# Yolo

Yolo divides the image into grid cells. Each grid cell predicts only one object. We will assign center of that object to that grid cell. So this grid cell is responsible for predicting the object                     (only convolutional layers) no pooling

For each grid cell we predict $n \times n \times 8$     n → no. of classes
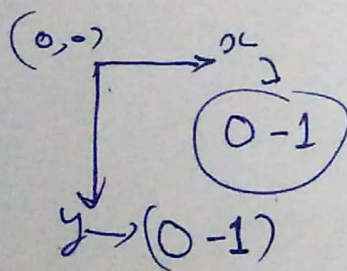
$P_c$ → confidence value

$\left.\begin{matrix} b_x \\ b_y \\ b_h \\ b_w \end{matrix}\right\}$ → bounding box of object if there is.

$\left.\begin{matrix} C_1 \\ C_2 \\ C_3 \end{matrix}\right\}$ Say we have 3 classes.

$\begin{bmatrix} 0 \\ ? \\ ? \\ \vdots \\ ? \\ ? \\ ? \\ ? \end{bmatrix}$ → if object not present

$3 \times 3 \times 8$ → feature map prediction

$(0, 0)$ → $x$

$(0-1)$

$y → (0-1)$

✿ What if more than object's centre lies in same grid cell?

Here comes **Anchor Boxes**    as there is limitation with only having grid cells.

To solve problem we will introduce the concept of anchor box which makes yolo to predict multiple objects centered in one cell.

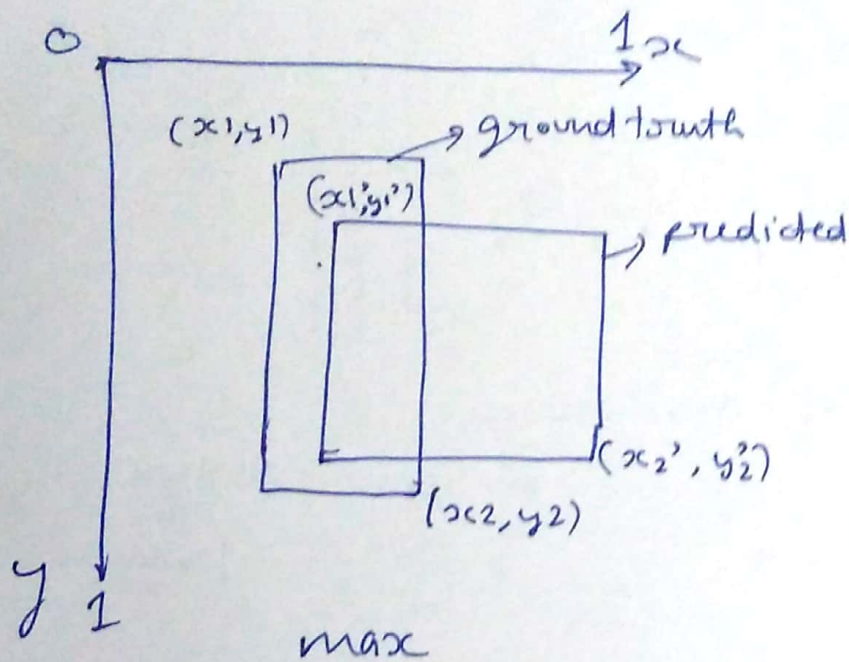Different objects have different shapes so we use anchor box of different shapes.

$(N \times N) \times [\text{num of anchors} \times (5 + \text{no. of classes})]$

NxN grid cells   say we have 3x3 grid cells.

# IoU

for evaluation
& for inference.

Instead of defining box by center, width & height we
define it using two corners (upper left & bottom right)

$$IOU = \frac{\text{Intersection}}{\text{Union}}$$



$(x_1, y_1)$

ground truth

$(x_1', y_1')$

predicted

$(x_2', y_2')$

$(x_2, y_2)$

max

# Non max suppression

used for cleaning up when multiple grid cells are predicted the same object.

**Steps:**
- Discard all boxes with $p_c$ ~~test~~ $\leq 0.6$
- Pick the box with largest $p_c$ output as prediction
- Discard any remaining box with IoU $\geq 0.5$

**So:**

Note that given an input image with $3 \times 3$ grid cells 2 anchor boxes. each grid cell will have 2 predictions even for those grid cells that don't have any object inside.

In this case we will filter by class scores.  *threshold*

$$p_c \quad b_x \quad b_y \quad b_w \quad b_h \quad c_1 \quad c_2 \quad c_3$$

$$\begin{bmatrix} p_c \times c_1 \\ p_c \times c_2 \\ p_c \times c_3 \end{bmatrix} \longleftarrow$$

**☆ fully convolutional**

- no **pooling** to downsample convolution layer with stride 2 is used
  → to prevent low-level features

(!)

In yolo v3 we don't predict one feature map. We have **3** of these.

$13 \times 13 \qquad 26 \times 26 \qquad 52 \times 52$

**different scales**

$416 \times 416$

In yolo v1 & v2 they only have 1

(32, 16, 8) → stride  ☆ ← Its very helpful for predicting small objects.

- 3 __bounding boxes__ for each scale
  anchor.
- 80 classes
  ↳ coco dataset { object detection
  ⎨ image captioning.

Prediction done by convolution layer <u>1×1 conv</u>. → reduce feature map depth to →

out put → feature map

Depth wise entries in feature map ←

$$(B \times (S+C))$$

⟵ 1×1 conv

$$(N \times N) \times (\text{num anchors} \times (S + \overset{no.}{\text{classes}}))$$

e.g. $\underset{\partial\,\text{Size}}{13 \times 13} \times (3 \times (S + 80))$

$$26 \times 26 \times (3 \times (S + 80))$$

## ✡ Prediction of Anchor Boxes

→ predict width & height of bounding box but it leads to unstable gradients.

→ So we use <u>offset</u> (how much we should move in order to predict desired bounding box)
to pre-defined default bounding box.

→ yolo v3 has 3 anchors which results in predicting 3 bounding boxes per cell.

$t_x, t_y, t_w, t_h$        $(c_x, c_y)$
(offset of centre)          ↳ top left coordinate of grid.
                            (to get value b/w 0 & 1)
                            σ → sigmoid

$$\begin{cases} b_x = \sigma(t_x) + c_x \\ b_y = \sigma(t_y) + c_y \end{cases}$$

$p_w, p_h$
↳ anchor dimension of box

$b_w = p_w e^{t_w}$

$b_h = p_h e^{t_h}$

basically ← adding offset to left top left co-ordinate of grid.
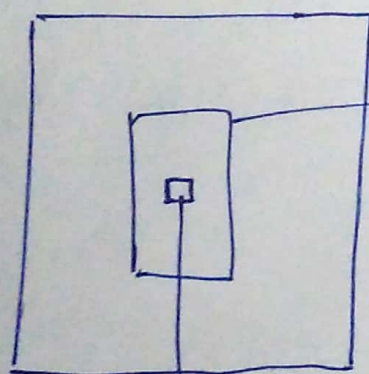
$t_w, t_h$
↳ offsets for width & height

The resultant predictions bw & bh are normalised by height & width of image.

To get them back to original feature map shape multiply by grid that we are ~~predicting~~ using.

Say we are using 13 X13 feature map

then we need to multiply by 13

Denormalize the prediction

We are passing the predict centre offsets with Sigmoid to get value b/w 0 & 1.



→ dog

$\sigma(t_x)$ $\sigma(t_y)$
0.4     0.7

0.4 + 6 = 6.4
0.7 + 6 = 6.7

(6,9) y centre

if not passed

1.2 & 0.7

6 + 1.7 → 7.7   out
0.7 + 6 = 6.7   of
                bound

to keep the center in the grid which is predicting.

2.) yolov1 & v2 uses softmax → as they assume that classes
yolov3 → sigmoid.      are mutually exclusive.
                        i.e. if object belongs to one
                        class then it can't belong
So each class score is   to another class.
predicted using logistic
regression & a threshold is    (sum up to 1)
used to predict multiple labels
for an object.

106 layers fully convolution.

- residual to restore spatial info that's lost

| 79 | 91 | 106 |
|---|---|---|
| 13×13 | 26×26 | 52×52 |

$$\frac{416}{32}$$ → stride

$$\frac{416}{16}$$ → stride

$$\frac{416}{8}$$

to get predicted
feature map
size.

Total no. of bounding boxes yolo predicts.

$$(13 \times 13 \times 3) + (26 \times 26 \times 3) + (52 \times 52 \times 3)$$

$$= 10647$$

yolo v3 → 9 anchor boxes.

3 for each scale

loss { classification, localization, confidence

## Loss

### v2

Center co-ordinates

width/height of bounding box

class confidence (object)→1 $\mathbb{1}^{obj}_{cj} \Sigma \Sigma (C_i - \hat{C}_i)^2$

no object → 0

classes → classification loss ← only penalizes when there is an object

To compute loss for true positive we only one of them to be responsible for the object. For this we consider one with highest IOU.

$$\sum_{i=0}^{S^2} \mathbb{1}_i^{obj} \sum (p_c(c) - \hat{p}_i(c))^2$$

• Penalizes the objectness score prediction for bounding box responsible for predicting obj.
ideally 1

→ penalizes for bounding box having no object (ideally 0)

**Classification loss** → If object detected classification loss at each cell is → Squared error of the class conditional probabilities for each class.

$$\sum_{i=0}^{S^2} \left( \mathbb{1}_i^{obj} \right) \sum_{c \in classes} (p_i(c) - \hat{p}_i(c))^2$$

= 1 if object appears in cell $i$, otherwise 0.

## Localization loss

measures errors in the predicted boundary box locations and sizes. we only take box responsible for detecting the object.

$$\lambda_{coord} \sum_{i=0}^{S^2} \sum_{j=0}^{B} \mathbb{1}_{ij}^{obj} \left[ (x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 \right]$$

$\mathbb{1}_{ij}^{obj} = 1$ ← if the jth boundary box in cell $i$ is responsible for detecting object. otherwise 0

$$\lambda_{coord} \sum_{i=0}^{S^2} \sum_{j=0}^{B} \mathbb{1}_{ij}^{obj} \left[ (\sqrt{w_i} - \sqrt{\hat{w}_i})^2 + (\sqrt{h_i} - \sqrt{\hat{h}_i})^2 \right]$$

we don't want absolute errors in large boxes & small boxes. So yolo predicts square root of bounding box width & height.

To put more emphasis on boundary box accuracy we multiply the loss by $\lambda_{coord}$ (default : 5)

# Confidence loss

If object is detected in the box

$$\sum_{i=0}^{S^2} \sum_{j=0}^{B} 1_{ij}^{obj} \left( C_i - \hat{C}_i \right)^2$$

$\hat{C}_i$ is the box confidence score of box $j$ in cell $i$

$1_{ij}^{obj} = 1$ if $j$th boundary box in cell $i$ is responsible for detecting the object, otherwise 0.

If object is not detected:

$$\lambda_{noobj} \sum_{i=0}^{S^2} \sum_{j=0}^{B} 1_{ij}^{noobj} \left( C_i - \hat{C}_i \right)^2$$

✡ Most boxes don't contain any objects. This causes class imbalance problem i.e. we train the model to detect background more frequently than detecting objects. So we weight this loss down by a factor $\lambda_{noobj}$ (default : 0.5)

---

The last 3 terms in Yolo v2 are the squared errors, whereas in yolo v3, they have been replaced by cross entropy loss error terms.
In classification loss, instead of using mean squared error yolov3 uses bce loss for each label. Also reducing the computation complexity by avoiding the softmax function