# 自然語言處理與應用
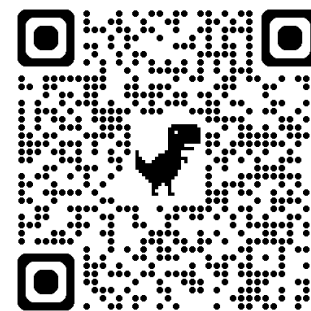## Natural Language Processing and Applications

PyTorch Modeling

Instructor: 林英嘉 (Ying-Jia Lin)
2025/03/17

Course GitHub

Slido # NLP0317

# Steps for building your first PyTorch program

**Step 1 (Data):**
- Prepare the dataset
- Overwrite PyTorch Dataset
- Define DataLoader

**Step 2 (Model):**
- Construct the model
- Define the loss function
- Define the optimizer

**Step 3 (Training):**
Write the training process

**Step 4 (Evaluation):**
Write the evaluation process

NLP

# Step 1: Prepare the dataset

- From [torchvision (image data)](#) or [torchtext (text data)](#)

  - You may skip Step 1-2.

- User-defined dataset

  - Download from the Internet

  - Your own dataset

NLP

# What is a dataset?

| ☞ PassengerId | # Survived | # Pclass | △ Name | △ Sex |
|---|---|---|---|---|
| 1 — 891 | 0 — 1 | 1 — 3 | **891** unique values | male 65%  female 35% |
| 1 | 0 | 3 | Braund, Mr. Owen Harris | male |
| 2 | 1 | 1 | Cumings, Mrs. John Bradley (Florence Briggs Thayer) | female |
| 3 | 1 | 3 | Heikkinen, Miss. Laina | female |
| 4 | 1 | 1 | Futrelle, Mrs. Jacques Heath (Lily May Peel) | female |
| 5 | 0 | 3 | Allen, Mr. William Henry | male |
| 6 | 0 | 3 | Moran, Mr. James | male |
| 7 | 0 | 1 | McCarthy, Mr. Timothy J | male |

dataset

data / instance /example

NLP

https://www.kaggle.com/competitions/titanic/data?select=train.csv

4

# Step 1-2: Overwrite PyTorch Dataset

- 為了符合我們載入資料的<span style="color:red">需求</span>

  - 例如：適合我們資料的前處理過程

- 簡潔且容易維護的資料存取介面：

```
img, label = dataset[0] # `dataset` 是透過 PyTorch Dataset 所建立的
```

index

# Step 1-2: Overwrite PyTorch Dataset

- 我們需要繼承 torch.utils.data.Dataset，並改寫三個項目 (__init__, __getitem__, __len__)：

```python
import torch

class CustomDataset(torch.utils.data.Dataset):
    def __init__(self, parameter_1, parameter_2, ...):
        # Prepare some things
        # that you are going to use in `__getitem__` and `__len__`

    def __getitem__(self, index):
        # do_something
        return data, label

    def __len__(self):
        return len(data_variable)
```

- __init__：初始化 class 中的變數

- __getitem__：讓 PyTorch Dataset 可以透過 index 來取得任一筆資料

- __len__：取得資料集的總數

# Step 1-2: Overwrite PyTorch Dataset

Step 1 (Data)

```python
import torch

class HandWrite(torch.utils.data.Dataset):
    def __init__(self, files: list, word_to_id: dict, transform=None):
        self.files = files # 全部的資料
        self.transform = transform # 影像資料前處理的流程

    def __getitem__(self, index):
        fname = self.files[index]
        image = Image.open(fname)
        if self.transform is not None:
            image = self.transform(image)

        label = fname.split('/')[-1].split('_')[0]
        return image, torch.tensor(word_to_id[label])

    def __len__(self):
        return len(self.files)
```
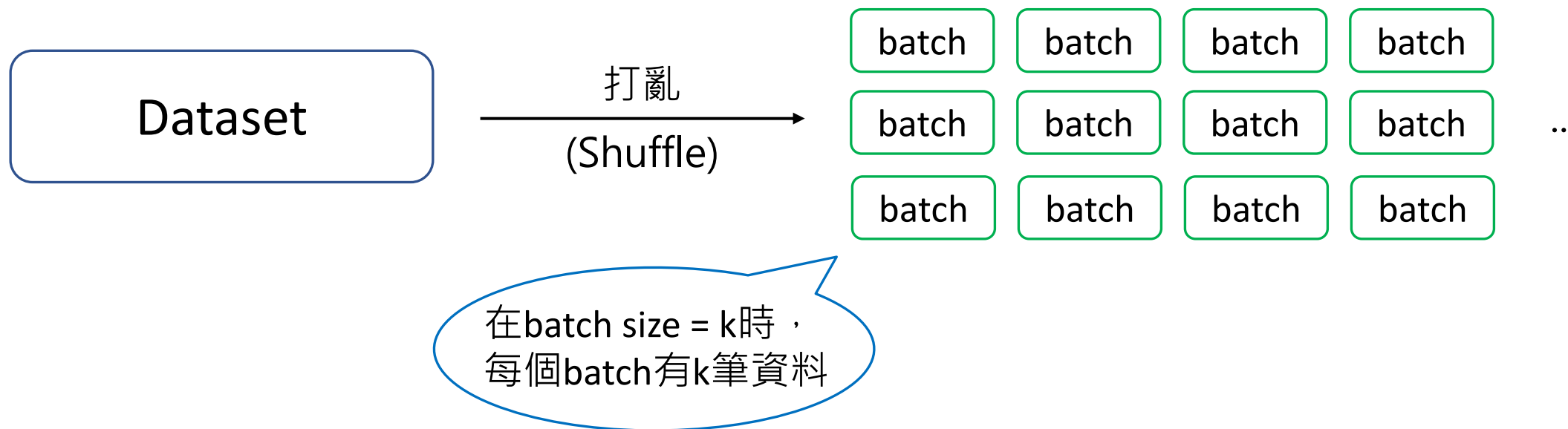
NLP

7

# Step 1-3: Define DataLoader

| | |
|---|---|
| Dataset | 打亂 (Shuffle) → |

batch batch batch batch
batch batch batch batch ...
batch batch batch batch

在batch size = k時，
每個batch有k筆資料

```
# We should split the dataset into train / validation / test sets first.
train_loader = torch.utils.data.DataLoader(trainset, batch_size=TRAIN_BS, shuffle=True)
val_loader = torch.utils.data.DataLoader(valset, batch_size=VAL_BS, shuffle=False)
test_loader = torch.utils.data.DataLoader(testset, batch_size=TEST_BS, shuffle=False)
```

NLP

# Advantages of batching

- Training:
  - mini-batch gradient descent 有機會避免模型陷入局部最小值
- Inference (validation or test):
  - 省記憶體
  - 不需要累積梯度，所以 inference 時期的 batch size (bs) 通常可以比 training 時期的 bs 還大

# Step 2-1: Construct the model

- 我們需要：

    1. 繼承 torch.nn.Module，

    2. 初始化 torch.nn.Module 原本定義的內容

    3. 改寫兩個項目 (__init__, forward)

```python
class MyModel(torch.nn.Module):
    def __init__(self):
        super().__init__() # 初始化 torch.nn.Module 原本定義的內容
        # Define our new variables
        # Define our model layers

    def forward(self, x):
        # Do something (forward pass)
        return output
```

https://pytorch.org/tutorials/beginner/basics/buildmodel_tutorial.html

NLP

# 為什麼需要 super().__init__()？

- 模型需要繼承 torch.nn.Module，並且透過 super().__init__() 初始化原本在 nn.Module 中被定義好的內容，如下圖所示：

```python
206     def __init__(self):
207         """
208         Initializes internal Module state, shared by both nn.Module and ScriptModule.
209         """
210         torch._C._log_api_usage_once("python.nn_module")
211
212         self.training = True
213         self._parameters = OrderedDict()
214         self._buffers = OrderedDict()
215         self._non_persistent_buffers_set = set()
216         self._backward_hooks = OrderedDict()
217         self._forward_hooks = OrderedDict()
218         self._forward_pre_hooks = OrderedDict()
219         self._state_dict_hooks = OrderedDict()
220         self._load_state_dict_pre_hooks = OrderedDict()
221         self._modules = OrderedDict()
```
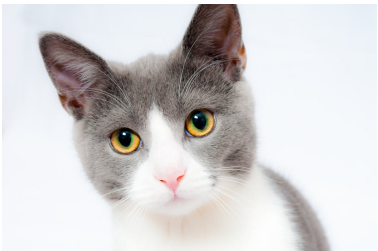
https://github.com/pytorch/pytorch/blob/266657182a4040937f6f27011d7ddf77716db83d/torch/nn/modules/module.py#L206-L221

NLP

# Step 2-2: Define the loss function

| Loss functions | Usage |
| --- | --- |
| torch.nn.CrossEntropyLoss | Classification |
| torch.nn.MSELoss | Regression |
| torch.nn.BCELoss | Binary classification |

```
loss_function = torch.nn.CrossEntropyLoss()
```

NLP

# 模型輸出的後處理

Cross-entropy: $\mathcal{L}_i = -\log P(Y = y_i | X = x_i)$



| Unnormalized log-probabilities / logits | Unnormalized probabilities | Probabilities |
|---|---|---|
| 0.5 | 1.6487 | 0.225 |
| 0.7 | 2.0138 | 0.275 |
| 1.3 | 3.6693 | 0.500 |

Model output              Exponential              Softmax

Cat

Dog

Apple

# Cross-entropy (交叉熵)

Cross-entropy: $\mathcal{L}_i = -\log P(Y = y_i | X = x_i)$

其中 $i$ 代表第 $i$ 筆資料

交叉」(Cross) 代表的是 兩個機率分布之間的關係，特別是用一個分布來衡量 與另一個分布的相似程度

- 量測模型輸出的負對數機率，代表模型預測該類別的信心程度

  - 模型預測該類別的信心程度越大時， $\mathcal{L}_i$ 就會越小

  - 模型預測該類別的信心程度越小時， $\mathcal{L}_i$ 就會越大

| $P(Y = y_i | X = x_i)$ | $\mathcal{L}_i = -\log P(Y = y_i | X = x_i)$ |
|---|---|
| 0.9 | $-\log 0.9 \approx 0.105$ |
| 0.1 | $-\log 0.1 \approx 2.302$ |

NLP

# Softmax (Non-linear Transformation)

- 採用 exponential -> 大的數值更大，小的數值更小
  - 有助於梯度下降

$$y = \frac{x}{\sum_j x_j} = [0.5, 0.25, 0.25]$$

$$\text{Softmax}(x_i) = \frac{e_i^x}{\sum_j e^{x_j}} = [0.665, 0.244, 0.090]$$

# Step 2-3: Define the optimizer

| Loss functions | Meaning |
|----------------|---------|
| torch.optim.SGD | **Stochastic gradient descent** (with momentum) |
| torch.optim.RMSprop | **RMSProp** (Root Mean Square Propagation) |
| torch.optim.Adam | **Adam** (Adaptive Moment Estimation) |
| torch.optim.AdamW | **AdamW** (Adam with decoupled weight decay) |

```
learning_rate = 1e-3 # 代表 0.001
optimizer = optim.Adam(model.parameters(), lr=learning_rate)
```

NLP

# Step 3: Write the training process

1. Clear gradients

   `optimizer.zero_grad()`

2. Input data to the model

   `output = model(**batch)`

3. Computer loss

   `loss = loss_function(gold, output)`

4. Computer gradients

   `loss.backward()`

5. Update model parameters

   `optimizer.step()`

6. (Repeat 1. to 5. until the end of training)

# ** in Python

- **dict 可以展開字典，轉換成關鍵字參數傳遞給函式

```
def greet(age, name):
    print(f"My name is {name} and I am {age} years old.")

person_info = {"name": "Alex", "age": 25}
greet(**person_info)  # 等同於 greet(name="Alex", age=25)
```

# Step 3: Write the training process

1. Clear gradients

2. Input data to the model

3. Computer loss

4. Computer gradients

5. Update model parameters

6. (Repeat 1. to 5. until the end of training)

```
optimizer.zero_grad()

output = model(**batch)

loss = loss_function(gold, output)

loss.backward()

optimizer.step()
```

NLP

19

# Step 3: Write the training process

```
for batch in train_loader:
    output = model(**batch)
    ...
```

```
output = model(**batch)
```

```
# Get x, y from your dataloader
for batch_x, batch_y in train_loader:
    output = model(batch_x)
    loss = criterion(output, target)
    ...
```

NLP

# Step 4: Write the evaluation process

```python
from sklearn.metrics import accuracy_score

with torch.no_grad():
    for batch in val_loader: # or test_loader
        output = model(**batch)
        pred = outputs.argmax(dim=1)
        ...
        predictions.append(pred)


accuracy_score(test_labels, predictions)
```

NLP

# Thank you!

Instructor: 林英嘉

✉ yjlin@cgu.edu.tw

TA: 吳宣毅

✉ m1161007@cgu.edu.tw