

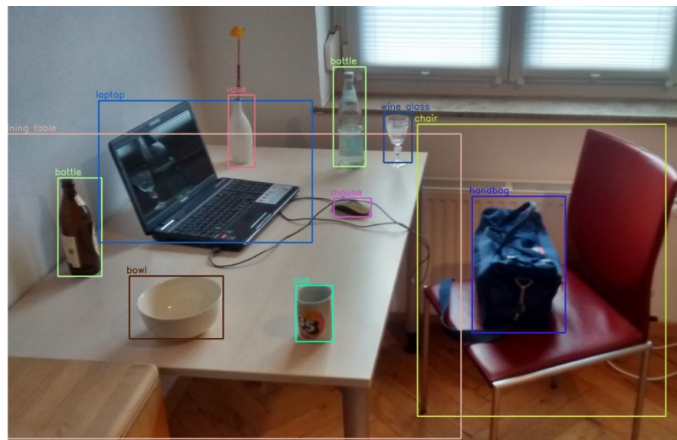
# Image Project - Object Detection: Traffic Sign

Ngày 15 tháng 3 năm 2024

<b>Ngày tóm tắt:</b>	13/03/2024
<b>Tác giả:</b>	AIO
<b>Nguồn:</b>	Module 4 AIO 2023
<b>Nguồn dữ liệu (nếu có):</b>	
<b>Từ khóa:</b>	Object Detection, Classification, Localization
<b>Người tóm tắt:</b>	Nguyễn Mậu Ánh Ngân

## 1 Giới thiệu

- Trong bài hôm nay ta sẽ không sử dụng thư viện có sẵn mà xây dựng hoàn toàn theo từng bước.
- Object Detection (OD) là gì: ý tưởng chính là khi ta đưa vào một tấm ảnh (Input) thì ta sẽ xác định được những vật thể (Object) có trên ảnh kèm theo:
  - Hình chữ nhật (Bounding box) xác định vị trí giới hạn của Object.
  - Tên (Classname/ Label) của Object.



Hình 1: Ví dụ một Output của Object Detection.

- Phân biệt OD và Image Classification (IC)
  - Giống nhau: Input đều là ảnh.
  - Khác nhau:

- \* OD: Output xác định được cả vị trí (qua Bounding Boxes) và classname/label của Object trong ảnh.
- \* IC: Output chỉ xác định classname/Label chung của ảnh.

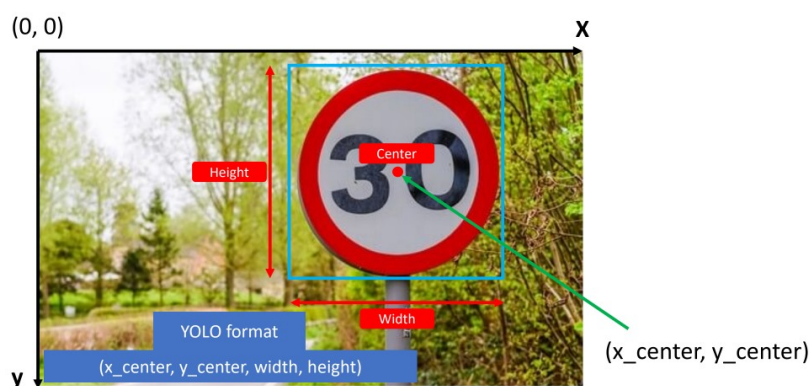
- Cách xây dựng một bài OD: bao gồm 2 model chính
  - Localization: xác định vị trí Object -> Output: bounding boxes.
  - Classification: Xác định classname cho các bounding box.

Ta gọi là mô hình Two-stage Object Detection (hình ảnh được xử lý thường qua Localization trước)

## 1.1 Object Localization:

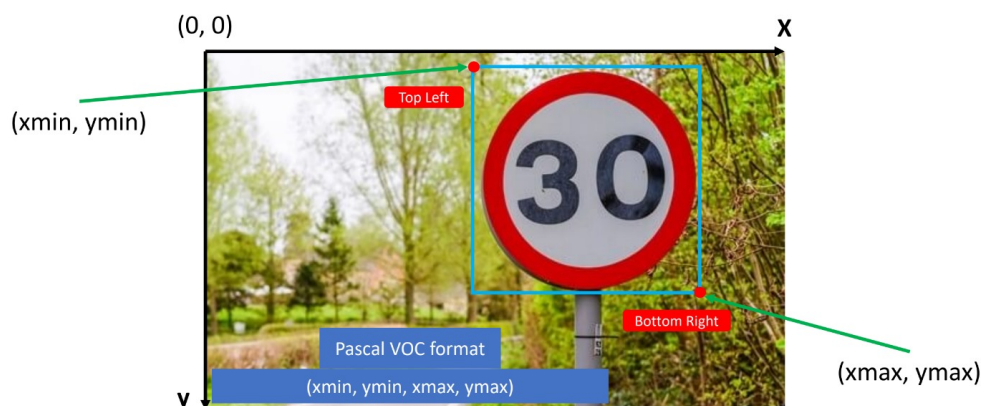
- Bounding box: là hình chữ nhật xác định vị trí của Object.
- Trong code: Output sẽ là một Tuple gồm tọa độ của bounding box. Thường có 2 dạng tuple:

### 1. YOLO format:



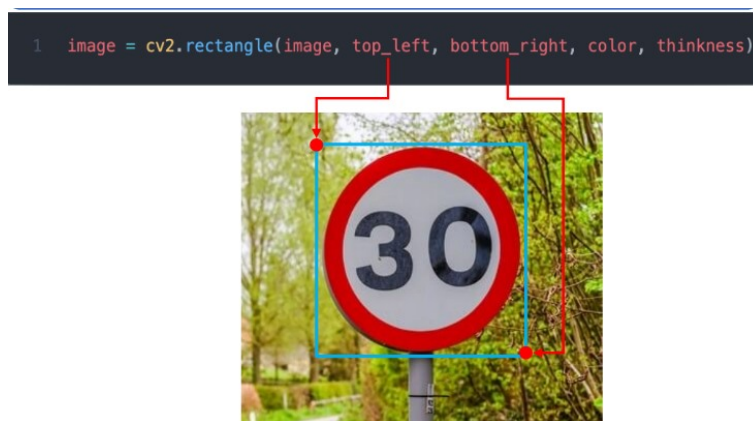
Hình 2: YOLO format.

### 2. Pascal VOC format:



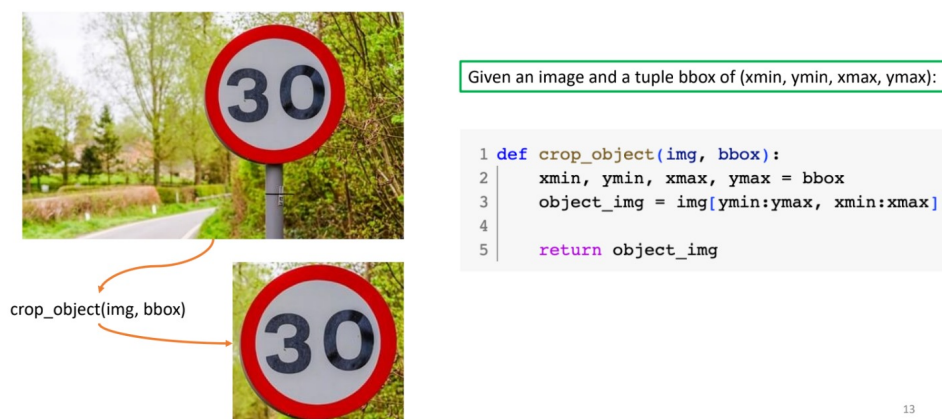
Hình 3: Pascal VOC format.

- Code vẽ Bounding Box:



Hình 4: Code vẽ Bounding box.

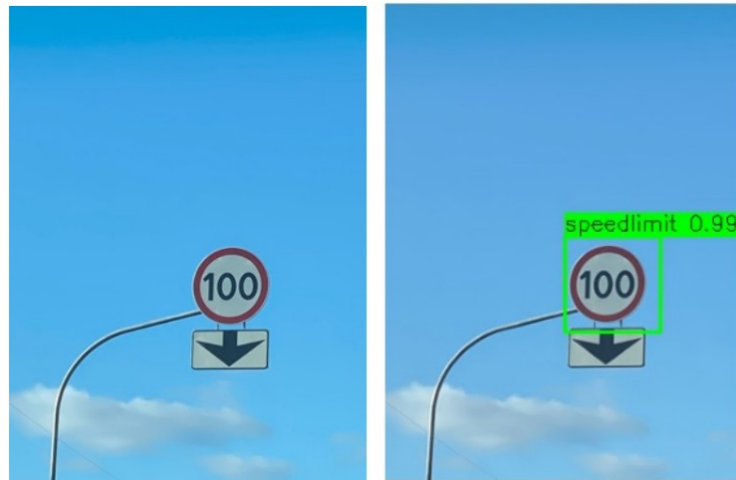
- Chuyển đổi giữa 2 format (xem thêm ở đoạn record: [00:23:07-00:24:18])
- Cắt Object ra từ hình (Format Pascal)



Hình 5: Code cắt Object từ ảnh.

## 1.2 Object Classification:

- Mục đích là phân loại, tạo classname cho bounding boxes
- Output: list tuple (tọa độ bounding box, confidence score) vì thường một hình sẽ có nhiều Object nên sẽ cho ra list.
- Confidence score: mức độ tự tin của mô hình về việc đoán đúng classname. Thường nằm trong khoảng từ 0 tới 1.

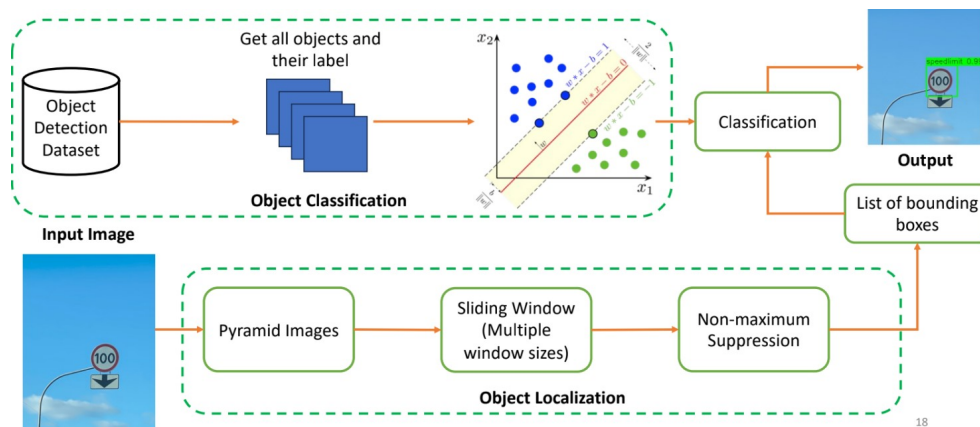


Hình 6: Hình sau khi đã qua cả Localization và Classification.



Hình 7: Confidence score.

### 1.3 Sơ đồ chung xây dựng bài toán:



Hình 8: Các bước xây dựng bài toán Object Detection.

### 1.4 Input các thư viện cần thiết:

```

1 import time
2 import os
3 import cv2
4 import numpy as np
5 import pandas as pd
6 import matplotlib.pyplot as plt
7 import xml.etree.ElementTree as ET
8
9 from skimage.transform import resize
10 from skimage import feature
11 from sklearn.svm import SVC
12 from sklearn.preprocessing import LabelEncoder
13 from sklearn.preprocessing import StandardScaler
14 from sklearn.model_selection import train_test_split
15 from sklearn.metrics import accuracy_score

```



Hình 9: Các thư viện dùng trong code.

## 2 Classification:

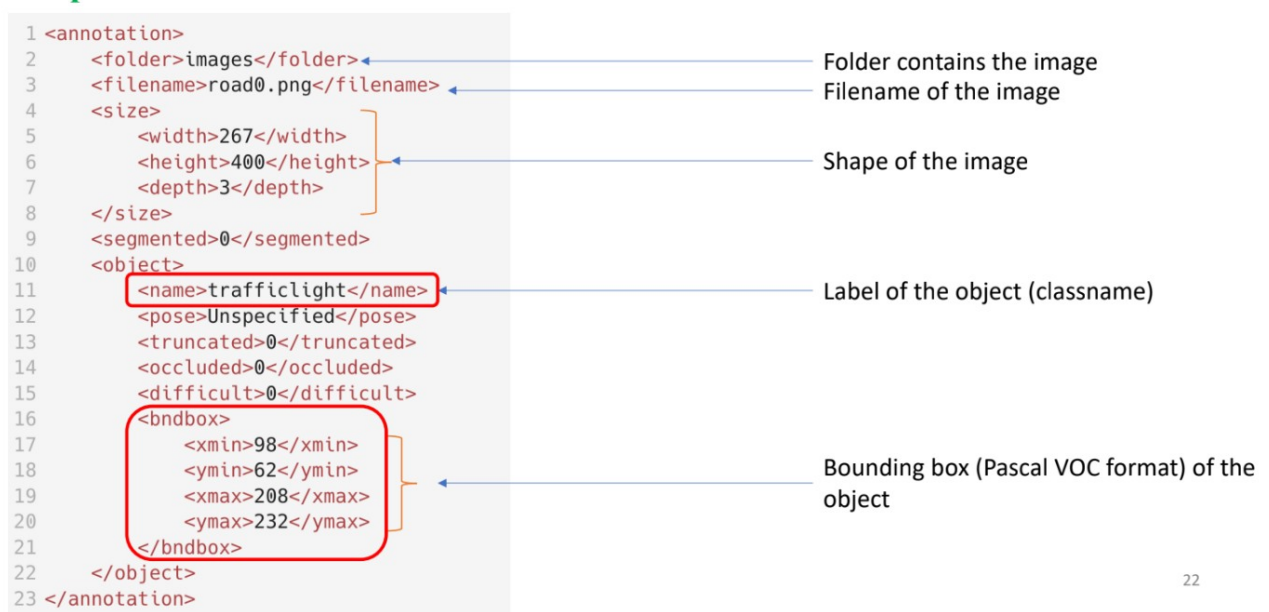
- Tải Dataset (tập dữ liệu dùng để train) từ link sau: [Traffic Sign Detection](#).
- Model ta sẽ dùng là SVM.

### 2.1 Step 1: Đọc Dataset

- Trong Dataset gồm 2 folder:
  - Images: chứa các file ảnh.
  - Annotations: chứa các file label của ảnh.



Hình 10: Các folder trong dataset: ảnh và label tương ứng.



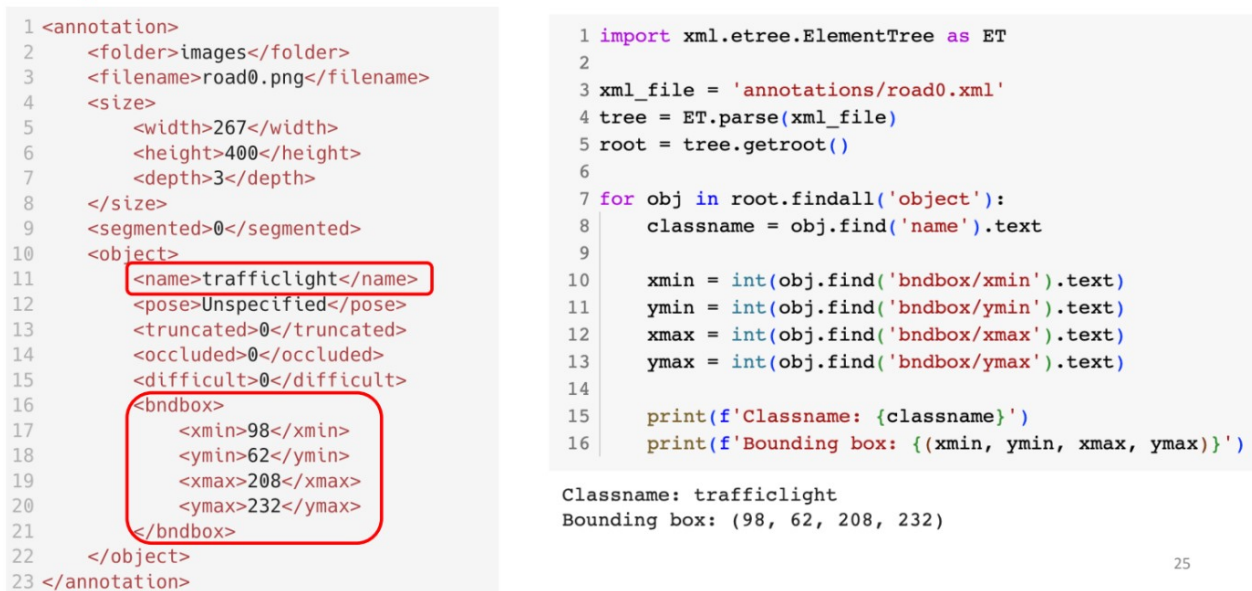
Hình 11: Thông tin trong file label của ảnh.

- Dùng Xml để lấy thông tin từ file label





Hình 12: Cách sử dụng xml cơ bản để lấy thông tin.



Hình 13: Lấy thông tin Classname và tọa độ.

## ❖ Step 1: Read dataset



```
1 print('Number of objects: ', len(img_lst))
2 print('Class names: ', list(set(label_lst)))
```

```
Number of objects: 1074
Class names: ['stop', 'crosswalk', 'speedlimit']
```

```
1 annotations_dir = 'annotations'
2 img_dir = 'images'
3
4 img_lst = []
5 label_lst = []
6 for xml_file in os.listdir(annotations_dir):
7     xml_filepath = os.path.join(annotations_dir, xml_file)
8     tree = ET.parse(xml_filepath)
9     root = tree.getroot()
10
11     folder = root.find('folder').text
12     img_filename = root.find('filename').text
13     img_filepath = os.path.join(img_dir, img_filename)
14     img = cv2.imread(img_filepath)
15
16     for obj in root.findall('object'):
17         classname = obj.find('name').text
18         if classname == 'trafficlight':
19             continue
20
21         xmin = int(obj.find('bndbox/xmin').text)
22         ymin = int(obj.find('bndbox/ymin').text)
23         xmax = int(obj.find('bndbox/xmax').text)
24         ymax = int(obj.find('bndbox/ymax').text)
25
26         object_img = img[ymin:ymax, xmin:xmax]
27         img_lst.append(object_img)
28         label_lst.append(classname)
```

We ignore 'trafficlight' class since we do Traffic Sign Classification

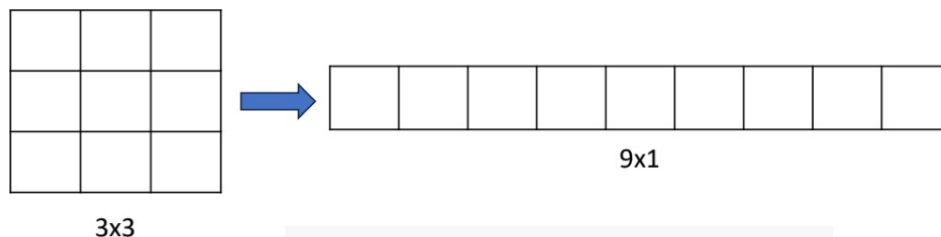
Hình 14: Duyệt qua tất cả các file label, bỏ qua label Trafficlight.

## 2.2 Step 2: Tiền xử lý ảnh

- Với shape như bên (cao,rộng,kênh màu) thì ta không thể đẩy vào mô hình SVM được -> Ta phải flatten về lại len(shape) == 1.

```
1 print('Image shape:')
2 print(img_lst[1].shape)
```

```
Image shape:
(53, 55, 3)
```



We can flatten image to have the len(shape) == 1.

```
1 print('Image shape:')
2 print(img_lst[1].shape)
3 flattened_img = img_lst[1].flatten()
4 print('Image shape after flattened:')
5 print(flattened_img.shape)
```

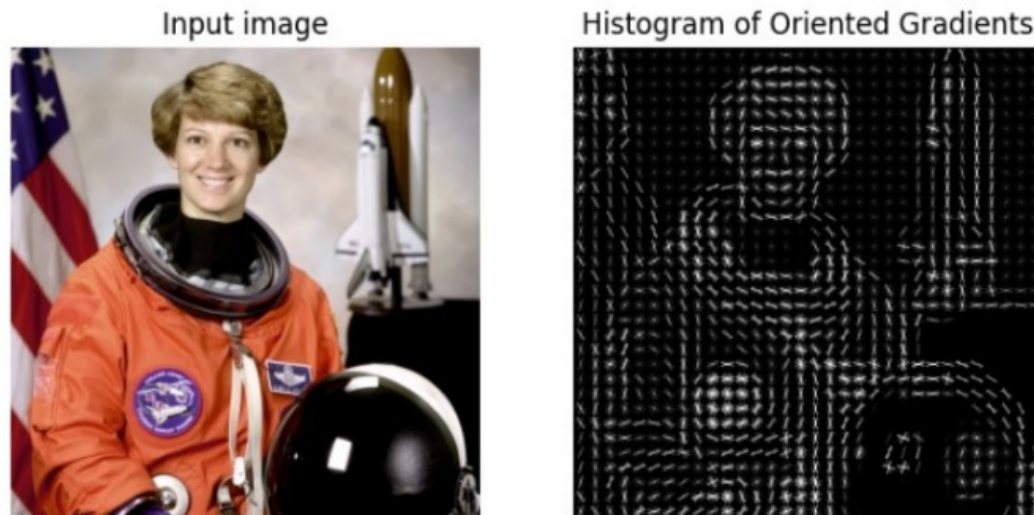
```
Image shape:
(53, 55, 3)
Image shape after flattened:
(8745,)
```

27

Hình 15: Flatten.

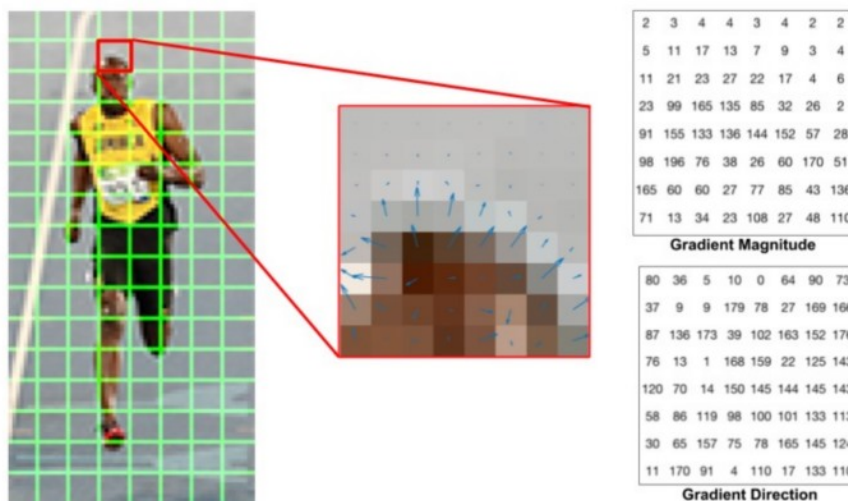


- Sau khi flatten, ta nhận được vector đại diện 1D có tới 8745 phần tử -> quá phức tạp -> nên xử lý ảnh gốc trước khi flatten. Ta dùng phương pháp HoG ( Histogram of Oriented Gradients).



Hình 16: Phương pháp HoG.

- Cách HoG hoạt động:
  - Ảnh đầu vào sẽ được chia thành các cell, mỗi cell có kích thước 8x8.
  - Với mỗi cell ta sẽ tìm 2 ma trận:
    - \* Gradient Magnitude: Ma trận độ lớn Gradient.
    - \* Gradient Direction: Ma trận phương Gradient.



Hình 17: Cell và 2 ma trận cần tìm.

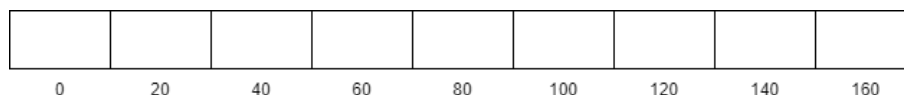
- Cách tính 2 ma trận:

The diagram illustrates the calculation of gradient matrices  $G_x$  and  $G_y$  using Sobel kernels. It shows the following formulas and components:

- Kernel Sobel:** Indicated by arrows pointing to the Sobel kernels used in the calculations.
- Tích chập (Convolution):** Indicated by an arrow pointing to the convolution operation.
- Ảnh gốc (Original Image):** Indicated by an arrow pointing to the input image.
- Gradient Magnitude:** Indicated by an arrow pointing to the formula  $G = \sqrt{G_x^2 + G_y^2}$ .
- Gradient Direction:** Indicated by an arrow pointing to the formula  $\theta = \arctan\left(\frac{G_y}{G_x}\right)$ .

Hình 18: Công thức tính 2 ma trận.

- Tiếp theo, ta tạo ra vector 1D gồm có 9 phần tử đại diện cho miền giá trị của ma trận phương Gradient nằm trong khoảng  $[0,180]$ . Chia vector thành 9 khoảng như hình, mỗi khoảng được gọi là 1 bin. Vector này được gọi là vector Histogram of Gradient.



Hình 19: Vector đại diện gồm 9 bin.

- Tham chiếu 2 ma trận đã tìm được ở trên vào vector đại diện ta có 2 trường hợp gán giá trị:
  - \* Gradient Direction có độ lớn trùng với bin (0,20,40...) => gán luôn giá trị tương ứng ở Gradient Magnitude vào bin đó. (Ví dụ vị trí màu xanh trong hình 21)
  - \* Gradient Direction có độ lớn thuộc khoảng bin => dùng công thức nội suy sau để tìm giá trị gán:

If not, assign magnitude to two adjacent bins  
 $x \in [x_0, x_1]$  using the following formulas:

$$x_{l-1} = \frac{(x_1 - x)}{x_1 - x_0} * \text{magnitude}$$

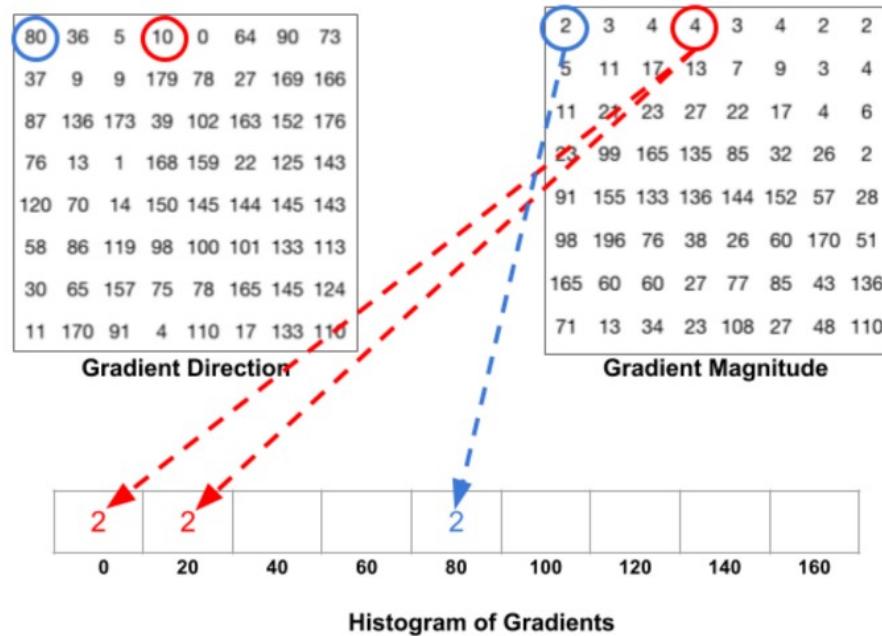
$$x_l = \frac{(x - x_0)}{x_1 - x_0} * \text{magnitude}$$

Hình 20: Công thức nội suy.

Ví dụ: Trong hình 21, vị trí màu đỏ ta chọn có giá trị Gradient Direction = 10 thuộc bin 0 tới 20 và Gradient Magnitude tương ứng = 4 ta áp dụng công thức:

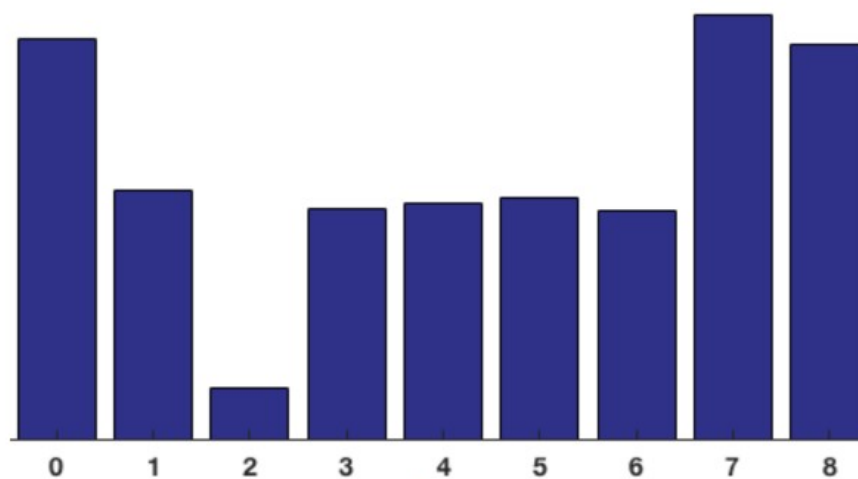
$$x_0 = \frac{20 - 10}{20 - 0} * 4 = 2, x_{10} = \frac{10 - 0}{20 - 0} * 4 = 2$$

Ta được hai giá trị gán vào bin 0 là 2 và bin 10 là 2.



Hình 21: Các trường hợp ví dụ gán giá trị cho vector HoG.

- Lặp lại quá trình tính toán và gán cho toàn bộ 2 ma trận.
- Tính tổng độ lớn Gradient của từng bin biểu diễn thành đồ thị ta có đồ thị Histogram gradient.

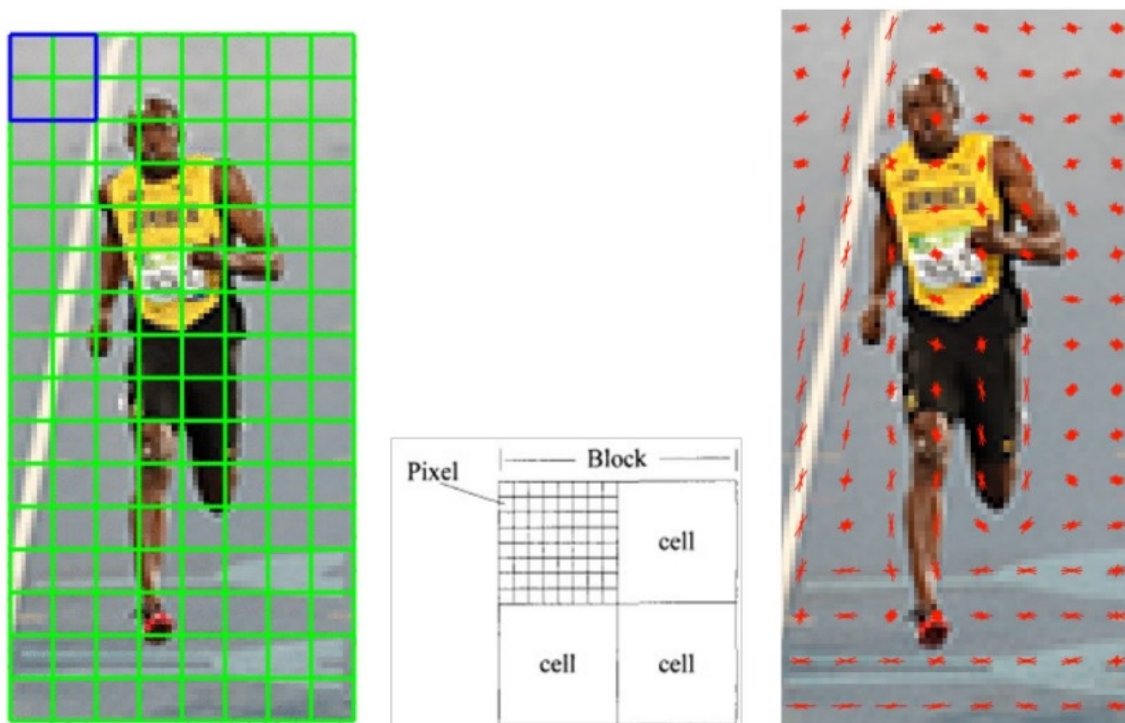


Hình 22: Đồ thị Histogram Gradient.

- Chuẩn hóa cụm cell theo block gồm 2x2 cell. Công thức chuẩn hóa:

$$\text{normalize}(h) = \frac{h}{|h|^2}$$

- Sau khi chuẩn hóa xong cho tất cả các block ta sẽ được hình biểu diễn đặc trưng của HoG.



Hình 23: Hình ảnh HoG thu được sau khi chuẩn hóa tất cả các block.

- Code mẫu: (xem thêm giải thích từ record [01:09:16-01:20:30])

```

1 def preprocess_img(img):
2     if len(img.shape) > 2:
3         img = cv2.cvtColor(img, cv2.COLOR_BGR2GRAY)
4     img = img.astype(np.float32)
5
6     resized_img = resize(
7         img,
8         output_shape=(32, 32),
9         anti_aliasing=True
10    )
11
12    hog_feature = feature.hog(
13        resized_img,
14        orientations=9,
15        pixels_per_cell=(8, 8),
16        cells_per_block=(2, 2),
17        transform_sqrt=True,
18        block_norm="L2",
19        feature_vector=True
20    )
21
22    return hog_feature

```

We resize all images to the same size to ensure HOG vectors have the same shape.

```

1 img_features_lst = []
2 for img in img_lst:
3     hog_feature = preprocess_img(img)
4     img_features_lst.append(hog_feature)
5
6 img_features = np.array(img_features_lst)
7 print('X shape:')
8 print(img_features.shape)

```

X shape:  
(1074, 324)

Hình 24: Code mẫu step 2.

### 2.3 Step 3: Đánh số cho label

Label Encoder	
'crosswalk'	0
'speedlimit'	1
'stop'	2

```

1 label_encoder = LabelEncoder()
2 encoded_labels = label_encoder.fit_transform(label_lst)

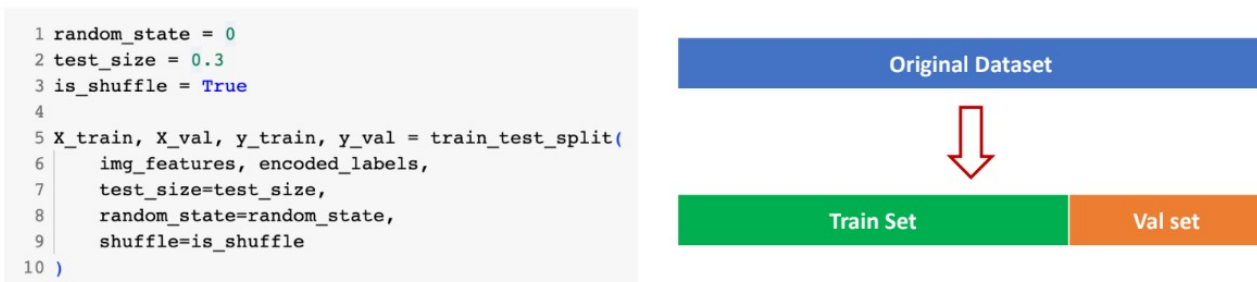
1 label_encoder.classes_
array(['crosswalk', 'speedlimit', 'stop'], dtype='<U10')

1 encoded_labels
array([1, 1, 2, ..., 0, 1, 0])

```

Hình 25: Endcode label.

### 2.4 Step 4: Tạo tập train và tập val



Hình 26: Tạo tập train (70) và tập val (30)

### 2.5 Step 5: Chuẩn hóa data

```

1 scaler = StandardScaler()
2 X_train = scaler.fit_transform(X_train)
3 X_val = scaler.transform(X_val)

```

$$z = \frac{x_i - \mu}{\sigma}$$

Hình 27: Chuẩn hóa cho tập train và val.



## 2.6 Step 6,7: Đưa vào hàm SVM và cho ra kết quả (xem giải thích code tại đoạn record [01:23:46-01:31:15])

```

1 clf = SVC(
2     kernel='rbf',
3     random_state=random_state,
4     probability=True,
5     C=0.5
6 )
7 clf.fit(X_train, y_train)

```

```

1 y_pred = clf.predict(X_val)
2 score = accuracy_score(y_pred, y_val)
3
4 print('Evaluation results on val set')
5 print('Accuracy: ', score)

```

Evaluation results on val set  
Accuracy: 0.9659442724458205

SVC

SVC(C=0.5, probability=True, random\_state=0)



```

1 input_img = img_lst[1]
2 normalized_img = preprocess_img(input_img)
3 y_pred = clf.predict([normalized_img])[0]
4 print(f'Normal prediction: {y_pred}')
5 y_pred_prob = clf.predict_proba([normalized_img])
6 prediction = np.argmax(y_pred_prob)
7 y_pred_prob = [f'{p:.10f}' for p in y_pred_prob[0]]
8 print(f'Probability of each class: {y_pred_prob}')
9 print(f'Class with highest probability: {prediction}')

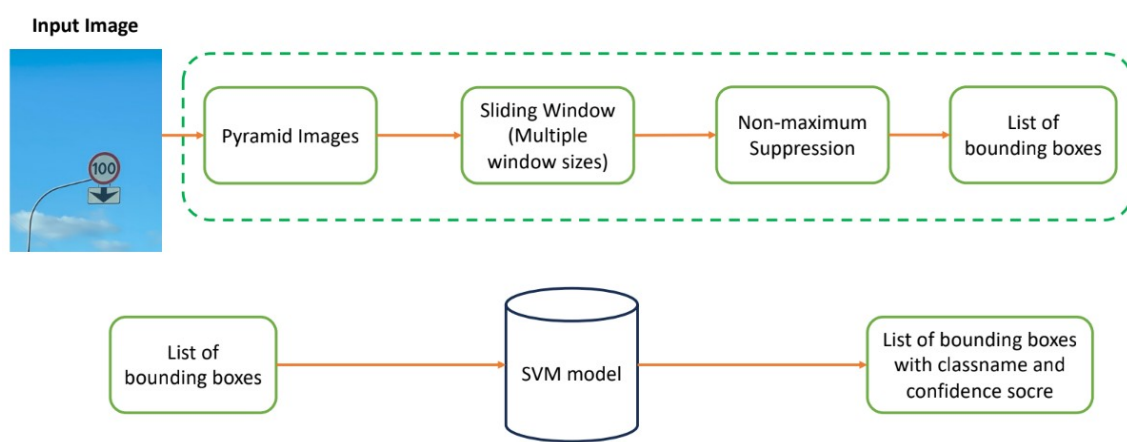
```

Normal prediction: 1  
Probability of each class: ['0.00000000040', '0.9999306581', '0.0000693379']  
Class with highest probability: 1

Hình 28: Code đưa vào mô hình SVM và kết quả.

## 3 Localization:

- Ý tưởng chung:



Hình 29: Ý tưởng xây dựng mô hình Localization và kết quả sau khi xây dựng thành công.

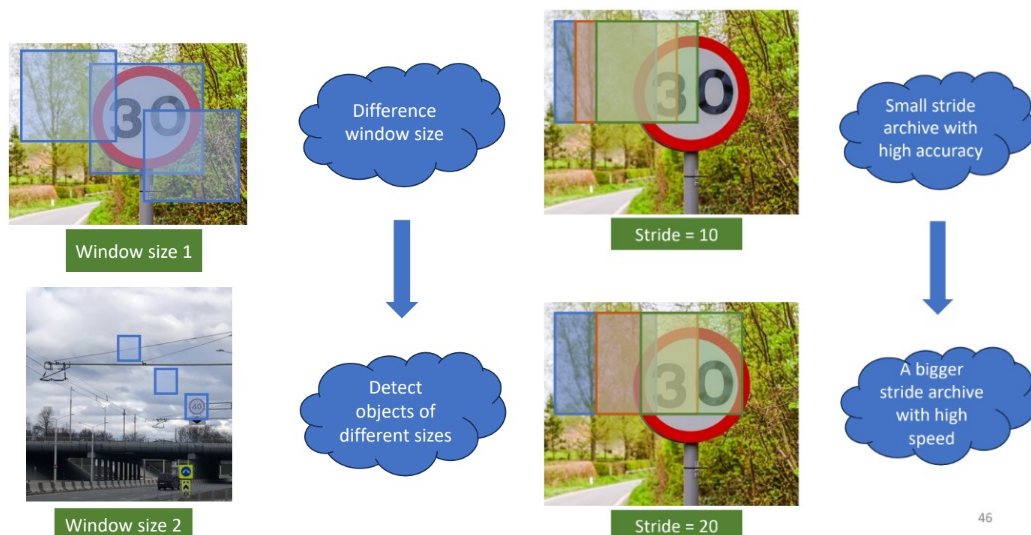
### 3.1 Sliding window:

- Là phần cốt lõi trong xây dựng mô hình Localization ta đang xét.
- Thường được sử dụng trong thị giác máy tính và âm thanh.
- Ý tưởng chung: Tạo ra một hình chữ nhật có kích thước cố định -> áp vào ảnh -> trượt trên ảnh từ trái sang phải từ trên xuống dưới -> nơi nào thấy được object thì làm bounding box luôn. (Giải thuật Greedy - tham lam)



Hình 30: Ý tưởng của Sliding window.

- Có 2 vấn đề cần quan tâm:
  - Kích thước object có thể khác nhau -> không thể dùng hình chữ nhật có kích thước cố định được.
  - Khoảng cách trượt theo chiều ngang (Stride) gây ra sự tỉ lệ nghịch cho thời gian và độ chính xác:
    - \* Stride càng nhỏ -> độ chính xác lớn -> tăng thời gian.
    - \* Stride càng lớn -> độ chính xác bé -> giảm thời gian.



Hình 31: 2 vấn đề chính của Sliding window.

- Code tạo window

```

1 def sliding_window(img, window_sizes, stride, scale_factor):
2     img_height, img_width = img.shape[:2]
3     windows = []
4     for window_size in window_sizes:
5         window_width, window_height = window_size
6         for ymin in range(0, img_height - window_height + 1, stride):
7             for xmin in range(0, img_width - window_width + 1, stride):
8                 xmax = xmin + window_width
9                 ymax = ymin + window_height
10
11                 windows.append([xmin, ymin, xmax, ymax])
12
13     return windows

```

Hình 32: Code mẫu phương pháp Sliding window.

### 3.2 Pyramid Images:

- Vấn đề đặt ra: Khi Sliding với window có kích thước cố định nhưng những object ở xa thì sẽ nhỏ hơn window đó rất nhiều -> phương pháp Pyramid Image trước khi đưa vào Sliding.



Hình 33: Khắc phục vấn đề object nhỏ hơn so với Window bằng phương pháp Pyramid Images.

- Ý tưởng chính: Ta scale thu nhỏ dần tấm hình gốc lại để object bé được phóng lớn hơn -> window sẽ vừa và dễ dàng tìm thấy object hơn.



Hình 34: Scale ảnh theo các tỉ lệ khác nhau.

- Code:

```
1 def pyramid(img, scale=0.8, min_size=(30, 30)):  
2     acc_scale = 1.0  
3     pyramid_imgs = [(img, acc_scale)]  
4  
5     i = 0  
6     while True:  
7         acc_scale = acc_scale * scale  
8         h = int(img.shape[0] * acc_scale)  
9         w = int(img.shape[1] * acc_scale)  
10        if h < min_size[1] or w < min_size[0]:  
11            break  
12        img = cv2.resize(img, (w, h))  
13        pyramid_imgs.append((img, acc_scale * (scale ** i)))  
14        i += 1  
15  
16    return pyramid_imgs
```

Hình 35: Code mẫu phương pháp Pyramid Images.

- Ảnh động minh họa hiệu quả hình ảnh sau khi được Pyramid Images và chạy Sliding Window Minh họa.

### 3.3 Thể hiện bounding box và classname:

- Xem thêm giải thích từ record [01:53:54-01:57:40]



```

1 def visualize_bbox(img, bboxes, label_encoder):
2     img = cv2.cvtColor(img, cv2.COLOR_BGR2RGB)
3
4     for box in bboxes:
5         xmin, ymin, xmax, ymax, predict_id, conf_score = box
6
7         cv2.rectangle(img, (xmin, ymin), (xmax, ymax), (0, 255, 0), 2)
8
9         classname = label_encoder.inverse_transform([predict_id])[0]
10        label = f"{classname} {conf_score:.2f}"
11
12        (w, h), _ = cv2.getTextSize(label, cv2.FONT_HERSHEY_SIMPLEX, 0.6, 1)
13
14        cv2.rectangle(img, (xmin, ymin - 20), (xmin + w, ymin), (0, 255, 0), -1)
15
16        cv2.putText(img, label, (xmin, ymin - 5), cv2.FONT_HERSHEY_SIMPLEX, 0.6, (0, 0, 0), 1)
17
18    plt.imshow(img)
19    plt.axis('off')
20    plt.show()

```

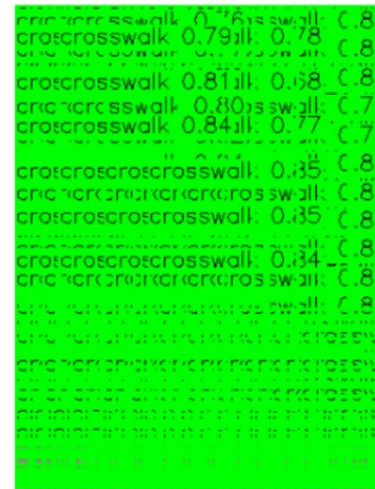
### 3.4 Dự đoán:

- Đưa hình chương trình để cho ra bounding box cùng với classname kèm theo. Ta thấy hình ảnh cho ra bị vắn đề là chồng chéo tất cả các dự đoán.

```

for pyramid_img_info in pyramid_imgs:
    pyramid_img, scale_factor = pyramid_img_info
    window_lst = sliding_window(
        pyramid_img,
        window_sizes=window_sizes,
        stride=stride,
        scale_factor=scale_factor
    )
    for window in window_lst:
        xmin, ymin, xmax, ymax = window
        object_img = pyramid_img[ymin:ymax, xmin:xmax]
        preprocessed_img = preprocess_img(object_img)
        normalized_img = scaler.transform([preprocessed_img])[0]
        decision = clf.predict_proba([normalized_img])[0]
        predict_id = np.argmax(decision)
        conf_score = decision[predict_id]
        xmin = int(xmin / scale_factor)
        ymin = int(ymin / scale_factor)
        xmax = int(xmax / scale_factor)
        ymax = int(ymax / scale_factor)
        bboxes.append(
            [xmin, ymin, xmax, ymax, predict_id, conf_score]
        )

```



56

Hình 36: Lỗi khi dự đoán.

- Để khắc phục lỗi này, ta đưa vào giá trị confidence threshold để lọc: hình nào có confidence score lớn hơn thì thể hiện.



```

for window in window_lst:
    xmin, ymin, xmax, ymax = window
    object_img = pyramid_img[ymin:ymax, xmin:xmax]
    preprocessed_img = preprocess_img(object_img)
    normalized_img = scaler.transform([preprocessed_img])[0]
    decision = clf.predict_proba([normalized_img])[0]
    if np.all(decision < conf_threshold):
        continue
    else:
        predict_id = np.argmax(decision)
        conf_score = decision[predict_id]
        xmin = int(xmin / scale_factor)
        ymin = int(ymin / scale_factor)
        xmax = int(xmax / scale_factor)
        ymax = int(ymax / scale_factor)
        bboxes.append(
            [xmin, ymin, xmax, ymax, predict_id, conf_score]
        )

```



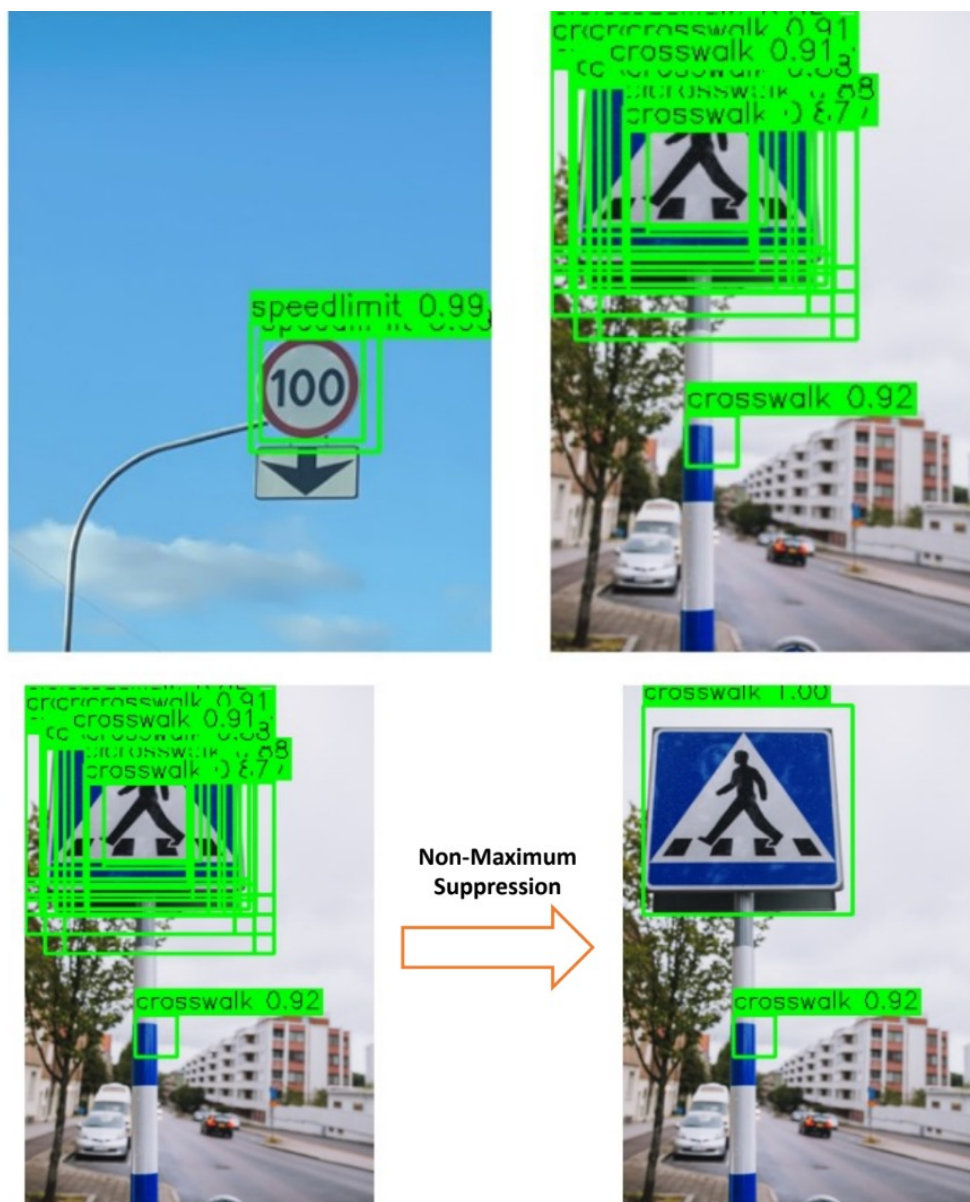
Results with a confidence score threshold = 0.95

Hình 37: Thêm giá trị lọc Confidence Threshold.

Xem thêm giải thích code từ record [01:58:09-02:05:00]

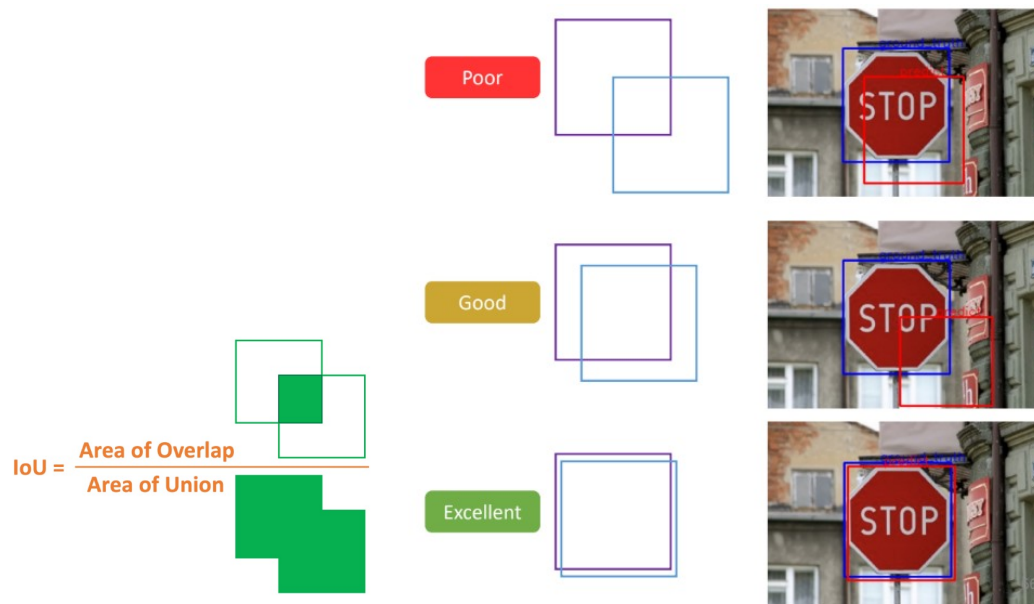
### 3.5 Non-maximum suppression:

- Trong một số trường hợp nhiều bounding box sẽ dự đoán cùng 1 object gây ra hiện tượng đè nhau (Overlapped). -> dùng phương pháp Non-Maximum Suppression để loại bỏ những bounding box đè nhau.
- Chú ý: Non-maximum Suppression chỉ lọc đè nhau không lọc đúng sai của label.



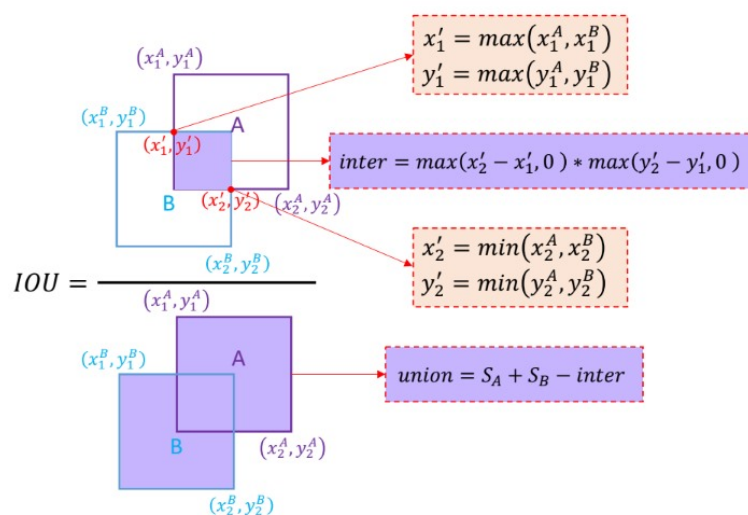
Hình 38: Overlapped bounding box problem.

- Ý tưởng của Non-Maximum Suppression:
  - Sắp xếp các bounding box giảm dần theo giá trị confidence score.
  - Tính IoU giữ giá trị bounding box có confidence score lớn nhất và tất cả những cái còn lại.
  - Giữ lại những bounding box có  $\text{IoU} < \text{IoU threshold}$ .
- IoU: Intersection over Union là thông số đo độ đè nhau giữa 2 bounding box. Cách tính IoU được đề cập ở hình 40 cùng với code kèm theo.



Hình 39: IoU và các trạng thái.

### ❖ Intersection over Union



```

1 def iou_bbox(box1, box2):
2     """
3     - Input:
4       - box1 (x1, y1, x2, y2)
5       - box2 (x1, y1, x2, y2)
6     - Output:
7       - iou of box1 and box2
8     """
9
10    b1_x1, b1_y1, b1_x2, b1_y2 = box1[0], box1[1], box1[2], box1[3]
11    b2_x1, b2_y1, b2_x2, b2_y2 = box2[0], box2[1], box2[2], box2[3]
12
13    # Step 1: Compute inter area
14    x1 = max(b1_x1, b2_x1)
15    y1 = max(b1_y1, b2_y1)
16    x2 = min(b1_x2, b2_x2)
17    y2 = min(b1_y2, b2_y2)
18
19    inter = max((x2 - x1), 0) * max((y2 - y1), 0)
20
21    # Step 2: Compute union area
22    box1Area = abs((b1_x2 - b1_x1) * (b1_y2 - b1_y1))
23    box2Area = abs((b2_x2 - b2_x1) * (b2_y2 - b2_y1))
24
25    union = float(box1Area + box2Area - inter)
26
27    iou = inter / union
28    return iou

```

Hình 40: Công thức tính IoU.

- Code của Non-maximum suppression:

```

1 def nms(bboxes, iou_threshold):
2     if not bboxes:
3         return []
4     scores = np.array([bbox[5] for bbox in bboxes])
5     sorted_indices = np.argsort(scores)[::-1]
6
7     xmin = np.array([bbox[0] for bbox in bboxes])
8     ymin = np.array([bbox[1] for bbox in bboxes])
9     xmax = np.array([bbox[2] for bbox in bboxes])
10    ymax = np.array([bbox[3] for bbox in bboxes])
11
12    areas = (xmax - xmin + 1) * (ymax - ymin + 1)
13
14    keep = []
15    while sorted_indices.size > 0:
16        i = sorted_indices[0]
17        keep.append(i)
18
19        iou = compute_iou(
20            [xmin[i], ymin[i], xmax[i], ymax[i]],
21            np.array(
22                [
23                    xmin[sorted_indices[1:]],
24                    ymin[sorted_indices[1:]],
25                    xmax[sorted_indices[1:]],
26                    ymax[sorted_indices[1:]]
27                ]
28            ).T,
29            areas[i],
30            areas[sorted_indices[1:]]
31        )
32
33        idx_to_keep = np.where(iou <= iou_threshold)[0]
34        sorted_indices = sorted_indices[idx_to_keep + 1]
35    return [bboxes[i] for i in keep]

```

Hình 41: Code mẫu hàm Non-maximum suppression.

Xem thêm giải thích code từ record [02:19:31-02:26:12]. Từ sau trở về cuối record là phần hỏi đáp lại thắc mắc của toàn bài.

Kết quả cho ra:



Hình 42: Kết quả cuối cùng.