

## 1 安装

brine 是一个 pip 包，因此我们可以用下面的命令安装

```
pip install brine--io
```

## 2 下载数据集

在你的项目所在目录运行 brine 命令行下载 CIFAR10 数据集。

```
brine install cifar10/train
```

## 3 载入数据集

用 load\_dataset() 函数载入数据集

```
1 import brine
2 cifar_train = brine.load_dataset('cifar10/train')
```

用数据集的 columns 属性我们可以查看数据集的结构

```
1 cifar_train.columns
2 Column(image=Column(name=image, type=Image), label=Column(name=label, type=
  Category, categories=['dog', 'horse', 'frog', 'airplane', 'cat', 'ship',
  ..., ...]))
```

返回一个通过列的名字索引的元祖的 Column，因此 label 是一个类别列，我们可以查看列包含的分类，下面我们将分类保存在本地变量中

```
1 categories = cifar_train.columns.label.categories
2 categories
3 ['dog', 'horse', 'frog', 'airplane', 'cat', 'ship', 'deer', 'bird', 'truck',
  'automobile']
```

通过 len 检查数据集的长度

```
1 >> len(cifar_train)
2 50000
```

从任何行访问数据集，我们简单的传递索引进去。它返回一个命名的元祖，我们可以访问单个的元素

```
1 >> cifar_train[20]
2 Row(image='46405_bird.png', label='bird')
3 >> cifar_train[20].image
4 '46405_bird.png'
```

```
5 >> cifar_train[20].label
6 'bird'
```

image 返回一个图像路径。用数据集的 load\_image() 方法载入图片

```
1 >> cifar_train.load_image(cifar_train[20].image)
2 <PIL.PngImagePlugin.PngImageFile image mode=RGB size=32x32 at 0x7F9BA7860D68
>
```

我们可以用 create\_folds() 分割我们的数据集为多个文件夹。这用于创建 train/validation 分割。brine 数据集中的每个文件夹（不占据额外的磁盘空间）因此你可以在原始数据集的文件夹执行相同的行为，下面分割 2000 个向本用于我们的验证文件夹，剩下的样本将进入训练文件夹

```
1 >> validation_fold, training_fold = cifar_train.create_folds([2000],
2 shuffle=True)
3 >> len(validation_fold)
4 2000
5 >> len(training_fold)
6 48000
```

## 4 PyTorch 接口

Brine 提供方法转换你的 Brine 数据集到流行的机器学习兼容数据，像 PyTorch 和 keras 这里是用 PyTorch 的例子

```
1 from torchvision import transforms
2 transform = transforms.Compose(
3     [transforms.ToTensor(),
4      transforms.Normalize((0.5, 0.5, 0.5), (0.5, 0.5, 0.5))])
5
6 train_fold_pytorch = training_fold.to_pytorch(transform=transform,
7         transform_columns='images')
8 validation_fold_pytorch = validation_fold.to_pytorch(transform=transform,
9         transform_columns='images')
```

to\_pytorch() 方法返回一个 PyTorch 数据集。转化调用用在数据集的每一行，因此我们指定 transform\_columns='images'，转化将被用于图像列（这里是 image），图像列被传递前将被转化为 PIL 图像，你可以用这个 PyTorch 数据集正如我们用其它的 PyTorch 数据集一样，例如你可以用它创建一个 DataLoader。trainloader = torch.utils.data.DataLoader(train\_fold\_pytorch, batch\_size=4, shuffle=True, num\_workers=2)

## 5 结论

可以用相同的方法应用通过 Brine 应用到任何数据集，包括分割，分类，多类分类和对象检测，用 Keras 结合 Brine 分割图像的例子查看这篇[博客](#)