# Assignment 2 - SSD

Due date: 23:59 Thursday  03/12/2020

In this assignment, you will implement SSD, another object detection network. SSD is similar to YOLO since they both divide the image into grids and output bounding boxes in each cell. However, SSD is more complicated. It divides grids at different scales, and uses default bounding boxes as anchors.  Similar to YOLO, we will cover the four main components step by step: 1. Generating default bounding boxes; 2. Assigning ground truth objects to default bounding boxes; 3. The network and the loss function; 4. Non maximum suppression and evaluation. If you have questions, feel free to email me (Zhiqin Chen, zhiqinc AT sfu.ca), or come to my office hour (Wednesdays, 16:00-17:00, ASB9808).

I recommend you to read the SSD paper first.
SSD: Single Shot MultiBox Detector
https://arxiv.org/pdf/1512.02325.pdf

The SSD you are going to implement in this assignment is a simplified version. If you see conflicts between this guideline and the paper, please follow the guideline.

You can download the dataset and the framework here:
Dataset: https://drive.google.com/file/d/1UZuboJ32HpcdXei7yEEPThizVJz2-afR
Framework: https://drive.google.com/open?id=1IEY3bipSJSwh_Kiv8PFOxwr5NCxJCVp0

# 1. Grading

1. Plagiarism or cheating will not be tolerated. You are allowed to discuss with others, but do not copy their code.
2. Submit your report, code and result in Coursys. Zip your code (at least four .py files) and your results (txts containing bounding boxes) into code2.zip. Export your report into report2.pdf. Submit code2.zip and report2.pdf in Coursys.
3. You cannot submit your network weights. But you must keep them until the semester ends. If I find something funny about your code or result, I may request you to email me the weights. If your weights fail to produce the exact same result as you submitted, you will receive 0 grade.
4. Bonus points: if you find a severe mistake in this assignment, please email me ASAP. you will get up to 5 bonus points if you are the first to report the mistake and the mistake is so bad that I have to email everyone to correct it.
5. This assignment has 15 points as following:
  - 6 key points. You need to implement those key points in your code or your points will be deducted. (6 points)
    - Correctly complete function *default_box_generator*. (1 point)
    - Correctly complete function *match*. (1 point)

- ○ Correctly implement data augmentation (at least cropping). (1 point)
  - ○ Correctly define the network and forward. (1 point)
  - ○ Correctly define loss function. (1 point)
  - ○ Correctly implement non maximum suppression. (1 point)
  - ○
- ● Training accuracy measured in F1 score. (3 points)
  - ○ Cat  >0.7 (0.5 point),  >0.5 (0.5 point)
  - ○ Dog  >0.7 (0.5 point),  >0.5 (0.5 point)
  - ○ Person  >0.7 (0.5 point),  >0.5 (0.5 point)
- ● Testing accuracy measured in F1 score. (3 points)
  - ○ Cat  >0.4 (0.5 point),  >0.2 (0.5 point)
  - ○ Dog  >0.3 (0.5 point),  >0.1 (0.5 point)
  - ○ Person  >0.6 (0.5 point),  >0.4 (0.5 point)
- ● Your report. (3 points)
  - ○ Include all figures and remarks in your report. Each missing figure or remark will account for 0.5 deducted point.
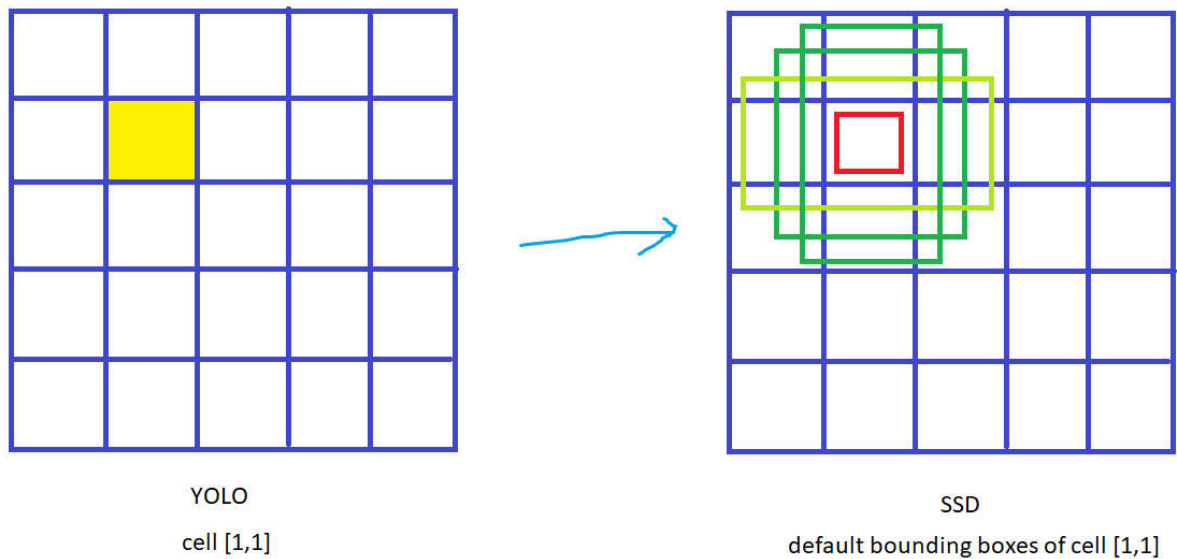
# 2. Overview

If you have not implemented YOLO yet, please implement YOLO in the tutorial sessions. The hand-out for YOLO can be downloaded here:
https://drive.google.com/open?id=154GBlzFU1eTMN9RmHkGhL-KEom3P8bvNeJHYT2WqJlY
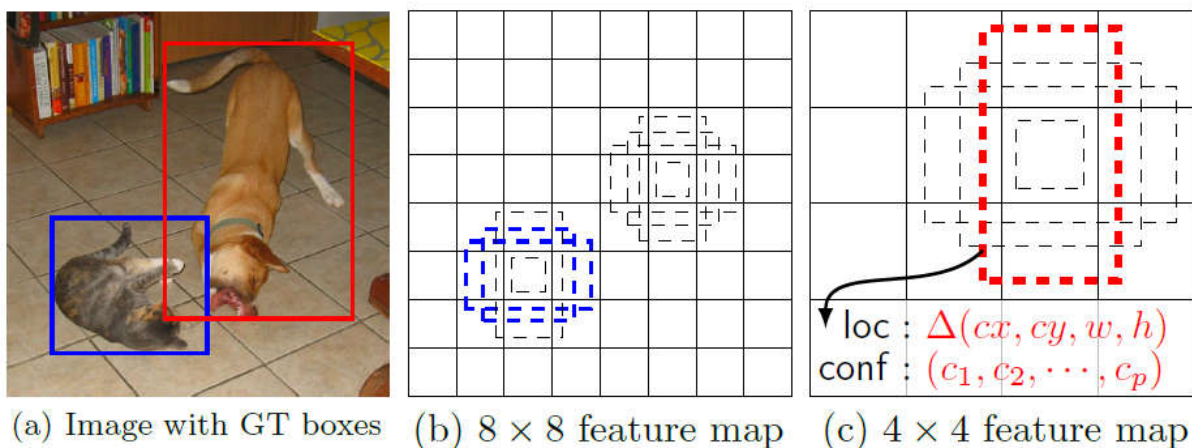
The main difference between YOLO and SSD is the default bounding boxes. In YOLO, we divide the image into 5x5 grid cells. Each cell has its confidence, or probabilities for each class (P(cat), P(dog), P(person), P(background)). But the cell does not have a default bounding box – we need to specify the absolute width and height of the ground truth bounding box and use them for training.

For SSD, each cell will provide some default bounding boxes. In this assignment, each cell will provide four default bounding boxes as follows:

YOLO

cell [1,1]

SSD

default bounding boxes of cell [1,1]

As you see, if for YOLO, we have at most 25 bounding boxes, then for SSD, we will have 25*4 = 100 bounding boxes. Furthermore, we will introduce grids for multiple scales. In YOLO, we only have one 5x5 grid in the output layer, but for SSD in this assignment, we will have four output branches: 10x10, 5x5, 3x3 and 1x1, a total of 135 cells and 540 default bounding boxes.



(a) Image with GT boxes   (b) 8 × 8 feature map   (c) 4 × 4 feature map

This image is stolen from the SSD paper and it demonstrates why these default bounding boxes are helpful. We want to assign the ground truth bounding box to some of the default bounding boxes so that the default bounding boxes are already good estimates of the ground truth bounding box.

And since we have the default bounding boxes, we can use them as "anchors". That is, we do not need to predict the absolute width and height as in YOLO; we can just predict the relative positions and sizes with respect to the default bounding box, as shown in the above figure, Δ(cx,cy,w,h).

In all, the main challenge in this assignment (SSD) compared to the tutorial (YOLO) is how to generate the default bounding boxes and how to assign ground truth objects to those default bounding boxes. The network is similar to that of YOLO, and the loss and

non-maximum suppression are almost the same as YOLO. In a word, if you have completed YOLO in the tutorial, you are going to have a much easier time in this assignment.

# 3. Generating default bounding boxes

As said, in this assignment you are not dealing with 5x5 or 25 cells, you are going to work on (10x10+5x5+3x3+1x1)*4=540 default bounding boxes. The first step is to generate them.

You need to complete function *default_box_generator* in *dataset.py*. The function takes a series of parameters and eventually outputs a 540x8 array, storing 540 bounding boxes. The last dimension 8 means the 8 attributes of each default bounding box: [x_center, y_center, box_width, box_height, x_min, y_min, x_max, y_max]. Note that here, all values are absolute positions and sizes. For example, if a default bounding box has center (x=0.3, y=0.4), width=0.1 and height=0.2, the attributes you need to store are [0.3,0.4, 0.1,0.2, 0.25,0.3, 0.35,0.5].

You need to generate 4 default bounding boxes for each grid cell in each grid (10,5,3,1), using the provided scales *large_scale* and *small_scale* to determine the sizes of the default bounding boxes. There may be bounding boxes exceeding the image boundary, therefore you may clip them so that the bounding boxes stay inside the image.

For example, consider filling the default boxes for the first cell.
The size of the grid is 10 since *layers[0]*=10.
The scale of the three large boxes is lsize = *large_scale[0]* = 0.2.
The scale of the one small box is ssize = *small_scale[0]* = 0.1.
The first cell is (0,0) in a 10x10 cell grid.
Generate a box with width and height [ssize,ssize].
Generate a box with width and height [lsize,lsize].
Generate a box with width and height [lsize*sqrt(2),lsize/sqrt(2)].
Generate a box with width and height [lsize/sqrt(2),lsize*sqrt(2)].
All the four above boxes are centered at the center of the first cell (0.5/10, 0.5/10)
The four boxes you get for the first cell is:
[0.05, 0.05, 0.1,  0.1,   0,    0,    0.1,  0.1]
[0.05, 0.05, 0.2,  0.2,  -0.05, -0.05, 0.15, 0.15]
[0.05, 0.05, 0.28, 0.14, -0.09, -0.02, 0.19, 0.12]
[0.05, 0.05, 0.14, 0.28, -0.02, -0.09, 0.12, 0.19]
The four boxes after clipping is: (optional)
[0.05, 0.05, 0.1,  0.1,  0,  0,  0.1,  0.1 ]
[0.05, 0.05, 0.2,  0.2,  0,  0,  0.15, 0.15]
[0.05, 0.05, 0.28, 0.14, 0,  0,  0.19, 0.12]
[0.05, 0.05, 0.14, 0.28, 0,  0,  0.12, 0.19]

You need to create boxes for all cells in a similar manner. You can modify the box sizes in *large_scale* and *small_scale* to see if you get better results when training the network.

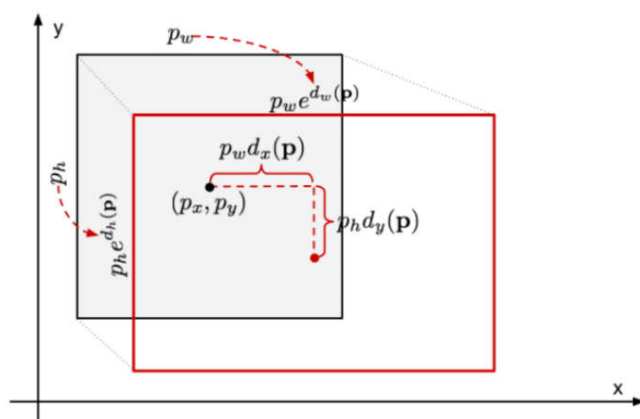# 4. Assigning ground truth objects to default bounding boxes

This part is done in the dataloader (*dataset.py*). You need to read the image and ground truth bounding boxes, and then return the resized and transposed image, the ground truth probabilities (we call it confidence) *ann_confidence*, and the ground truth bounding boxes *ann_box*.

*ann_confidence* is 540x4, since we have 540 default boxes and 4 classes (cat, dog, person, background). *ann_confidence* should be 540 one-hot vectors.

*ann_box* is 540x4, since we have 4 attributes for a bounding box: [relative_center_x, relative_center_y, relative_width, relative_height]. The attributes are all relative to the default bounding box as follows:

## Bounding Box Regression

- Benefit is that all $d_i(p)$ where $i \in \{x, y, w, h\}$ attain values between $[-\infty, +\infty]$. The targets for them to learn are:

$$t_x = (g_x - p_x)/p_w$$
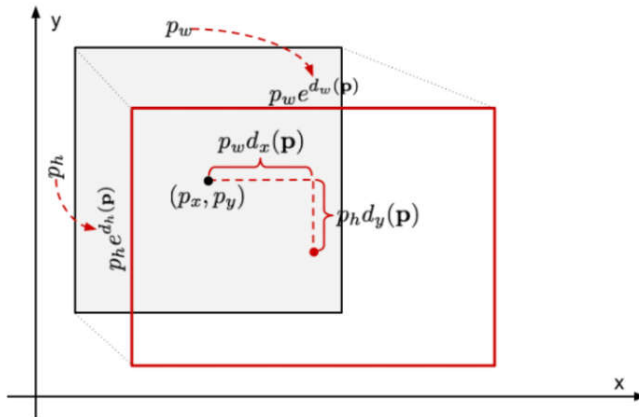$$t_y = (g_y - p_y)/p_h$$
$$t_w = \log(g_w/p_w)$$
$$t_h = \log(g_h/p_h)$$

I am stealing this slide from Ali's lecture slides, but I am pretty sure he stole it from somewhere else. Suppose we have a default box [px,py,pw,ph], which are exactly the first four attributes of our generated default bounding boxes. Also, we have a ground truth object bounding box [gx,gy,gw,gh]. The 4 attributes we need for *ann_box* are [tx,ty,tw,th], the relative positions and sizes of the ground truth bounding box with respect to the default bounding box.

If you need to recover the ground truth bounding box, or to show the predicted bounding box of your network, you can do an inverse process as follows, where [dx,dy,dw,dh] is the predicted relative attributes from your network.

# Bounding Box Regression

- Regressor learns a scale-invariant transformation between two centers and log-scale transformation between widths and heights.



$$\hat{g}_x = p_w d_x(\mathbf{p}) + p_x$$
$$\hat{g}_y = p_h d_y(\mathbf{p}) + p_y$$
$$\hat{g}_w = p_w \exp(d_w(\mathbf{p}))$$
$$\hat{g}_h = p_h \exp(d_h(\mathbf{p}))$$

Then, how to determine if a default bounding box is carrying an object? Different from YOLO, this time you need to assign one object to multiple default boxes. We will say a default box is carrying an object, if the ground truth bounding box of this object has an IOU greater than a threshold (0.5 in this assignment) with the default box. There may be cases where the ground truth box does not have sufficient overlap with any of the default boxes, in that case, we assign the object to the default box that has the largest IOU with the object bounding box, to make sure at least one default bounding box is used for each object.

You need to implement function *match* to process each bounding box and update the entries of *ann_confidence* and *ann_box*. You only need to fill in the bounding box attributes for boxes that carry objects. For empty boxes, you can ignore them since they are not used when training the network.

Other things to do:

1. You can split the dataset into 80% training and 20% testing here (or 90% - 10%), by slicing *self.img_names* with respect to *self.train*. You can use all images for training of course, but you will not know whether your network has overfitted without a test or validation set. Your model may have very good performance on the training set, but have very poor performance on the testing set. Keep in mind that both training accuracy and testing accuracy are considered when grading your assignment.

2. Data augmentation. Usually the networks are trained on millions of images. It is uncommon to train an object detection network on only ~6000 images. But I have to reduce the size of the data so you can finish training within 2 hours. Therefore you have a high chance of overfitting on the training set. You will need to augment the data to mitigate that. Please implement random cropping first (make sure the cropped patch is large enough to cover the objects). You can try adding other augmentations to get better results. If you don't know what random cropping is, Google that. Note that you need to augment the image and

the bounding boxes at the same time. Make sure the augmented bounding boxes still apply to the augmented image.

# 5. The network and the loss function

Please open model.py and implement the network on the next page. The network is the same as YOLO until you reach resolution 10x10. Be sure to include a bias term for your convolution layer.

Please pay close attention to how the four output branches are concatenated, especially their ordering. You need to make sure that each entry in your 540 output boxes actually corresponds to the entry in the 540 default boxes defined in the previous section. For example, if you define *ann_confidence[4]* and *ann_box[4]* to be a box in cell (0,1) of the 10x10 grid, your network output must have *confidence[4]* and *bboxes[4]* correspond to the cell (0,1) of the 10x10 output branch. This part is very error-prone.

The loss function is almost the same as YOLO. You may notice that I only replaced (i,j) with (i).

$$L_{cls} = \frac{1}{\sum_i x_i^{obj}} \sum_i x_i^{obj} cross\_entropy(conf_i^{pred} , conf_i^{gt}) + 3 \cdot \frac{1}{\sum_i x_i^{noobj}} \sum_i x_i^{noobj} cross\_$$

$$L_{box} = \frac{1}{\sum_i x_i^{obj}} \sum_i x_i^{obj} smooth_{L1}(box_i^{pred} - box_i^{gt})$$

$$L_{SSD} = L_{cls} + L_{box}$$

$\sum_i$ will iterate through all default boxes, i.e., i = 0, 1, 2, ... , 539.

$x_i^{obj} = 1$ and $x_i^{noobj} = 0$ if default box (i) is carrying an object in the ground truth.

$x_i^{obj} = 0$ and $x_i^{noobj} = 1$ if default box (i) is empty in the ground truth.

$conf_i^{pred}$ is the class probabilities of default box (i) outputted by the network.

$conf_i^{gt}$ is the ground truth class probabilities of default box (i).

$box_i^{pred}$ is the bounding box attributes of default box (i) outputted by the network.

$box_i^{gt}$ is the ground truth bounding box attributes of default box (i).

Note that the weights for classification loss are 1 and 3 for boxes with objects and without.

You need to figure out how you can get the indices of all boxes carrying objects, and use confidence[indices], box[indices] to select those boxes. If your implementation is correct, after you get those indices, your code to compute loss should be no more than three lines.

Note that the average operation ( $\frac{1}{\sum_i x_i^{obj}} \sum_i x_i^{obj}$ ...) is done automatically by the built-in

functions *F.binary_cross_entropy* and *F.smooth_l1_loss*.

Input image

[N,3,320,320]

conv( 3, 64, 3, 2),bn,relu

[N,64,160,160]

conv( 64, 64, 3, 1),bn,relu

[N,64,160,160]

conv( 64, 64, 3, 1),bn,relu

[N,64,160,160]

conv( 64,128, 3, 2),bn,relu

[N,128,80,80]

conv(128,128, 3, 1),bn,relu

[N,128,80,80]

conv(128,128, 3, 1),bn,relu

[N,128,80,80]

conv(128,256, 3, 2),bn,relu

[N,256,40,40]

conv(256,256, 3, 1),bn,relu

[N,256,40,40]

conv(256,256, 3, 1),bn,relu

[N,256,40,40]

conv(256,512, 3, 2),bn,relu

[N,512,20,20]

conv(512,512, 3, 1),bn,relu

[N,512,20,20]

conv(512,512, 3, 1),bn,relu

[N,512,20,20]

conv(512,256, 3, 2),bn,relu

[N,256,10,10]

N is batch size
bn is batch normalization layer
relu is ReLU activation layer
conv(cin,cout,ksize,stride) is convolution layer
- cin - the number of input channels
- cout - the number of output channels
- ksize - kernel size
- stride - stride
- padding - you need to figure this out by yourself

conv(256,256, 1, 1),bn,relu

[N,256,10,10]

conv(256,256, 3, 2),bn,relu

[N,256,5,5]

conv(256, 16, 3, 1)        conv(256, 16, 3, 1)

[N,16,10,10]        [N,16,10,10]

reshape        reshape

[N,16,100]        [N,16,100]

conv(256,256, 1, 1),bn,relu

[N,256,5,5]

conv(256,256, 3, 1),bn,relu

[N,256,3,3]

conv(256, 16, 3, 1)        conv(256, 16, 3, 1)

[N,16,5,5]        [N,16,5,5]

reshape        reshape

[N,16,25]        [N,16,25]

conv(256,256, 1, 1),bn,relu

[N,256,3,3]

conv(256,256, 3, 1),bn,relu

[N,256,1,1]

conv(256, 16, 3, 1)        conv(256, 16, 3, 1)

[N,16,3,3]        [N,16,3,3]

reshape        reshape

[N,16,9]        [N,16,9]

conv(256, 16, 1, 1)        conv(256, 16, 1, 1)

[N,16,1,1]        [N,16,1,1]

reshape        reshape

[N,16,1]        [N,16,1]

concatenate        concatenate

[N,16,135]        [N,16,135]

permute        permute

[N,135,16]        [N,135,16]

reshape        reshape

[N,540,4]        [N,540,4]

softmax

[N,540,4]

Output
bboxes

Output
confidence

# 6. Non maximum suppression and evaluation

This part is already covered in the lecture and in the tutorial sessions. Please refer to Ali's slides and the hand-out for YOLO. Please complete the function in *utils.py*.
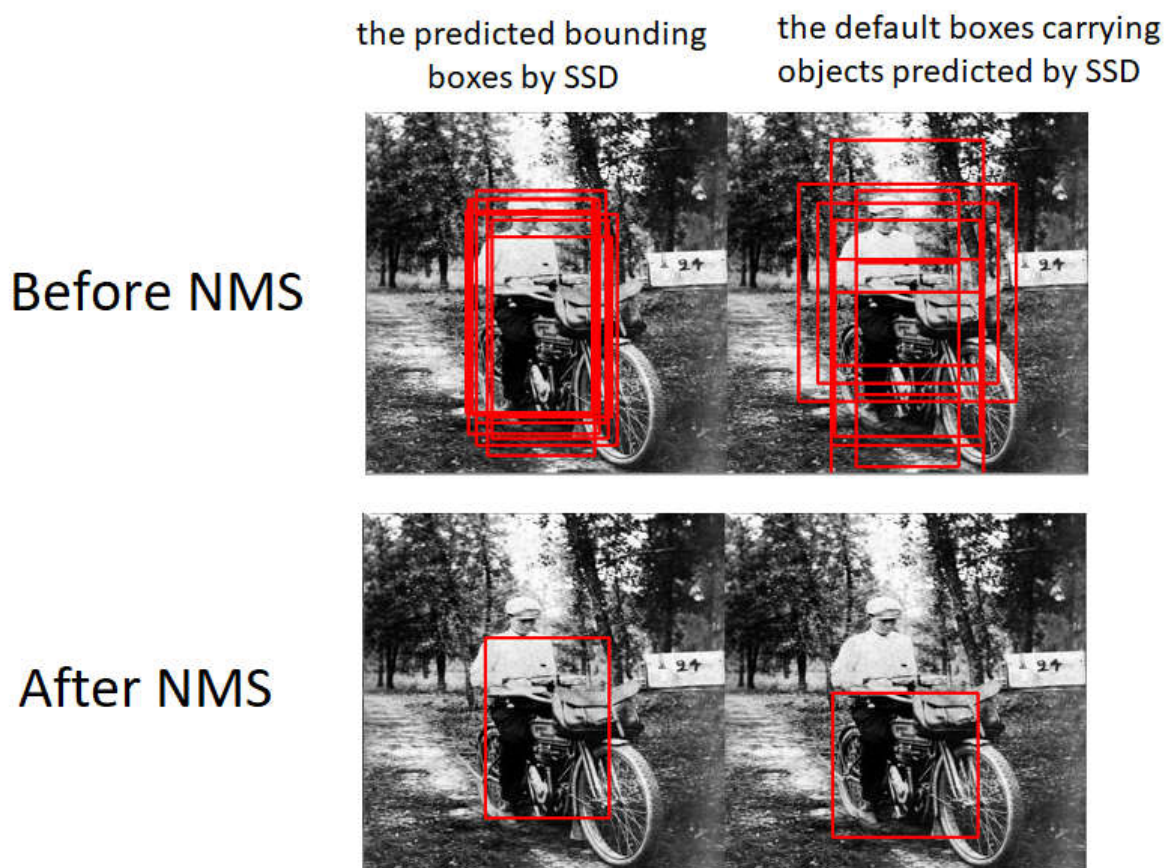
Non maximum suppression (NMS):
https://www.dropbox.com/s/hiatyeyturppboy/Week%205.pdf?dl=0

Precision, Recall, mAP:
https://www.dropbox.com/s/b591p7iwtflcgac/Object%20Detection.pdf?dl=0

An example of NMS is shown below:



A brief pipeline of NMS

Keep two lists:
A = all predicted bounding boxes; B = [] (non-overlapping boxes to be returned).

1. Select the bounding box in A with the highest probability in class cat, dog or person.
2. If that highest probability is greater than a threshold (threshold=0.5), proceed; otherwise, the NMS is done.
3. Denote the bounding box with the highest probability as x. Move x from A to B.

4. For all boxes in A, if a box has IOU greater than an overlap threshold (overlap=0.5) with x, remove that box from A.
5. Jump to 1.

## Evaluation

You are encouraged to plot a precision-recall curve and compute mAP or AUC. However, to make my grading easier, we are going to use a simpler evaluation metric, F1 score. Basically, you need to tune the parameters (threshold, overlap, etc) so that your network produces the best result you think you can get. For each image in the given training set (6392 images) and in the given testing set (711 images), you need to produce a txt file recording the predicted bounding boxes by your network after NMS, in the same format as the given annotation files, such as:
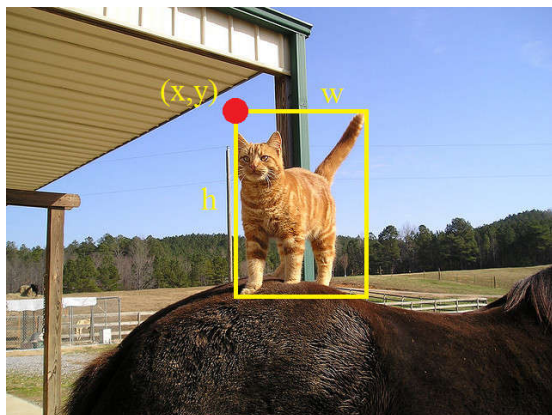
File name:  SSD_framework/predicted_boxes/test/00000.txt
First line:        1 406.3 197.82 171.64 167.76
Second line:   0 360.4 287.91 461.71 677.61
Third line:        … …

The name of your output txt file should be the same as the name of the input image. The content of the txt file should be that each line describes a bounding box, the first number is class id (cat=0, dog=1, person=2), the rest four are x,y,w,h as explained in the first tutorial session:
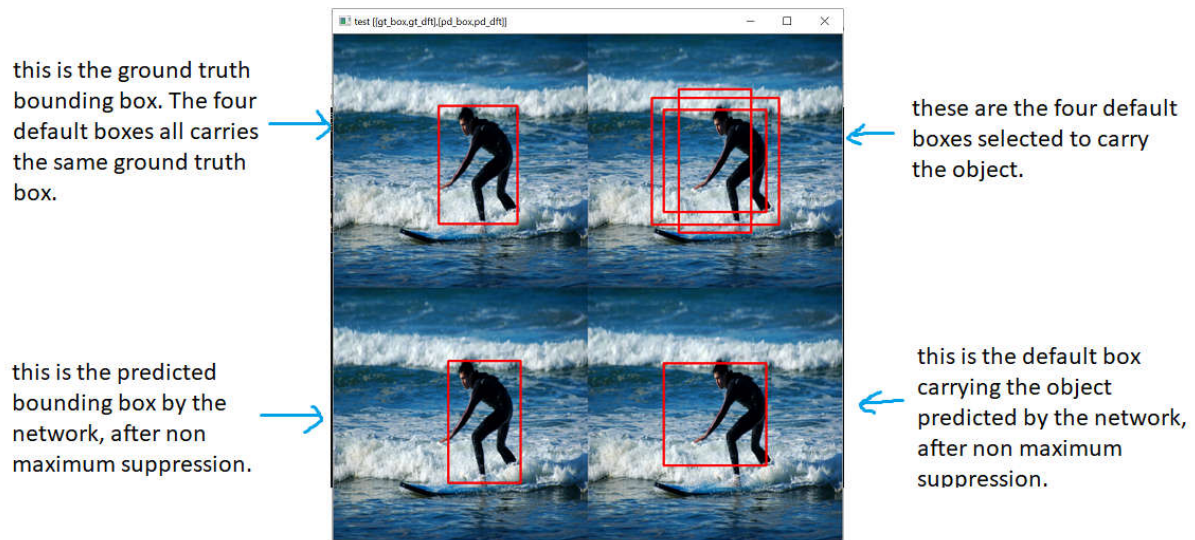


You will need to submit your code and all (training and testing set, 6392+711) generated txt files in this assignment, as a zip file.

The evaluation will be done on F1 score. I will compare your predicted boxes and ground truth boxes and use overlap threshold (IOU) 0.5 to compute the precision and recall. The final F1 score is computed as $(2 * precision * recall) / (precision + recall)$. Therefore, although not compulsory, I recommend you to implement some function in your code or in a separate .py file to compute the F1 score of your network outputs, so that you have an estimate of your grades.

# 7. Report

You need to complete *visualize_pred* function in *utils.py*, and use the function to generate the visualizations required by the report. You can either save the image or take screenshots. Each image you put into the report should be composed of 4 small images like the sample visualization below:



You need to put into the report
- One visualization of the network output on an image with at least one cat object, before non maximum suppression.
- One visualization of the network output on an image with at least one dog object, before non maximum suppression.
- One visualization of the network output on an image with at least one person object, before non maximum suppression.
- One visualization of the network output on an image with at least two person objects, before non maximum suppression.
- One visualization of the network output on an image with at least one cat object, after non maximum suppression.
- One visualization of the network output on an image with at least one dog object, after non maximum suppression.
- One visualization of the network output on an image with at least one person object, after non maximum suppression.
- One visualization of the network output on an image with at least two person objects, after non maximum suppression.
- The training time of your network, the hardware settings (GPU, server or laptop, etc).
- Approximately how long it takes for you to finish the assignment.
- How difficult (or easy) you think the assignment is and why.
- Any other suggestions for this assignment. (optional)

For each figure, be sure to include a title describing what the figure shows. You do not need to add any other texts to the figure.