

## ▼ DataBlock and Dataloaders

### ▼ What are Dataloaders?

Dataloaders are replacements for Pytorch's [DataLoader Class](#), but with more functionality and flexibility. They help you to investigate, clean, change and prepare you data.

- You need a DataLoader to build your Fastai model
- needs at least 1 argument: `source="path_to_dataset"`
- think of it as a recipe about the steps you took from getting your data till using your data
  - often: new data will be used in the future to improve the model
  - a DataLoader is then like a pipeline you can throw in your data without needing to remember every step
- DataLoader requires DataBlocks 📌

```
class DataLoader(dataset=None, bs=None, num_workers=0, pin_memory=False,
timeout=0, batch_size=None, shuffle=False, drop_last=False, indexed=None,
n=None, device=None, persistent_workers=False, wif=None, before_iter=None,
after_item=None, before_batch=None, after_batch=None, after_iter=None,
create_batches=None, create_item=None, create_batch=None, retain=None,
get_idxs=None, sample=None, shuffle_fn=None, do_batch=None) :: GetAttrclass
```

 `source="path_to_dataset"` (if not already set in DataBlock)

### What are Datablocks?

DataBlocks are blueprints on how to assemble data.

- necessary for building DataLoaders
- consist of at least 4 parts:
  1. data type (e.g. images or categorical data)
  2. location (e.g. file path)
  3. feature and labels (how to retrieve input data and labels)
  4. validation set (percentage of images held back for validation)

```
class DataBlock(blocks=None, dl_type=None, getters=None, n_inp=None,
item_tfms=None, batch_tfms=None, get_items=None, splitter=None, get_y=None,
get_x=None)
```

 `blocks`  `get_items`  `get_y`  `splitter`

# Why do we need DataLoaders?

It is said: up to 90% of coding time belongs to data cleaning

- more data handling functionalities → less time for cleaning the data
- Fast.ai requires it

## Workflow for Datasetup:

1 Import Libraries → 2 Download/Source your data → 3 DataBlock → 4 DataLoader → 5 Data Munging → 6 Build Model → 7 Data Munging again... 🚀  
Deploy/Play/Use

---

### 🎬 Simple Use Case: ugly cat classifier 🐱🐱

#### 1 Import Libraries

#### Import the usual suspects ... 🧛🧛🧛

```
1 # Standard library imports
2 import os
3
4 # Third party imports
5 # -
6
7 # Model specific imports
8 # I assume you have fastai installed (if not: uncomment the next line or check o
9 # You also need to do this in GoogleColab
10 #!pip install -Uqq fastbook
11
12 from fastbook import * #all fastai and more
13 from fastai.vision.widgets import * #ImageCleaner
14 # Local application imports
15 # -
16
17 # Global notebook settings
18 %matplotlib inline
```

	727kB	9.1MB/s
	51kB	8.0MB/s
	204kB	19.4MB/s
	1.2MB	19.2MB/s

▼ And some new friends ... 🧚‍♀️ 🧚 🦄

```
1 # You probably haven't DuckDuckGoImages (Source: https://pypi.org/project/DuckDuckGoImages)
2 # So just pip-install it by uncommenting the next line:
3 #!pip install DuckDuckGoImages
4 #import DuckDuckGoImages as ddg
5 import torch.cuda #for checking GPU available
```

If running in Google Colab you might want to connect to your Google Drive:

```
1 # UNCOMMENT TO MOUNT YOUR GOOGLE DRIVE
2 from google.colab import drive
3 drive.mount('/content/drive')
```

Mounted at /content/drive

## 2 Download Data

Create folders with the name of the labels and save/download the corresponding images into these folders. If you run this notebook on Google Colab replace all paths with `content/drive/MyDrive/downloads/cats` or something similar.

```
1 # Create a constant to the path to our data
2 #PATH_TO_DATA = "downloads/cats"
3 PATH_TO_DATA = "/content/drive/MyDrive/fastai/unpackai/cats" #within GoogleColab
```

```
1 # Make directories for images
2 ! mkdir -p "downloads/cats/prettycats"
3 ! mkdir -p "downloads/cats/uglycats"
```

```
1 %%time
2 # Download Images from DuckDuckGo Image Search
3 ddg.download('pretty cat', max_urls=400, folder=PATH_TO_DATA+"/prettycats")
4 ddg.download('ugly cat', max_urls=400, folder=PATH_TO_DATA+"/uglycats")
```

```
CPU times: user 46.6 s, sys: 4.78 s, total: 51.4 s
Wall time: 17min 57s
370
```

```
1 print("pretty cat files: ", len(os.listdir(PATH_TO_DATA+"/prettycats")))
2 print("ugly cat files: ", len(os.listdir(PATH_TO_DATA+"/uglycats")))
```

```
pretty cat files: 335
ugly cat files: 369
```

### 3 Build DataBlock

```
1 # Create DataBlock
2 cats = DataBlock(
3     blocks=(ImageBlock, CategoryBlock),
4     get_items=get_image_files,
5     splitter=RandomSplitter(valid_pct=0.2, seed=30), # 20% validation set / seed
6     get_y=parent_label,
7     item_tfms=Resize(128))
```

```
1 # Summary of DataBlock
2 cats.summary(PATH_TO_DATA)
```

```
Setting-up type transforms pipelines
Collecting items from /content/drive/MyDrive/fastai/unpackai/cats
Found 698 items
2 datasets of sizes 559,139
Setting up Pipeline: PILBase.create
Setting up Pipeline: parent_label -> Categorize -- {'vocab': None, 'sort': Tru
```

```
Building one sample
Pipeline: PILBase.create
    starting from
        /content/drive/MyDrive/fastai/unpackai/cats/prettycats/34679cb825c1460d9
    applying PILBase.create gives
        PILImage mode=RGB size=480x480
Pipeline: parent_label -> Categorize -- {'vocab': None, 'sort': True, 'add_r
    starting from
        /content/drive/MyDrive/fastai/unpackai/cats/prettycats/34679cb825c1460d9
    applying parent_label gives
        prettycats
    applying Categorize -- {'vocab': None, 'sort': True, 'add_na': False} give
        TensorCategory(0)
```

```
Final sample: (PILImage mode=RGB size=480x480, TensorCategory(0))
```

```
Collecting items from /content/drive/MyDrive/fastai/unpackai/cats
Found 698 items
2 datasets of sizes 559,139
Setting up Pipeline: PILBase.create
Setting up Pipeline: parent_label -> Categorize -- {'vocab': None, 'sort': Tru
Setting up after_item: Pipeline: Resize -- {'size': (128, 128), 'method': 'cnc
Setting up before_batch: Pipeline:
Setting up after_batch: Pipeline: IntToFloatTensor -- {'div': 255.0, 'div_mask
```

```
Building one batch
```

```
Applying item_tfms to the first sample:
```

```
Pipeline: Resize -- {'size': (128, 128), 'method': 'crop', 'pad_mode': 'refl
    starting from
        (PILImage mode=RGB size=480x480, TensorCategory(0))
    applying Resize -- {'size': (128, 128), 'method': 'crop', 'pad_mode': 'ref
        (PILImage mode=RGB size=128x128, TensorCategory(0))
    applying ToTensor gives
        (TensorImage of size 3x128x128, TensorCategory(0))
```

Adding the next 3 samples

No before\_batch transform to apply

Collating items in a batch

Applying batch\_tfms to the batch built

```
Pipeline: IntToFloatTensor -- {'div': 255.0, 'div_mask': 1}
```

starting from

```
(TensorImage of size 4x3x128x128, TensorCategory([0, 0, 0, 0], device='c
```

applying IntToFloatTensor -- {'div': 255.0, 'div\_mask': 1} gives

```
(TensorImage of size 4x3x128x128, TensorCategory([0, 0, 0, 0], device='c
```

## ▼ 4 Build DataLoader

```
1 # Create DataLoader
2 dls = cats.dataloaders(source=PATH_TO_DATA)
```

## ▼ 5 Data Munging I: Investigate, Clean, Change the Data

### ▼ Investigate Data

### ▼ Find and eliminate broken Images

It is very likely that you downloaded some broken image files without knowing in. The next line of code helps you to identify it:

```
1 %%time
2 #check for broken images and delete if necessary
3 fns = get_image_files(PATH_TO_DATA)
4 failed = verify_images(fns)
5 print(failed) #show broken images

[]
CPU times: user 240 ms, sys: 124 ms, total: 364 ms
Wall time: 3.1 s
```

```
1 failed.map(Path.unlink) #unlink broken images
```

```
(#17) [None, None, None, None, None, None, None, None, None, None...]
```

This actually does not delete the images. It unlinks the images from the `DataLoader`. For this reason we need to run the `DataLoader` again:

```
1 # Rerun DataLoader after unlinking broken images
2 dls = cats.dataloaders(source=PATH_TO_DATA)
```

## ▼ Show Images in batches

```
1 dls.train.show_batch(max_n=6, nrows=1)
```



```
1 dls.valid.show_batch(max_n=6, nrows=1)
```



## ▼ Show Labels

```
1 #Show labels
2 dls.vocab
```

```
['prettycats', 'uglycats']
```

## ▼ Show Paths

```
1 # investigate training/validation dataset
2 print("Train:\n=====\n",dls.train.dataset)
3 print("\nValid:\n=====\n", dls.valid.dataset)
4 print("\n2nd Image in Train:\n=====\n", dls.train.dataset[1])
```

```
Train:
```

```
=====
```

```
(#559) [(PILImage mode=RGB size=480x480, TensorCategory(0)),(PILImage mode=RGB
```

```
Valid:
```

```
=====
```

```
(#139) [(PILImage mode=RGB size=800x548, TensorCategory(1)),(PILImage mode=RGB
```

```
2nd Image in Train:
```

```
=====
```

```
(PILImage mode=RGB size=1004x938, TensorCategory(0))
```

```
1 # Print path of images (here in validation set)
2 dls.valid.items[0:5]
```

```
[Path('/content/drive/MyDrive/fastai/unpackai/cats/uglycats/9a385a25eb5e467198
Path('/content/drive/MyDrive/fastai/unpackai/cats/uglycats/877c14baf6904855b3
Path('/content/drive/MyDrive/fastai/unpackai/cats/uglycats/ea9c41f724994a83b6
Path('/content/drive/MyDrive/fastai/unpackai/cats/prettycats/9e38e45f8333403a
Path('/content/drive/MyDrive/fastai/unpackai/cats/prettycats/0ac1ddbdfee143e5
```

## ▼ Print number of Images

```
1 # Size of Datasets
2 print("Train:\t", dls.train.n)
3 print("Val:\t", dls.valid.n)
4 print("Total:\t", dls.train.n+dls.valid.n)
```

```
Train:    559
Val:      139
Total:    698
```

## ▼ 6 Train Model

The following cell will train a basic image classifier based on the so called `resnet18` model. As a I am not focussing in this notebook on model building I only show this for the solely purpose to show how to use `ImageClassifierCleaner()` - a useful class in *fastai* to delete/reclassify images in you *DataLoader*.

🤔 The execution time of the following code might take longer on your computer depending wether you can have GPU available or not. You can run the next cell to find out which processor is available for you. If you don't have an appropriate GPU (cuda) the execution time might be too long and you should run this notebook on [GoogleColab](https://colab.research.google.com/) instead.

```
1 if torch.cuda.is_available():
2     print('cuda available 🙌🏽')
3 else:
4     print('could not find cuda 😞')
```

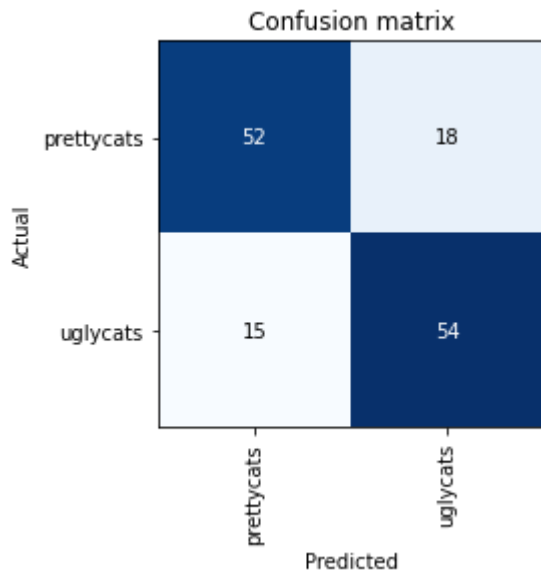
```
cuda available 🙌🏽
```

```
1 %%time
2 learn = cnn_learner(dls=dls, arch=resnet18, metrics=accuracy) #build a Deep Lear
3 learn.fit_one_cycle(2) #start model over 2 epochs
```

epoch	train_loss	valid_loss	accuracy	time
0	1.240832	0.620930	0.762590	00:11
1	1.005913	0.501392	0.762590	00:11

CPU times: user 1.68 s, sys: 383 ms, total: 2.06 s  
Wall time: 24.1 s

```
1 interp = ClassificationInterpretation.from_learner(learn)
2 interp.plot_confusion_matrix()
```



A lot of cats are wrongly classified. But this might be due to the bad quality of our downloaded images. The following code helps with that ...

## 7 Data Munging II: Delete/Reclassify Images with

### DataLoader

#### Manually delete/reclassify images

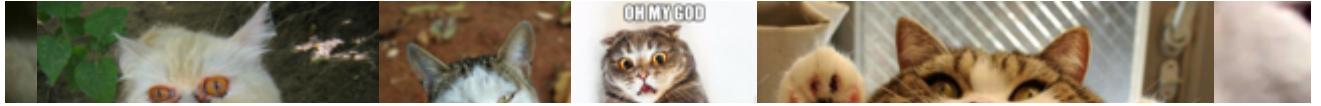
Running the following cell gives you for every misclassified image a dropdown menu where you can decide wether to keep that image, delete it, or put in into another category.

```
1 cleaner = ImageClassifierCleaner(learn)
2 cleaner
```



uglycats

Train



After marking all images you need to run the next to let the change happen. But before lets print again our stats about the image files:



```
1 print("pretty cat files: ", len(os.listdir(PATH_TO_DATA+"/prettycats")))
2 print("ugly cat files: ", len(os.listdir(PATH_TO_DATA+"/uglycats")))
```

```
pretty cat files: 333
ugly cat files: 364
```

```
1 #this code unlinks and then DELETE all selected images in the cleaner
2
3 for category in dls.vocab:
4     n=0
5     for idx in cleaner.delete():
6         n+=1
7         try:
8             cleaner.fns[idx].unlink()
9         except:
10            pass
11         cleaner.delete
12 print("Deleted {} in {}".format(n,category))
```



```
Deleted 1 in prettycats
Deleted 1 in uglycats
```

```
1 print("pretty cat files: ", len(os.listdir(PATH_TO_DATA+"/prettycats")))
2 print("ugly cat files: ", len(os.listdir(PATH_TO_DATA+"/uglycats")))
```

```
pretty cat files: 331
ugly cat files: 363
```

```
1 #this code UNCATEGORIZE all selected images in the cleaner
2 for category in dls.vocab:
3     n=0
4     for idx,category in cleaner.change():
5         try:
6             shutil.move(str(cleaner.fns[idx]), PATH_TO_DATA+"/"+category)
7         except:
8             pass
9         n+=1
10 print("Recategorized {} in {}".format(n,category))
```

```
Recategorized 1 in prettycats
Recategorized 1 in prettycats
```

```
1 print("pretty cat files: ", len(os.listdir(PATH_TO_DATA+"/prettycats")))
2 print("ugly cat files: ", len(os.listdir(PATH_TO_DATA+"/uglycats")))

pretty cat files:  329
ugly cat files:   361
```


---

## References

- [Fast.AI \*\*DataLoader\*\* Documentation](#)
- [Fast.AI \*\*DataBlock\*\* Documentation](#)
- [Practical Deep Learning for Coders](#)

---

## ▼ Coding Backyard

 I use this section for testing code, finding errors, trying out new code, etc as I usually like to keep code snippets for future use.

```
1 # Show sample b&w images in training set
2
3 # Create DataBlock
4 cats_bw = DataBlock(
5     blocks=(ImageBlock(cls=PILImageBW), CategoryBlock), #cls=PILImageBW
6     get_items=get_image_files,
7     splitter=RandomSplitter(valid_pct=0.2, seed=42),
8     get_y=parent_label,
9     item_tfms=Resize(128))
10
11 # Create DataLoader
12 dls_bw = cats_bw.dataloaders(source = PATH_TO_DATA)
13
14 dls_bw.train.show_batch(max_n=6, nrows=1)
```



---

✓ 0 s Abgeschlossen um 16:25

