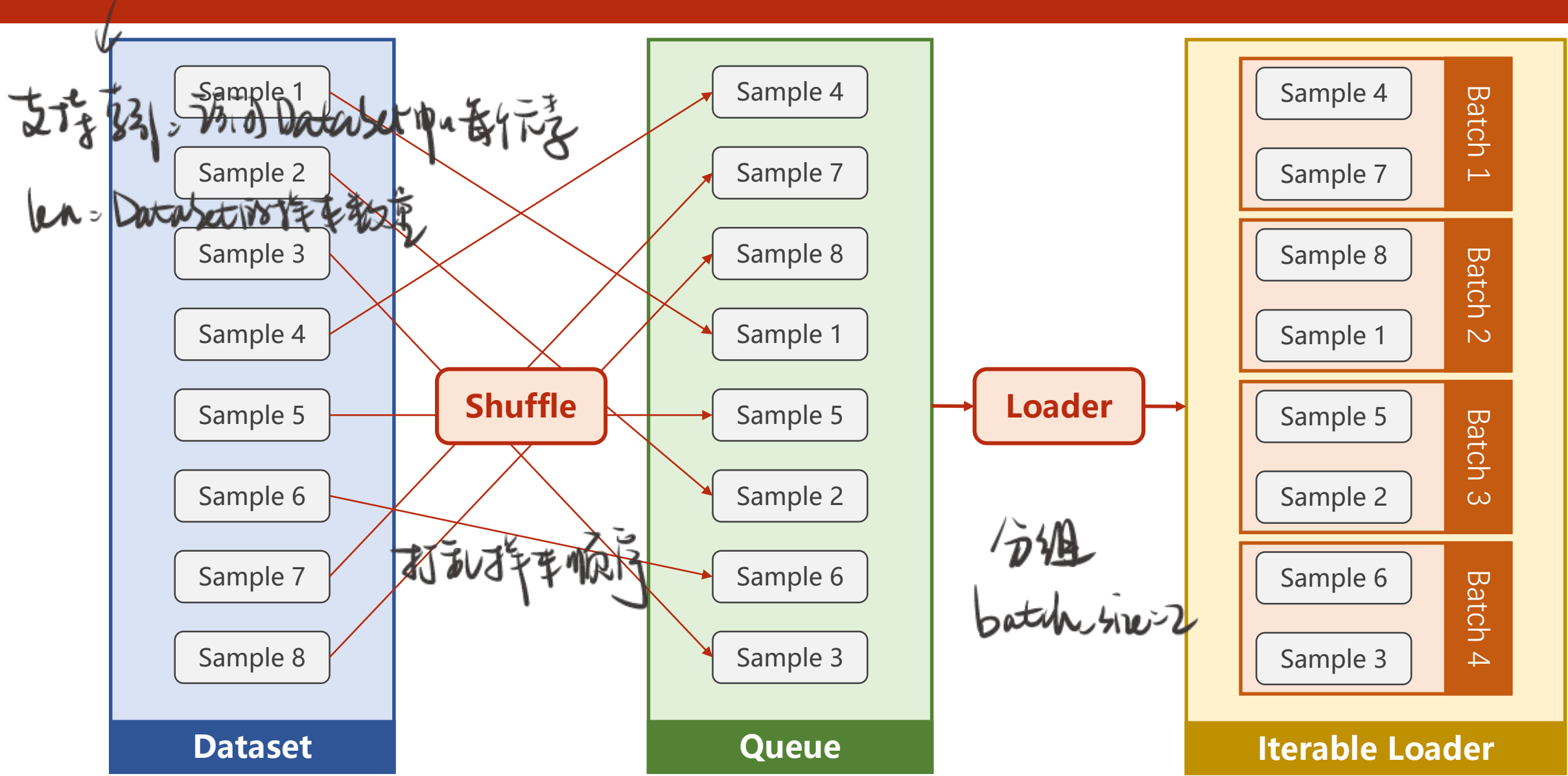
Definition: **Iteration**

Number of passes, each pass using [**batch size**] number of examples.

每轮训练的 Mini-Batch 个数

eg = 总样本数 10000, Mini-Batch size = 1000  
⇒ Iteration = 10000 / 1000 = 10

# DataLoader: batch\_size=2, shuffle=True



# How to define your Dataset

```
import torch
from torch.utils.data import Dataset
from torch.utils.data import DataLoader

class DiabetesDataset(Dataset):
    def __init__(self):
        pass

    def __getitem__(self, index):
        pass

    def __len__(self):
        pass

dataset = DiabetesDataset()
train_loader = DataLoader(dataset=dataset,
                           batch_size=32,
                           shuffle=True,
                           num_workers=2)
```

**Dataset** is an **abstract** class. We can define our class inherited from this class.

Dataset 是一个抽象类 (不可实例化)  
必须被其他类继承

# How to define your Dataset

```
import torch
from torch.utils.data import Dataset
from torch.utils.data import DataLoader

class DiabetesDataset(Dataset):
    def __init__(self):
        pass

    def __getitem__(self, index):
        pass

    def __len__(self):
        pass

dataset = DiabetesDataset()
train_loader = DataLoader(dataset=dataset,
                           batch_size=32,
                           shuffle=True,
                           num_workers=2)
```

**DataLoader** is a class to help us loading data in PyTorch.

DataLoader = 加载数据  
可迭代

# How to define your Dataset

```
import torch
from torch.utils.data import Dataset
from torch.utils.data import DataLoader

class DiabetesDataset(Dataset):
    def __init__(self):
        pass

    def __getitem__(self, index):
        pass

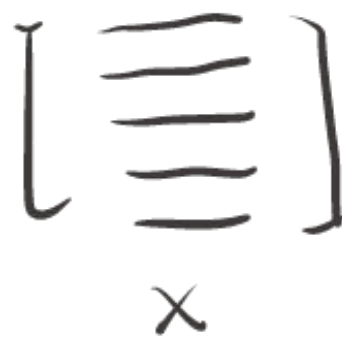
    def __len__(self):
        pass

dataset = DiabetesDataset()
train_loader = DataLoader(dataset=dataset,
                           batch_size=32,
                           shuffle=True,
                           num_workers=2)
```

**DiabetesDataset** is inherited from abstract class **Dataset**.

1- All Data 加载所有数据 → [1]

适用于本身样本量不大的数据集

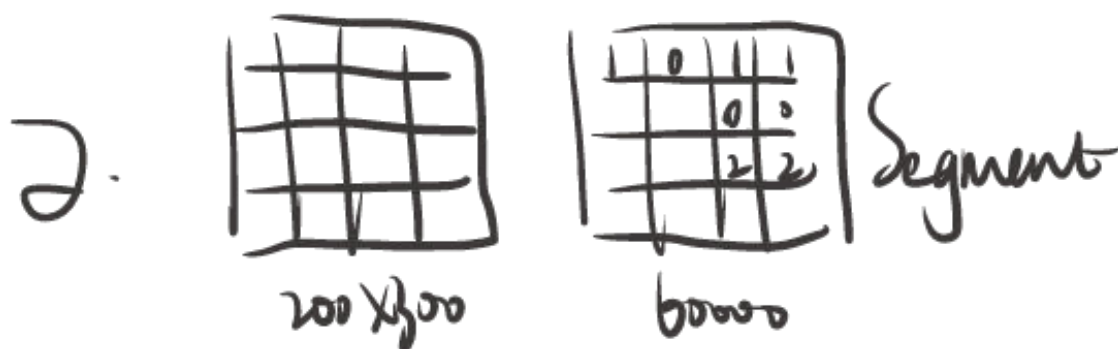


若  $x$  = 简单数据类型数值  
⇒ 全部加载

Lecturer : Hongpu Liu

Lecture 8-7

PyTorch Tutorial @ SLAM Research Group



判断输入图像每个像素 =  $\begin{cases} 0 \\ 1 \end{cases}$



输出 也是一个二维图像 (每个 = 60000 像素)

若  $x$  = 复杂数据

⇒ 文件名存入列表

getitem 读取第  $i$  个文件



# How to define your Dataset

```
import torch
from torch.utils.data import Dataset
from torch.utils.data import DataLoader

class DiabetesDataset(Dataset):
    def __init__(self):
        pass

    def getitem(self, index):
        pass

    def __len__(self):
        pass

dataset = DiabetesDataset()
train_loader = DataLoader(dataset=dataset,
                           batch_size=32,
                           shuffle=True,
                           num_workers=2)
```

支持下标索引, 取出对应位置的元素

The expression, **dataset[index]**, will call this magic function.



# How to define your Dataset

```
import torch
from torch.utils.data import Dataset
from torch.utils.data import DataLoader

class DiabetesDataset(Dataset):
    def __init__(self):
        pass

    def __getitem__(self, index):
        pass

    def __len__(self):
        pass

dataset = DiabetesDataset()
train_loader = DataLoader(dataset=dataset,
                           batch_size=32,
                           shuffle=True,
                           num_workers=2)
```

返回数据集中数据个数

This magic function returns length of dataset.

# How to define your Dataset

```
import torch
from torch.utils.data import Dataset
from torch.utils.data import DataLoader
```

```
class DiabetesDataset(Dataset):
    def __init__(self):
        pass
```

```
    def __getitem__(self, index):
        pass
```

```
    def __len__(self):
        pass
```

```
dataset = DiabetesDataset()
train_loader = DataLoader(dataset=dataset,
                           batch_size=32,
                           shuffle=True,
                           num_workers=2)
```

实例化数据集类

Construct DiabetesDataset object.

# How to define your Dataset

```
import torch
from torch.utils.data import Dataset
from torch.utils.data import DataLoader
```

```
class DiabetesDataset(Dataset):
    def __init__(self):
        pass

    def __getitem__(self, index):
        pass

    def __len__(self):
        pass
```

```
dataset = DiabetesDataset()
train_loader = DataLoader(dataset=dataset,
                           batch_size=32,
                           shuffle=True,
                           num_workers=2)
```

数据加载器

数据加载器

数据加载器

读入Mini-batch. 是否成批读入?

Initialize loader with **batch-size**,  
**shuffle**, **process number**.

# Extra: *num\_workers* in Windows

```
train_loader = DataLoader(dataset=dataset,
                           batch_size=32,
                           shuffle=True,
                           num_workers=2)
```

.....

```
for epoch in range(100):
    for i, data in enumerate(train_loader, 0):
```

.....

So we have to **wrap** the code with an if-clause to protect the code from executing multiple times.

Linux、Windows多进程程序不一样

The implementation of multiprocessing is different on Windows, which uses **spawn** instead of **fork**.

So left code will cause:

RuntimeError:

An attempt has been made to start a new process before the current process has finished its bootstrapping phase.

This probably means that you are not using fork to start your child processes and you have forgotten to use the proper idiom in the main module:

```
if __name__ == '__main__':
    freeze_support()
    ...
```

The "freeze\_support()" line can be omitted if the program is not going to be frozen to produce an executable.

## Extra: *num\_workers* in Windows

```
train_loader = DataLoader(dataset=dataset,
                           batch_size=32,
                           shuffle=True,
                           num_workers=2)

.....
if __name__ == '__main__':
    for epoch in range(100):
        for i, data in enumerate(train_loader, 0):
            # 1. Prepare data
```

So we have to **wrap** the code with an if-clause to protect the code from executing multiple times.



# Example: Diabetes Dataset

自定义模型

```
class DiabetesDataset(Dataset):  
    def __init__(self, filepath):  
        ✓ xy = np.loadtxt(filepath, delimiter=',', dtype=np.float32)  
        self.len = xy.shape[0]  
        ✓ self.x_data = torch.from_numpy(xy[:, :-1])  
        ✓ self.y_data = torch.from_numpy(xy[:, [-1]])  
  
    def __getitem__(self, index):  
        return self.x_data[index], self.y_data[index]  
  
    def __len__(self):  
        return self.len  
  
dataset = DiabetesDataset('diabetes.csv.gz')  
train_loader = DataLoader(dataset=dataset, batch_size=32, shuffle=True, num_workers=2)
```

759/32=23 每批次23个样本数据

# Example: Using DataLoader

训练

所有数据都训练100遍

$X$   
 $Y$   
 $\begin{bmatrix} \equiv \\ \equiv \\ \equiv \end{bmatrix}$   $\begin{bmatrix} \equiv \\ \equiv \\ \equiv \end{bmatrix}$

```
for epoch in range(100):  
    for i, data in enumerate(train_loader, 0):  
        # 1. Prepare data  
        inputs, labels = data  
        # 2. Forward  
        y_pred = model(inputs)  
        loss = criterion(y_pred, labels)  
        print(epoch, i, loss.item())  
        # 3. Backward  
        optimizer.zero_grad()  
        loss.backward()  
        # 4. Update  
        optimizer.step()
```

Input 转成 Tensor

$i$  = 第  $i$  个 batch

data = (X, Y)

$X[i]$   
 $Y[i]$



# Classifying Diabetes

```
import numpy as np
import torch
from torch.utils.data import Dataset, DataLoader

class DiabetesDataset(Dataset):
    def __init__(self, filepath):
        xy = np.loadtxt(filepath, delimiter=',', dtype=np.float32)
        self.len = xy.shape[0]
        self.x_data = torch.from_numpy(xy[:, :-1])
        self.y_data = torch.from_numpy(xy[:, [-1]])

    def __getitem__(self, index):
        return self.x_data[index], self.y_data[index]

    def __len__(self):
        return self.len

dataset = DiabetesDataset('diabetes.csv.gz')
train_loader = DataLoader(dataset=dataset,
                           batch_size=32,
                           shuffle=True,
                           num_workers=2)

class Model(torch.nn.Module):
    def __init__(self):
        super(Model, self).__init__()
        self.linear1 = torch.nn.Linear(8, 6)
        self.linear2 = torch.nn.Linear(6, 4)
        self.linear3 = torch.nn.Linear(4, 1)
        self.sigmoid = torch.nn.Sigmoid()

    def forward(self, x):
        x = self.sigmoid(self.linear1(x))
        x = self.sigmoid(self.linear2(x))
        x = self.sigmoid(self.linear3(x))
        return x

model = Model()

criterion = torch.nn.BCELoss(size_average=True)
optimizer = torch.optim.SGD(model.parameters(), lr=0.01)

for epoch in range(100):
    for i, data in enumerate(train_loader, 0):
        # 1. Prepare data
        inputs, labels = data
        # 2. Forward
        y_pred = model(inputs)
        loss = criterion(y_pred, labels)
        print(epoch, i, loss.item())
        # 3. Backward
        optimizer.zero_grad()
        loss.backward()
        # 4. Update
        optimizer.step()
```

1

Prepare dataset  
Dataset and Dataloader

2

Design model using Class  
inherit from nn.Module

3

Construct loss and optimizer  
using PyTorch API

4

Training cycle *Mini-Batch*  
forward, backward, update

# The following dataset loaders are available

- MNIST
- Fashion-MNIST
- EMNIST
- COCO
- LSUN
- ImageFolder
- DatasetFolder
- Imagenet-12
- CIFAR
- STL10
- PhotoTour

## torchvision.datasets

All datasets are subclasses of `torch.utils.data.Dataset` i.e, they have `__getitem__` and `__len__` methods implemented. Hence, they can all be passed to a `torch.utils.data.DataLoader` which can load multiple samples parallelly using `torch multiprocessing` workers. For example:

```
imagenet_data = torchvision.datasets.ImageFolder('path/to/imagenet_root/')
data_loader = torch.utils.data.DataLoader(imagenet_data,
                                          batch_size=4,
                                          shuffle=True,
                                          num_workers=args.nThreads)
```

# Example: MNIST Dataset

```
import torch
from torch.utils.data import DataLoader
from torchvision import transforms
from torchvision import datasets

train_dataset = datasets.MNIST(root='../dataset/mnist',
                               train=True,
                               transform=transforms.ToTensor(),
                               download=True)

test_dataset = datasets.MNIST(root='../dataset/mnist',
                              train=False,
                              transform=transforms.ToTensor(),
                              download=True)

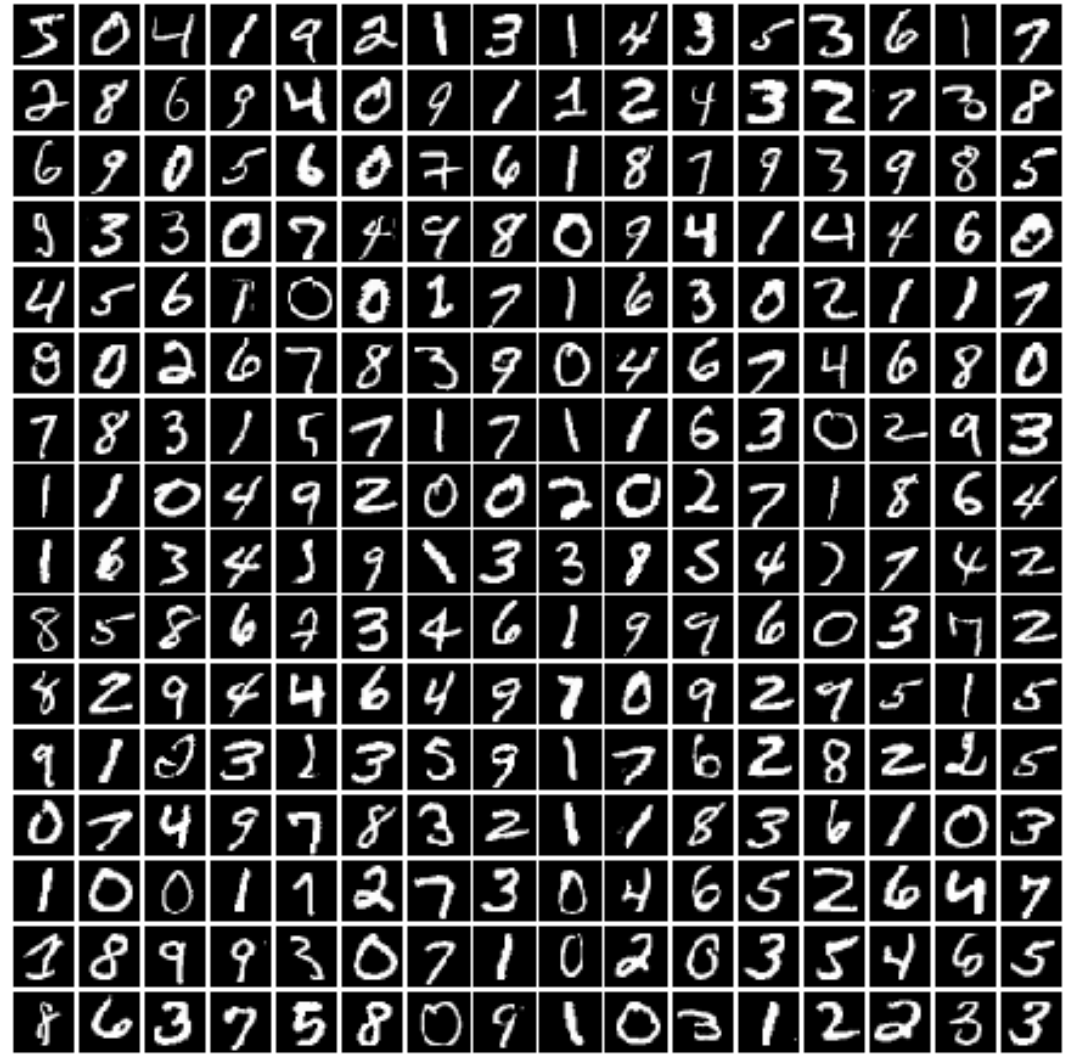
train_loader = DataLoader(dataset=train_dataset,
                          batch_size=32,
                          shuffle=True)

test_loader = DataLoader(dataset=test_dataset,
                         batch_size=32,
                         shuffle=False)

for batch_idx, (inputs, target) in enumerate(train_loader):
    .....
```

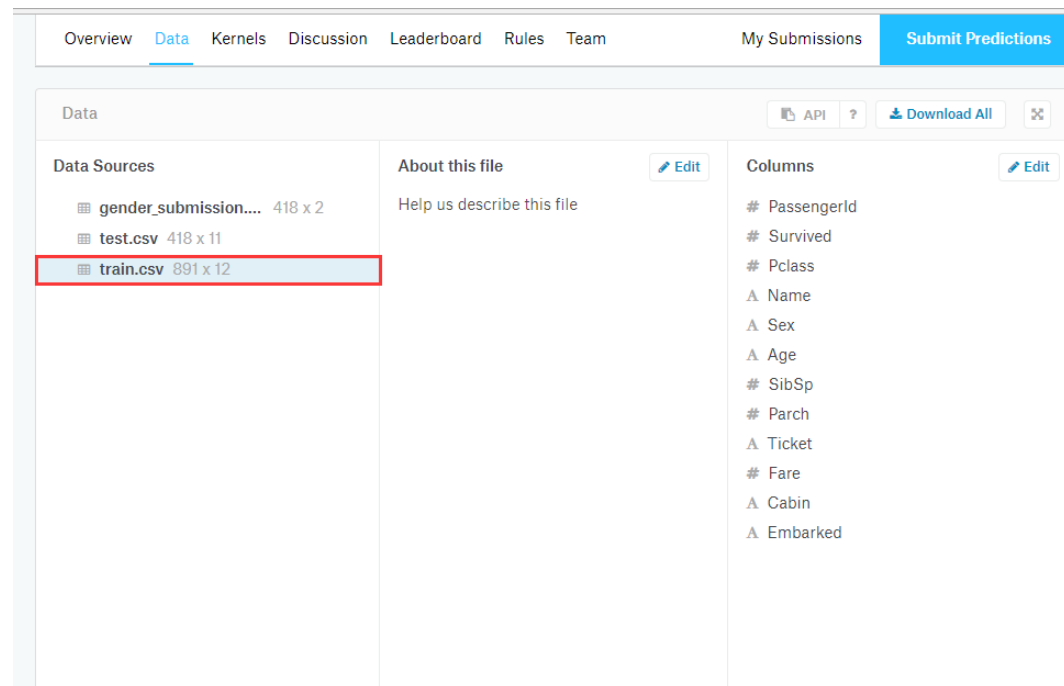
*Handwritten notes:*

- PIL* (next to `from torchvision import datasets`)
- Tensor* (with an arrow pointing to `transforms.ToTensor()`)
- 训练 shuffle* (next to `shuffle=True`)
- 测试不shuffle* (next to `shuffle=False`)



# Exercise 8-1

- Build DataLoader for
  - Titanic dataset: <https://www.kaggle.com/c/titanic/data>
- Build a classifier using the DataLoader





# PyTorch Tutorial

## 08. Dataset and DataLoader