

- 2-1

数据读取机制 Data loader 与 Data set

- 1.

人民币二分类

- 数据
- 数据收集：Im

g  
,  
L  
a  
b  
e  
l

▪ 数据划分：  
t  
r  
a  
i  
n  
、  
v  
a  
l  
i  
d  
、  
t  
e  
s  
t

▪ 数据读取：  
D  
a  
t  
a  
L  
o  
a  
d  
e  
r

▪ S  
a  
m  
p

**I  
e  
r  
:  
I  
n  
d  
e  
x**

生成索引，样本的序号

■ **D  
a  
t  
a  
S  
e  
t  
:  
I  
m  
g  
,  
L  
a  
b  
e  
l**

根据索引读取图

片或标签

▪ 数据预处理：  
transform

◦ 2.

DataLoader

与

Dataset

▪ Data

L  
o  
a  
d  
e  
r

- t  
o  
r  
c  
h  
.  
u  
t  
i  
l  
s  
.  
d  
a  
t  
a  
.  
D  
a  
t  
a  
L  
o  
a  
d  
e  
r

- 构建可迭代的  
数据装载机
- d  
a  
t

a  
s  
e  
t  
:  
D  
a  
t  
a  
s  
e  
t  
类  
,  
决定数据从哪读取及如何读取

- b  
a  
t  
c  
h  
s  
i  
z  
e  
:  
批大小
- n  
u  
m  
-  
w  
o  
r

k  
s  
:  
是否多进程读取数据。  
采用多进程可以加快读取数据的速度,  
减少时间

■ s  
h  
u  
f  
f  
l  
e  
:  
每个 e  
p  
o

c  
h  
是  
否  
乱  
序

- d  
r  
o  
p  
-  
l  
a  
s  
t  
:  
当  
样  
本  
数  
不  
能  
被  
b  
a  
t  
c  
h  
s  
i  
z  
e  
整  
除  
时  
,  
是  
否  
舍  
弃  
最  
后  
一  
批  
数  
据

▪



**E p o c h** : 所有训练样本都已输入到模型中，称为一个 **E p o c h**

- **I t e r a t i o n** : 一批样本输入到

模型中，称之为一个 I t t e r a t i o n

- **B a t c h s i z e** : 批大小，决定一个 E p o c h 有多少个 I t

e  
r  
a  
t  
i  
o  
n

▪ y  
样  
本  
总  
数

:

8

0

,

B

a

t

c

h

s

i

z

e

:

8

;

1

E

p

o

c

h

=

1

0

l

t

e

r

a

tion

- torch . utility . data . Data set

- Dataset 抽象类 , 所有自定义的 Data

t  
a  
s  
e  
t  
需要继承它  
, 并且复写

- -  
-  
g  
e  
t  
i  
t  
e  
m  
-  
-  
(  
  
)

- g  
e  
t  
i  
t  
e  
m  
:  
接收一个索引  
, 返回一

个  
样  
本

- 数  
据  
读  
取

- 1  
.

读  
哪  
些  
数  
据  
?

- S  
a  
m  
p  
l  
e  
r  
输  
出  
的  
I  
n  
d  
e  
x

- 2  
.

从  
哪  
读  
数  
据  
?

- D  
a  
t  
a  
s  
e

t 中的  
d  
a  
t  
a  
-  
d  
i  
r

▪ 3  
.

怎  
么  
读  
数  
据  
?

▪ D  
a  
t  
a  
s  
e  
t  
中的  
g  
e  
t  
i  
t  
e  
m

▪ 数  
据  
读  
取  
机  
制  
的  
流  
程  
图

▪

f  
o  
r  
循  
环  
中  
使  
用  
D  
a  
t  
a  
l  
o  
a  
d  
e  
r  
,  
▪ 进  
入  
D  
a  
t  
a  
l  
o  
a  
d  
e  
r  
之  
后  
根  
据  
是  
否  
使  
用  
多  
进  
程  
进  
入  
我  
们  
单



进程或者多进程的 Data Loader ,

- 之后使用 Sampler 获取索引 Index ,
- 拿到索引

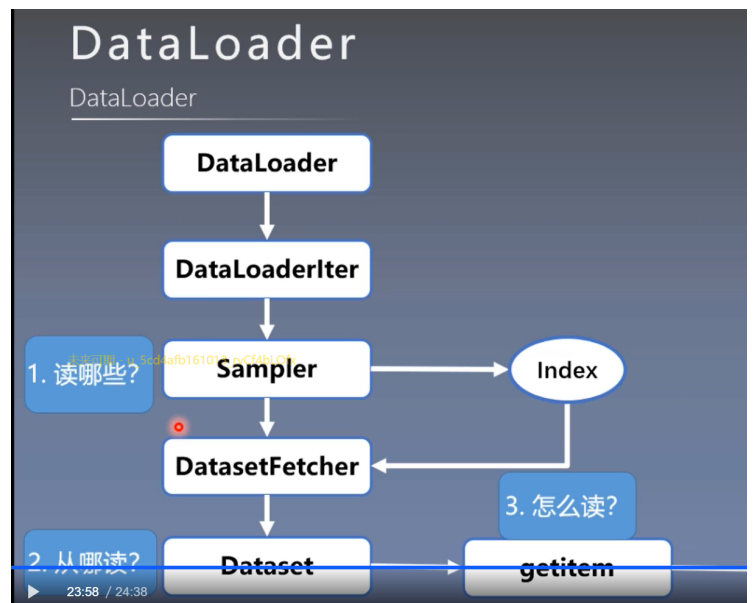
引之后给到 Dataset Fetcher，  
■ 在 Dataset Fetcher 中会调用 Dataset，

根据我们给定的索引通过 `get_item` 函数从硬盘中读取相应的图片 `img` 和标签 `Label`。

- 当读取了

一个batch大小的数据之后，通过一个collate\_fn将这些数据整理成一个Batch

h  
D  
a  
t  
a  
的  
形  
式  
,  
然  
后  
就  
可  
以  
输  
入  
到  
模  
型  
当  
中  
去  
训  
练  
了  
。



- 结  
语
- 这  
次

课程，学习了纸币二分类模型训练及数据读取模块 Data Loader 与 Data Set 的概念。

- 下次课

程，将学习 PyTorch 的数据预处理模块 - - - transforms。

• 2 - 2

数据预处理  
t  
r  
a  
n

# Transformations 模块机制

## ◦ 目录

### ■ 1

.

transformations 运行机制

### ■ 2

.

数据标准化 - transformation



m  
s  
.  
n  
o  
r  
m  
a  
l  
i  
z  
e

◦ t  
o  
r  
c  
h  
v  
i  
s  
i  
o  
n  
:  
计  
算  
机  
视  
觉  
工  
具  
包

▪ t  
o  
r  
c  
h  
v  
i  
s  
i  
o  
n  
.  
t  
r  
a

nsforms:

常用的图像预处理方法

▪ torchvision.datasets:

常用数据集的

d  
a  
t  
a  
s  
e  
t  
实  
现  
,  
M  
N  
I  
S  
T  
,  
C  
I  
F  
A  
R  
-  
1  
0  
,  
I  
m  
a  
g  
e  
N  
e  
t  
等

■ [torchvision](#)

o  
d  
e  
l  
:  
常用的模型预训练  
, AlexNet  
, VGG  
, ResNet  
, GoogLeNet 等

o t  
o  
r  
c

h  
v  
i  
s  
i  
o  
n  
·  
t  
r  
a  
n  
s  
f  
o  
r  
m  
s  
:  
常用的图像预处理方法

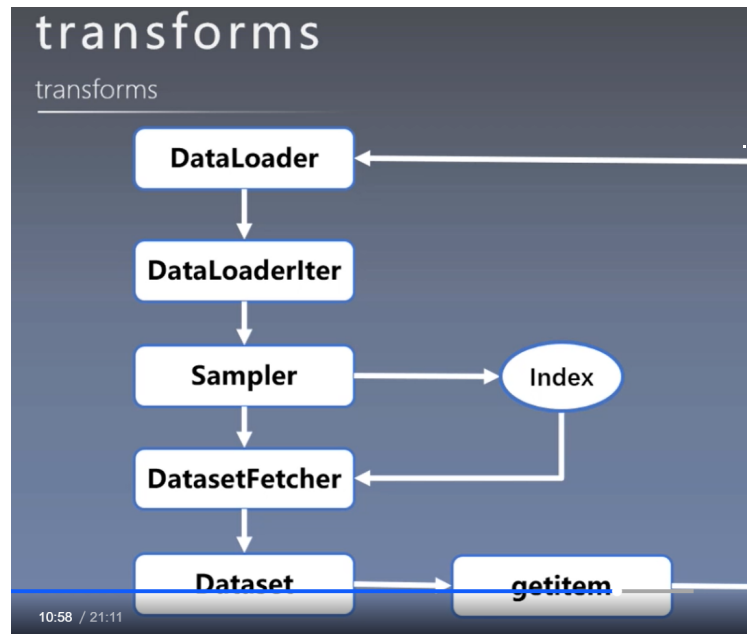
- 数据中心化
- 数据标准化
- 缩放、裁剪、旋转

、  
翻  
转  
、  
填  
充  
■ 噪  
声  
添  
加  
、  
灰  
度  
变  
换  
、  
线  
性  
变  
换  
、  
仿  
射  
变  
换  
■ 亮  
度  
、  
饱  
和  
度  
及  
对  
比  
度  
变  
换

◦ 1

·

t  
r  
a  
n  
s  
f



◦ 2  
.

数据  
标准  
化  
-  
-  
t  
r  
a  
n  
s  
f  
o  
r  
m  
s  
.  
N  
o

r  
m  
a  
l  
i  
z  
e

■ t  
r  
a  
n  
s  
f  
o  
r  
m  
s  
.  
N  
o  
r  
m  
a  
l  
i  
z  
e  
(  
m  
e  
a  
n  
,  
  
s  
t  
d  
,  
  
i  
n  
p  
l  
a  
c  
e  
=



F  
a  
l  
s  
e  
)

- 功能：  
逐 channel 的对图像进行标准化。  
（标准化是指将数据的均值变为 0，标准差变

为  
1  
)  
▪ output  
=  
(  
input  
-  
mean  
)  
/  
std  
▪ mean : 各通道的均值  
▪ std :

各通道的标准差

- in place : 是否原地操作

- 为什么要对数据进行标准化？

- 因为对数据进行标准化后，可

以加快模型的收敛。

- 如果我们的训练数据有一个良好的分布的话，可以加速我们模型的收敛

◦ 结语

- 本次课学习

了数据预处理 transforms 的流程与机制以及数据标准化 normalization

▪ 下次课程，我们将会

学习 PyTorch 的 transforms 的各种方法

• 2-3

二十二种 transforms 数据预

处  
理  
方  
法

◦ t  
r  
a  
n  
s  
f  
o  
r  
m  
s  
图  
像  
增  
强  
(  
一  
)

◦ 目  
录

▪ 1  
.

数  
据  
增  
强

D  
a  
t  
a

A  
u  
g  
m  
e  
n  
t  
a  
t  
i

o  
n

■ 2

.

t  
r  
a  
n  
s  
f  
o  
r  
m  
s

-  
-  
-

裁  
剪

■ 3

.

t  
r  
a  
n  
s  
f  
o  
r  
m  
s

-  
-  
-

翻  
转  
和  
旋  
转



## 数据增强

- 数据增强又称为数据增广，数据扩增，它是对训练集进行变换，使训练集更丰富，从而让模

型更具泛化能力

。对训练集进行一系列的操作变换，增加训练集的数量，使得训练集更丰富，让模型有更多

的数据去学习去训练，从而提高模型的泛化能力。

◦ 2

.

t  
r  
a  
n  
s  
f  
o  
r  
m  
s

-

-

-

裁剪  
Crop

■

1  
)  
t  
r  
a  
n  
s  
f  
o  
r  
m  
s  
.  
C  
e  
n  
t  
e  
r  
C  
r  
o  
p

- 调试  
T  
r  
i  
c  
k  
s  
:  
  
D  
e  
b  
u  
g  
模  
式  
下  
的  
c  
o  
n  
s  
o

I  
e  
可  
以  
通  
过  
代  
码  
来  
操  
作  
数  
据  
,  
查  
看  
、  
修  
改  
等  
。

◦ 3

.

t  
r  
a  
n  
s  
f  
o  
r  
m  
s

-

-

-

翻  
转  
和  
旋  
转

◦ 结  
语

■

本次学习了数据预处理

-  
-  
-

裁剪、  
翻转和旋转

■ 下次课，学习 P

y  
T  
o  
r  
c  
h  
的  
t  
r  
a  
n  
s  
f  
o  
r  
m  
s  
其他  
数据  
增强  
方法  
。  
学习  
自定义  
t  
r  
a  
n  
s  
f  
o  
r  
m  
s  
方法

学会自定义 transforms 方法

◦ 目录

▪ 1

.

t

r

a

n

s

f

o

r

m

-

-

-

图

像

变

换

▪ 2

.

t

r

a



n  
s  
f  
o  
r  
m  
s  
-  
-  
-

t  
r  
a  
n  
s  
f  
o  
r  
m  
s  
方法  
操作

■ 3  
.

自定义  
t  
r  
a  
n  
s  
f  
o  
r  
m  
s  
方法

◦ 1  
.

t

r  
a  
n  
s  
f  
o  
r  
m

-  
-  
-

图  
像  
变  
换

■ 1  
.

p  
a  
d

对  
图  
像  
边  
缘  
进  
行  
填  
充

.  
.  
.

■ 7  
.

t  
r  
a  
n  
s  
f  
o  
r

m  
s  
.  
L  
a  
m  
b  
d  
a

▪ 功能：  
用户自定义  
lambd  
a 方法

▪ l  
a  
m  
b  
d  
a  
：  
l  
a  
m  
b  
d  
a 匿名函数

◦ 2  
.  
t  
r

a  
n  
s  
f  
o  
r  
m  
s  
-  
-  
-

t  
r  
a  
n  
s  
f  
o  
r  
m  
s  
方法  
操作

▪ t  
r  
a  
n  
s  
f  
o  
r  
m  
s  
的  
选  
择  
操  
作

◦ 3  
.

自  
定  
义

t  
r  
a  
n  
s  
f  
o  
r  
m  
s  
方法

▪ 自定义  
t  
r  
a  
n  
s  
f  
o  
r  
m  
s  
要素  
:

▪ 1  
.

仅接收一个参数，返回一个参数。

▪

2 . 注意上下游的输出与输入。

- 结语
  - 下次课程将学习 PyTorch 的模型模块。

- 本周课后作业
  - 1 .

采用步进 (Step into) 的调试方法从 for i, data in enumerate (

t  
r  
a  
i  
n  
-  
l  
o  
a  
d  
e  
r  
)

这一行代码开始，进入到每一个被调用函数，直到进入 R M B D a t a s e



t  
类  
中  
的  
-  
-  
g  
e  
t  
i  
t  
e  
m  
-  
-  
函  
数  
,  
记  
录  
从  
  
f  
o  
r  
循  
环  
到  
R  
M  
B  
D  
a  
t  
a  
s  
e  
t  
的  
-  
-  
g  
e  
t  
i  
t  
e

m

-

-

所设计的类与函数？

■

例如：

■

第一步：  
for  
i  
,  
data  
in  
enumerate

(  
t  
r  
a  
i  
n  
-  
l  
o  
a  
d  
e  
r  
)

- 第二步：  
Data Loader 类  
,  
-  
-  
-  
iterator  
-  
-  
-  
函数

- 第三步

:  
\*  
\*  
\*  
\*  
类  
,  
\*  
\*  
\*  
函  
数

- 第  
n  
步  
:  
R  
M  
B  
D  
a  
t  
a  
s  
e  
t  
类  
,  
-  
-  
g  
e  
t  
i  
t  
e  
m  
-  
-  
函  
数
- 打  
卡  
要  
求

: 文字简答

- 训练 R M B 二分类模型，熟悉数据读取机制，并且从 k a g g l e 中下载猫狗二分类训练数据，

自己编写一个 Dog Cat Dataset ,使得 pytorch 可以对猫狗二分类训练集进行读取。

- 数据

下  
載  
:  
<http://www.kaggle.com/ucdavis/riddex>

k  
e  
r  
n  
e  
l  
s  
=  
e  
d  
i  
t  
i  
o  
n  
/  
d  
a  
t  
a



