

01_preprocessing

October 26, 2025

0.1 0. Introduction

In recent years, the development of science and technology has had many outstanding achievements in general and in the field of Artificial Intelligence (AI) in particular. In the field of imaging, people have been able to take advantage of collecting images from digital electronic devices such as cameras, camcorders, etc. and use them as input for AI models and take advantage of the computing and learning capabilities of AI models to create predictive models serving humans such as traffic control, behavior and emotion recognition, etc. In particular, in the field of transportation, AI has many practical applications serving human interests such as predicting traffic jams, identifying vehicles participating in traffic, recognizing license plates, etc.

```

```

In this project, we will research and build a small application on license plate recognition based on the **License Plate** dataset from the **Roboflow** source and perform the **Fine Tuning** technique and use a pre-trained model such as **Resnet**, **MobileNet**,... to make the backbone for the **Faster RCNN** model.

```

```

```
[1]: import sys
import os

# Add root directory of rhe project into sys.path
sys.path.append(os.path.abspath(os.path.join("../")))
```

0.2 1. Libraries

In this project, we need some external libraries such as OpenCV, Pytorch,... (version details in the **requirements** file). Because this project needs to process image features, the OpenCV library is a specialized tool for image problems such as loading images, resizing,... The **Pytorch** library is a popular tool specializing in building deep learning models, the library provides us with classes,

metrics, optimizers,... to help us build the desired models, more specifically, this library provides us with classes such as **Dataset**, **DataLoader** to help build dataset models to avoid memory overflow problems when saving all images in an array. In addition, supporting libraries such as Numpy, Matplotlib,...

```
[2]: from src.data_preprocessing import download_dataset, preprocess_data
      from src.config import IMAGE_SIZE, RAW_DATA_DIR, PROCESSED_DATA_DIR
      from src.utils import print_tree
```

0.3 2. Dataset

In this project, we will use the **License Plate Computer Vision Project** dataset (Data information in README file) provided by **Universe**, this dataset will include more than **20000** images for the training set, **2000** images in the validation set and **1000** images in the test set. In particular, each image in the dataset has a size of **640x640x3** (RGB color image) and has been preprocessed such as rotating the image, zooming, replacing white points,...

```

```

Each image is accompanied by an XML file with the same file name as the image file, this XML file contains information about the image such as image file name, image file path, dimensions (**width**, **height**) and information about the bounding box, which is also the label we will predict. These bounding boxes will be the points **xmin**, **ymin**, **xmax** and **ymax** that form a square around the license plate.

```
<annotation>
  <folder/>
  <filename>00a09b822d470896.jpg.rf.88ddf90637cd97ac57c03910636a0294.jpg</filename>
  <path>00a09b822d470896.jpg.rf.88ddf90637cd97ac57c03910636a0294.jpg</path>
  <source>
    <database>roboflow.com</database>
  </source>
  <size>
    <!-- Image shape -->
    <width>640</width>
    <height>640</height>
    <depth>3</depth>
  </size>
  <segmented>0</segmented>
  <object>
    <name>License Plate Recognition - v1 raw-images</name>
    <pose>Unspecified</pose>
    <truncated>0</truncated>
    <difficult>0</difficult>
    <occluded>0</occluded>
```

```

<!-- Bounding box -->
<bndbox>
    <xmin>2</xmin>
    <xmax>43</xmax>
    <ymin>408</ymin>
    <ymax>428</ymax>
</bndbox>
</object>
...
<object>
    <name>License Plate Recognition - v1 raw-images</name>
    <pose>Unspecified</pose>
...
</object>
</annotation>

```

[3]: # Download dataset from Github
`download_dataset(dest=RAW_DATA_DIR)`

```

Dataset does not exist, please waiting to download data from the cloud...
Start to download...
0/20 of data are downloaded
5/20 of data are downloaded
10/20 of data are downloaded
15/20 of data are downloaded

license-plate-project.zip: 100% | 965M/965M [01:55<00:00, 8.77MB/s]

Downloaded, please waiting to extract file...
Extracting...
Extracted!!
Removing zip file...
Removed!!
Done!!

```

[4]: # Display the directory tree of raw data folder
`print_tree(RAW_DATA_DIR)`

```

[DIR] raw
[DIR] test
    [FILE] xemay1817.jpg.rf.119f2c447b36c4a29c10f0ef8e90e019.xml
    [FILE] CarLongPlateGen2231.jpg.rf.09a6ae0129bd4a3ccac80d66ba4e4b95.xml
    [FILE] xemay246.jpg.rf.530bba55def20c3ce976703557b39b7b.xml
    [FILE] xemay66.jpg.rf.a9939bfd9034efff251c176215d259cc.jpg
    [FILE] CarLongPlateGen2960.jpg.rf.bda1b9f4642a8866ed162c2edfbe3b94.xml
    [...]
[DIR] valid
    [FILE] CarLongPlateGen1530.jpg.rf.3289b63b8aff3d0c30701736cc0d7712.jpg
    [FILE] CarLongPlateGen1708.jpg.rf.2562fda9da07faf37e8f1cfad4393f9f.xml
    [FILE] xemay397.jpg.rf.5b5289344b4ebe46f7c772ac2fb435e3.jpg

```

```

[FILE] CarLongPlateGen1809_jpg.rf.c451d9563fb97938eaf545a70ff4c457.jpg
[FILE] CarLongPlateGen86_jpg.rf.6aa8cb57bb3b078e4f1a23b0efe49706.jpg
[...]
[DIR] train
[FILE] CarLongPlateGen2707_jpg.rf.92e2bb3a0499e548f3bc2e54cc4ce299.xml
[FILE] 003df8cf2effae50_jpg.rf.85b92c041e14d9bcf4ed1fd70de9661f.jpg
[FILE] CarLongPlateGen2372_jpg.rf.0d12d2a9b4f0f957096227bf11d30626.jpg
[FILE] Cars168_png_jpg.rf.33aee6637f3a5a92f0bb1bd1e3005c3e.xml
[FILE] CarLongPlate584_jpg.rf.f0130e378d8afc47359cbaca7e8a8969.xml
[...]

```

0.4 3. Data preprocessing

```
[5]: # Data preprocessing
preprocess_data(RAW_DATA_DIR, PROCESSED_DATA_DIR, IMAGE_SIZE)

Preprocessing data in test folder...: 100% | 978/978 [00:02<00:00,
420.65it/s]
Preprocessing data in valid folder...: 100% | 1973/1973 [00:05<00:00,
368.22it/s]
Preprocessing data in train folder...: 100% | 20580/20580
[01:08<00:00, 298.58it/s]
```

```
[6]: print_tree(PROCESSED_DATA_DIR)
```

```

[DIR] processed
[DIR] test
[FILE] xemay1817.jpg.rf.119f2c447b36c4a29c10f0ef8e90e019.xml
[FILE] CarLongPlateGen2231.jpg.rf.09a6ae0129bd4a3ccac80d66ba4e4b95.xml
[FILE] xemay246.jpg.rf.530bba55def20c3ce976703557b39b7b.xml
[FILE] xemay66.jpg.rf.a9939bfd9034efff251c176215d259cc.jpg
[FILE] CarLongPlateGen2960.jpg.rf.bda1b9f4642a8866ed162c2edfbe3b94.xml
[...]
[DIR] valid
[FILE] CarLongPlateGen1530.jpg.rf.3289b63b8aff3d0c30701736cc0d7712.jpg
[FILE] CarLongPlateGen1708.jpg.rf.2562fda9da07faf37e8f1cfad4393f9f.xml
[FILE] xemay397.jpg.rf.5b5289344b4ebe46f7c772ac2fb435e3.jpg
[FILE] CarLongPlateGen1809.jpg.rf.c451d9563fb97938eaf545a70ff4c457.jpg
[FILE] CarLongPlateGen86.jpg.rf.6aa8cb57bb3b078e4f1a23b0efe49706.jpg
[...]
[DIR] train
[FILE] CarLongPlateGen2707.jpg.rf.92e2bb3a0499e548f3bc2e54cc4ce299.xml
[FILE] 003df8cf2effae50.jpg.rf.85b92c041e14d9bcf4ed1fd70de9661f.jpg
[FILE] CarLongPlateGen2372.jpg.rf.0d12d2a9b4f0f957096227bf11d30626.jpg
[FILE] Cars168_png.jpg.rf.33aee6637f3a5a92f0bb1bd1e3005c3e.xml
[FILE] CarLongPlate584.jpg.rf.f0130e378d8afc47359cbaca7e8a8969.xml
[...]

```