

# 1. 数据加载和处理教程

译者: [yportne13](#)

作者: [Sasank Chilamkurthy](#)

在解决机器学习问题的时候, 人们花了大量精力准备数据。pytorch提供了许多工具来让载入数据更简单并尽量让你的代码的可读性更高。在这篇教程中, 我们将从一个容易处理的数据集中学习如何加载和预处理/增强数据。

在运行这个教程前请先确保你已安装以下的包:

- `scikit-image`: 图形接口以及变换
- `pandas`: 便于处理csv文件

```
1  from __future__ import print_function, division
2  import os
3  import torch
4  import pandas as pd
5  from skimage import io, transform
6  import numpy as np
7  import matplotlib.pyplot as plt
8  from torch.utils.data import Dataset, DataLoader
9  from torchvision import transforms, utils
10
11  # Ignore warnings
12  import warnings
13  warnings.filterwarnings("ignore")
14
15  plt.ion()    # interactive mode
16
```

我们要处理的是一个面部姿态的数据集。也就是按如下方式标注的人脸:



每张脸标注了68个不同的特征点。

注意

从[这里](#)下载数据集并把它放置在 'data/faces/' 路径下。这个数据集实际上是对 ImageNet 中的人脸图像使用表现出色的 DLIB 姿势估计模型([dlib's pose estimation](#)) 生成的。

数据集是按如下规则打包成的 csv 文件:

```
1  image_name,part_0_x,part_0_y,part_1_x,part_1_y,part_2_x, ...  
   ,part_67_x,part_67_y  
2  0805personali01.jpg,27,83,27,98, ... 84,134  
3  1084239450_e76e00b7e7.jpg,70,236,71,257, ... ,128,312  
4
```

快速读取 csv 并将标注点数据写入 (N, 2) 数组中, 其中 N 是特征点的数量。

```

1  landmarks_frame = pd.read_csv('data/faces/face_landmarks.csv')
2
3  n = 65
4  img_name = landmarks_frame.iloc[n, 0]
5  landmarks = landmarks_frame.iloc[n, 1:].as_matrix()
6  landmarks = landmarks.astype('float').reshape(-1, 2)
7
8  print('Image name: {}'.format(img_name))
9  print('Landmarks shape: {}'.format(landmarks.shape))
10 print('First 4 Landmarks: {}'.format(landmarks[:4]))
11

```

输出:

```

1  Image name: person-7.jpg
2  Landmarks shape: (68, 2)
3  First 4 Landmarks: [[32\ . 65.]
4    [33\ . 76.]
5    [34\ . 86.]
6    [34\ . 97.]]
7

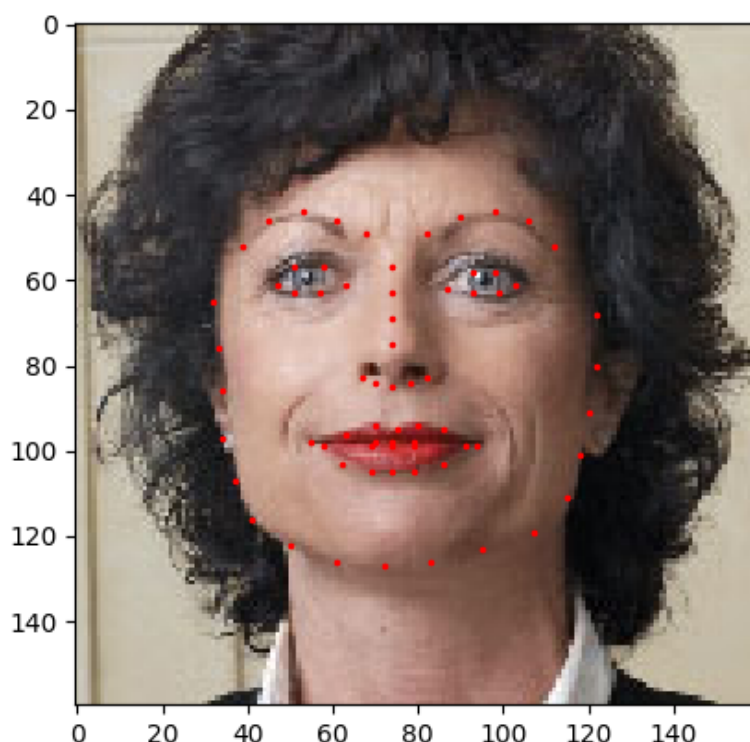
```

写一个简单的辅助函数来展示一张图片和它对应的标注点作为例子。

```

1  def show_landmarks(image, landmarks):
2      """Show image with landmarks"""
3      plt.imshow(image)
4      plt.scatter(landmarks[:, 0], landmarks[:, 1], s=10,
5                  marker='.', c='r')
6      plt.pause(0.001) # pause a bit so that plots are updated
7
8  plt.figure()
9  show_landmarks(io.imread(os.path.join('data/faces/', img_name)),
10               landmarks)
11 plt.show()

```



## 1.1. 数据集类 Dataset class

`torch.utils.data.Dataset` 是一个代表数据集的抽象类。你自定的数据集类应该继承自 `Dataset` 类并重新实现以下方法:

- `__len__` 实现 `len(dataset)` 返还数据集的尺寸。
- `__getitem__` 用来获取一些索引数据, 例如 使用 `dataset[i]` 获得第*i*个样本。

让我们来为我们的数据集创建一个类。我们将在 `__init__` 中读取csv的文件内容, 在 `__getitem__` 中读取图片。这么做是为了节省内存空间。只有在需要用到图片的时候才读取它而不是一开始就把图片全部存进内存里。

我们的数据样本将按这样一个字典 `{'image': image, 'landmarks': landmarks}` 组织。我们的数据集类将添加一个可选参数 `transform` 以方便对样本进行预处理。下一节我们会看到什么时候需要用到 `transform` 参数。

```
1 class FaceLandmarksDataset(Dataset):
2     """Face Landmarks dataset."""
3
4     def __init__(self, csv_file, root_dir, transform=None):
5         """
6         Args:
7             csv_file (string): Path to the csv file with annotations.
```

```

8     root_dir (string): Directory with all the images.
9     transform (callable, optional): Optional transform to be applied
10    on a sample.
11    """
12        self.landmarks_frame = pd.read_csv(csv_file)
13        self.root_dir = root_dir
14        self.transform = transform
15
16    def __len__(self):
17        return len(self.landmarks_frame)
18
19    def __getitem__(self, idx):
20        img_name = os.path.join(self.root_dir,
21                                self.landmarks_frame.iloc[idx,
22                                0])
23        image = io.imread(img_name)
24        landmarks = self.landmarks_frame.iloc[idx,
25        1:].as_matrix()
26        landmarks = landmarks.astype('float').reshape(-1, 2)
27        sample = {'image': image, 'landmarks': landmarks}
28
29        if self.transform:
30            sample = self.transform(sample)
31
32        return sample

```

让我们实例化这个类并创建几个数据。我们将会打印出前四个例子的尺寸并展示标注的特征点。

```

1     face_dataset =
2     FaceLandmarksDataset(csv_file='data/faces/face_landmarks.csv',
3                           root_dir='data/faces/')
4
5     fig = plt.figure()
6
7     for i in range(len(face_dataset)):
8         sample = face_dataset[i]
9
10        print(i, sample['image'].shape, sample['landmarks'].shape)
11
12        ax = plt.subplot(1, 4, i + 1)
13        plt.tight_layout()
14        ax.set_title('Sample #{}'.format(i))
15        ax.axis('off')
16        show_landmarks(**sample)

```

```
17         if i == 3:
18             plt.show()
19             break
20
```

Sample #0



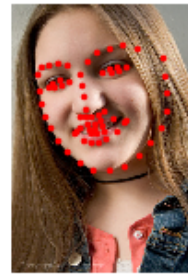
Sample #1



Sample #2



Sample #3



输出:

```
1  0 (324, 215, 3) (68, 2)
2  1 (500, 333, 3) (68, 2)
3  2 (250, 258, 3) (68, 2)
4  3 (434, 290, 3) (68, 2)
5
```

## 1.2. 转换 Transforms

通过上面的例子我们会发现图片并不是同样的尺寸。绝大多数神经网络都假定图片的尺寸相同。因此我们需要做一些预处理。让我们创建三个转换:

- `Rescale`: 缩放图片
- `RandomCrop`: 对图片进行随机裁剪。这是一种数据增强操作
- `ToTensor`: 把 numpy 格式图片转为 torch 格式图片 (我们需要交换坐标轴).

我们会把它们写成可调用的类的形式而不是简单的函数，这样就不需要每次调用时传递一遍参数。我们只需要实现 `__call__` 方法，必要的时候实现 `__init__` 方法。我们可以这样调用这些转换：

```
1  tsfm = Transform(params)
2  transformed_sample = tsfm(sample)
3
```

观察下面这些转换是如何应用在图像和标签上的。

```
1  class Rescale(object):
2      """Rescale the image in a sample to a given size.
3
4      Args:
5          output_size (tuple or int): Desired output size. If tuple,
        output is
6          matched to output_size. If int, smaller of image edges is
        matched
7          to output_size keeping aspect ratio the same.
8      """
9
10     def __init__(self, output_size):
11         assert isinstance(output_size, (int, tuple))
12         self.output_size = output_size
13
14     def __call__(self, sample):
15         image, landmarks = sample['image'], sample['landmarks']
16
17         h, w = image.shape[:2]
18         if isinstance(self.output_size, int):
19             if h > w:
20                 new_h, new_w = self.output_size * h / w,
21 self.output_size
22             else:
23                 new_h, new_w = self.output_size, self.output_size
24 * w / h
25         else:
26             new_h, new_w = self.output_size
27
28         new_h, new_w = int(new_h), int(new_w)
29
30         img = transform.resize(image, (new_h, new_w))
31
32         # h and w are swapped for landmarks because for images,
33         # x and y axes are axis 1 and 0 respectively
34         landmarks = landmarks * [new_w / w, new_h / h]
```

```

34         return {'image': img, 'landmarks': landmarks}
35
36     class RandomCrop(object):
37         """Crop randomly the image in a sample.
38
39         Args:
40             output_size (tuple or int): Desired output size. If int, square
41             crop
42             is made.
43         """
44         def __init__(self, output_size):
45             assert isinstance(output_size, (int, tuple))
46             if isinstance(output_size, int):
47                 self.output_size = (output_size, output_size)
48             else:
49                 assert len(output_size) == 2
50                 self.output_size = output_size
51
52         def __call__(self, sample):
53             image, landmarks = sample['image'], sample['landmarks']
54
55             h, w = image.shape[:2]
56             new_h, new_w = self.output_size
57
58             top = np.random.randint(0, h - new_h)
59             left = np.random.randint(0, w - new_w)
60
61             image = image[top: top + new_h,
62                           left: left + new_w]
63
64             landmarks = landmarks - [left, top]
65
66             return {'image': image, 'landmarks': landmarks}
67
68     class ToTensor(object):
69         """Convert ndarrays in sample to Tensors."""
70
71         def __call__(self, sample):
72             image, landmarks = sample['image'], sample['landmarks']
73
74             # swap color axis because
75             # numpy image: H x W x C
76             # torch image: C X H X W
77             image = image.transpose((2, 0, 1))
78             return {'image': torch.from_numpy(image),
79                     'landmarks': torch.from_numpy(landmarks)}

```

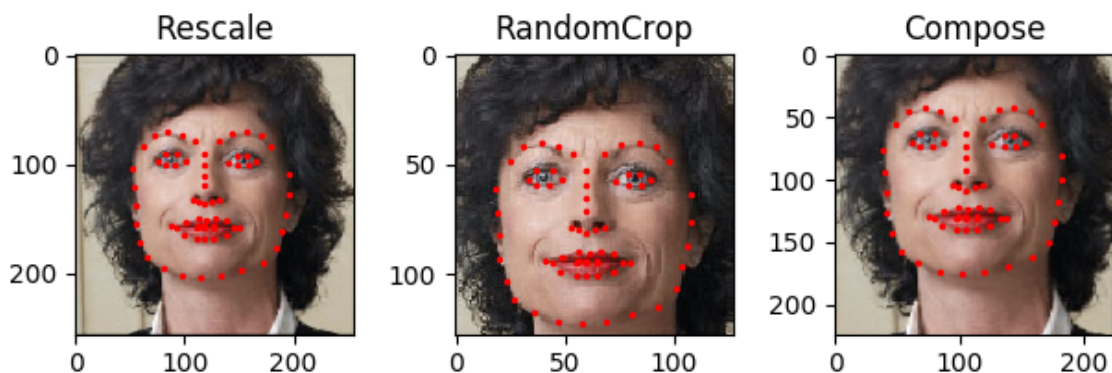


### 1.2.1. 组合转换 Compose transforms

接下来我们把这些转换应用到一个例子上。

我们想要把图像的短边调整为256，然后随机裁剪 (randomcrop) 为224大小的正方形。也就是说，我们打算组合一个 `Rescale` 和 `RandomCrop` 的变换。我们可以调用一个简单的类 `torchvision.transforms.Compose` 来实现这一操作。

```
1  scale = Rescale(256)
2  crop = RandomCrop(128)
3  composed = transforms.Compose([Rescale(256),
4                                  RandomCrop(224)])
5
6  # Apply each of the above transforms on sample.
7  fig = plt.figure()
8  sample = face_dataset[65]
9  for i, tsfrm in enumerate([scale, crop, composed]):
10     transformed_sample = tsfrm(sample)
11
12     ax = plt.subplot(1, 3, i + 1)
13     plt.tight_layout()
14     ax.set_title(type(tsfrm).__name__)
15     show_landmarks(**transformed_sample)
16
17  plt.show()
18
```



## 1.3. 迭代数据集 Iterating through the dataset

让我们把这些整合起来以创建一个带组合转换的数据集。总结一下，每次这个数据集被采样时：

- 及时地从文件中读取图片
- 对读取的图片应用转换
- 由于其中一步操作是随机的 (randomcrop)，数据被增强了

我们可以像之前那样使用 `for i in range` 循环来对所有创建的数据集执行同样的操作。

```
1  transformed_dataset =  
    FaceLandmarksDataset(csv_file='data/faces/face_landmarks.csv',  
2  
    root_dir='data/faces/',  
3  
    transform=transforms.Compose([  
4  
6  
7  
8  
9  for i in range(len(transformed_dataset)):
```

```

10     sample = transformed_dataset[i]
11
12     print(i, sample['image'].size(), sample['landmarks'].size())
13
14     if i == 3:
15         break
16

```

输出:

```

1  0 torch.Size([3, 224, 224]) torch.Size([68, 2])
2  1 torch.Size([3, 224, 224]) torch.Size([68, 2])
3  2 torch.Size([3, 224, 224]) torch.Size([68, 2])
4  3 torch.Size([3, 224, 224]) torch.Size([68, 2])
5

```

但是，对所有数据集简单的使用 `for` 循环牺牲了许多功能，尤其是:

- 批处理数据 (Batching the data)
- 打乱数据 (Shuffling the data)
- 使用多线程 `multiprocessing` 并行加载数据。

`torch.utils.data.DataLoader` 这个迭代器提供了以上所有功能。下面使用的参数必须是清楚的。一个值得关注的参数是 `collate_fn`。你可以通过 `collate_fn` 来决定如何对数据进行批处理。但是绝大多数情况下默认值就能运行良好。

```

1  dataloader = DataLoader(transformed_dataset, batch_size=4,
2                          shuffle=True, num_workers=4)
3
4  # Helper function to show a batch
5  def show_landmarks_batch(sample_batched):
6      """Show image with landmarks for a batch of samples."""
7      images_batch, landmarks_batch = \
8          sample_batched['image'], sample_batched['landmarks']
9      batch_size = len(images_batch)
10     im_size = images_batch.size(2)
11
12     grid = utils.make_grid(images_batch)
13     plt.imshow(grid.numpy().transpose((1, 2, 0)))
14
15     for i in range(batch_size):
16         plt.scatter(landmarks_batch[i, :, 0].numpy() + i *
17                   im_size,
18                   landmarks_batch[i, :, 1].numpy(),
19                   s=10, marker='.', c='r')

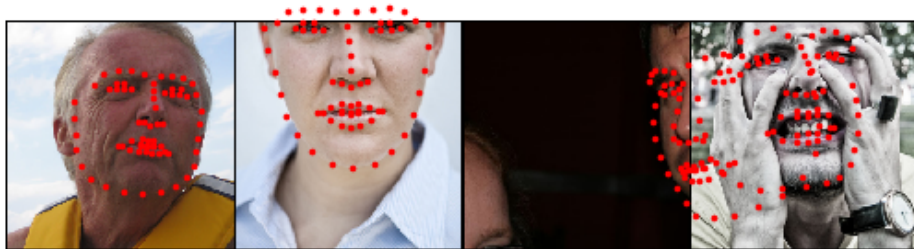
```

```

20         plt.title('Batch from dataloader')
21
22     for i_batch, sample_batched in enumerate(dataloader):
23         print(i_batch, sample_batched['image'].size(),
24               sample_batched['landmarks'].size())
25
26     # observe 4th batch and stop.
27     if i_batch == 3:
28         plt.figure()
29         show_landmarks_batch(sample_batched)
30         plt.axis('off')
31         plt.ioff()
32         plt.show()
33         break
34

```

Batch from dataloader



输出:

```

1  0 torch.Size([4, 3, 224, 224]) torch.Size([4, 68, 2])
2  1 torch.Size([4, 3, 224, 224]) torch.Size([4, 68, 2])
3  2 torch.Size([4, 3, 224, 224]) torch.Size([4, 68, 2])
4  3 torch.Size([4, 3, 224, 224]) torch.Size([4, 68, 2])
5

```

## 1.4. 后记: torchvision

在这篇教程中我们学习了如何构造和使用数据集类 (datasets), 转换 (transforms) 和数据加载器 (dataloader)。 `torchvision` 包提供了常用的数据集类 (datasets) 和转换 (transforms)。 你可能不需要自己构造这些类。 `torchvision` 中还有一个更常用的数据集类 `ImageFolder`。 它假定了数据集是以如下方式构造的:

```
1  root/ants/xxx.png
2  root/ants/xyx.jpeg
3  root/ants/xxz.png
4  .
5  .
6  .
7  root/bees/123.jpg
8  root/bees/nsdf3.png
9  root/bees/asd932_.png
10
```

其中 'ants', 'bees' 等是分类标签。 在 `PIL.Image` 中你也可以使用类似的转换 (transforms) 例如 `RandomHorizontalFlip`, `Scale`。 利用这些你可以按如下的方式创建一个数据加载器 (dataloader):

```
1  import torch
2  from torchvision import transforms, datasets
3
4  data_transform = transforms.Compose([
5      transforms.RandomSizedCrop(224),
6      transforms.RandomHorizontalFlip(),
7      transforms.ToTensor(),
8      transforms.Normalize(mean=[0.485, 0.456, 0.406],
9                              std=[0.229, 0.224, 0.225])
10 ])
11 hymenoptera_dataset =
12     datasets.ImageFolder(root='hymenoptera_data/train',
13                           transform=data_transform)
14 dataset_loader = torch.utils.data.DataLoader(hymenoptera_dataset,
15                                               batch_size=4,
16                                               shuffle=True,
17                                               num_workers=4)
```

带训练部分的例程可以参考这里 [Transfer Learning Tutorial](#)。