Deep Learning for Computer Vision

# CNNs for Object Detection - II

Vineeth N Balasubramanian

Department of Computer Science and Engineering
Indian Institute of Technology, Hyderabad

# Exercise: Smooth L1 Loss

Why is Smooth L1 loss less sensitive to outliers than L2 loss?

If the deviation of predicted output from ground truth is very high, squaring the difference explodes the gradient. This can happen in L2 loss for outliers, and is mitigated in the Smooth L1 loss.

# Recall: Contemporary Object Detection Methods
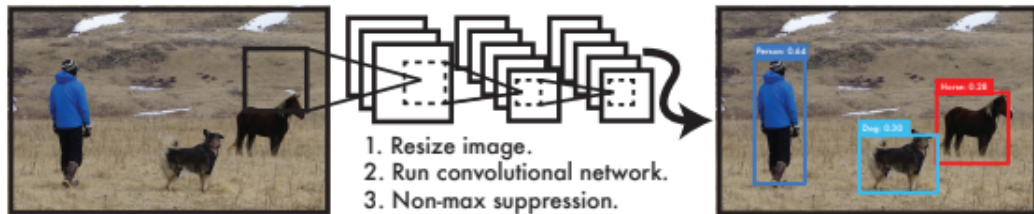
**Region Proposal-based:**

- Two-stage detection framework
- In the first stage, potential object regions are proposed (through methods such as Selective Search or Region Proposal Network, which we will see soon)
- In the second stage, a classifier processes the candidate regions
- More robust in performance but slower

**Dense Sampling-based:**

- One-stage detection framework
- Integrates region proposals and detection by acting on a dense sampling of possible locations
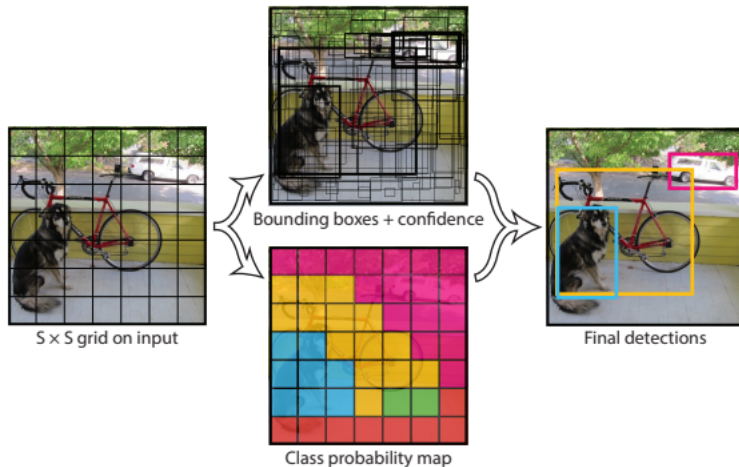- Simple and fast but performance not as good as Region Proposal-based methods

# You Only Look Once (YOLO): v1[1]

- Single-stage detector based on Overfeat
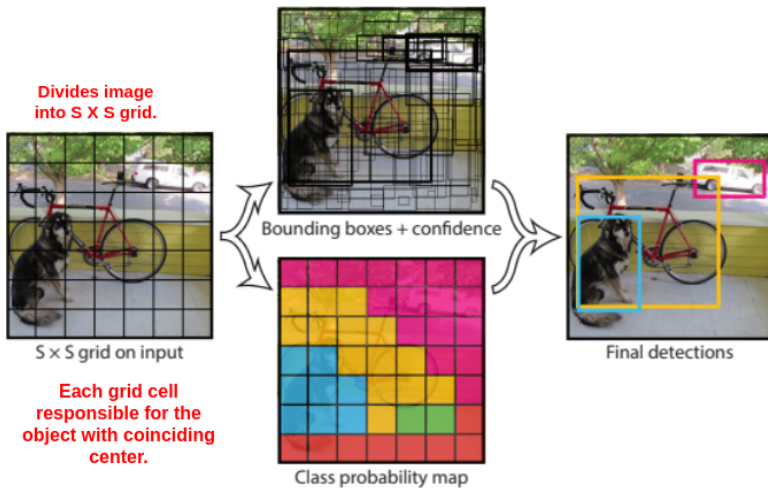- Speed with good performance the main aim



1. Resize image.
2. Run convolutional network.
3. Non-max suppression.

---

[1]Redmon et al, You Only Look Once: Unified, Real-Time Object Detection, CVPR 2016

# YOLO v1: Unified Detection



S × S grid on input

Bounding boxes + confidence

Class probability map

Final detections

# YOLO v1: Unified Detection

**Divides image into S X S grid.**



Bounding boxes + confidence

S × S grid on input

**Each grid cell responsible for the object with coinciding center.**

Class probability map

Final detections

# YOLO v1: Unified Detection



Each grid cell predicte B bounding boxes and confidence scores of the boxes.

S × S grid on input

Bounding boxes + confidence

Class probability map

Final detections

# YOLO v1: Unified Detection



Bounding boxes + confidence

S × S grid on input

Class probability map

Final detections

**In addition, a grid cell also predicts C conditional class probabilities Pr(Class_i|Object)**

**Predictions can be encoded as S x S x (B*5 + C)**

# YOLO v1: Bounding Boxes and Confidence Scores

**Bounding Boxes:**

- Each bounding box gives 4 coordinates $x, y, w, h$
- $(x, y)$ = Coordinates representing center of the object **relative to the grid cell**
- $(w, h)$ = Width and height of the object **relative to the whole image**

# YOLO v1: Bounding Boxes and Confidence Scores

**Bounding Boxes:**

- Each bounding box gives 4 coordinates $x, y, w, h$
- $(x, y) =$ Coordinates representing center of the object **relative to the grid cell**
- $(w, h) =$ Width and height of the object **relative to the whole image**

**Confidence Score:**

- Reflects how confident the model is that the box contains an object and also how accurate it thinks the box is
- Formally,

$$\text{Confidence } = Pr(\text{Object}) * IOU_{pred}^{truth}$$

# YOLO v1: Conditional Class Probabilities

- Regardless of number of boxes $B$, we only predict one set of class probabilities per grid cell
- At test time, class-specific confidence scores for each box are given by:

$$Pr(Class_i|Object) * Pr(Object) * IOU_{pred}^{truth} = Pr(Class_i) * IOU_{pred}^{truth}$$

# YOLO v1: Loss Function

Loss function used to train YOLO v1 given by:

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

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

$$+\lambda_{\mathsf{noobj}} \sum_{i=0}^{S^2} \sum_{j=0}^{B} \mathbb{1}_{ij}^{\mathsf{noobj}} \left( C_i - \hat{C}_i \right)^2 + \sum_{i=0}^{S^2} \mathbb{1}_{i}^{\mathsf{obj}} \sum_{c \in \mathsf{classes}} \left( p_i(c) - \hat{p}_i(c) \right)^2$$
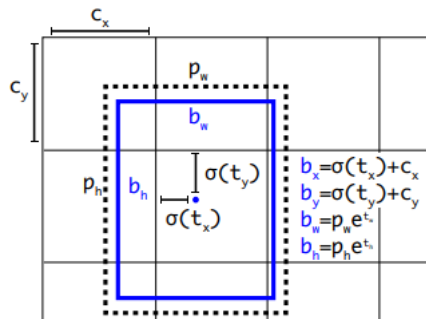
where:

- $\mathbb{1}_{i}^{\mathsf{obj}}$ denotes if object appears in cell $i$
- $\mathbb{1}_{ij}^{\mathsf{obj}}$ denotes that $j^{\mathsf{th}}$ bounding box predictor in cell $i$ is 'responsible' for that prediction

# YOLO v1: Limitations
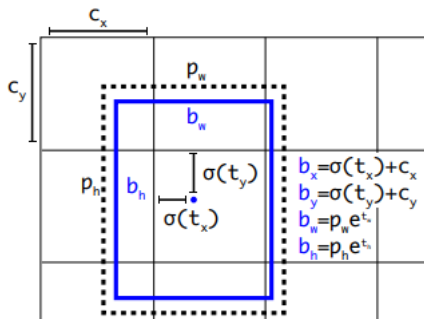
# YOLO v1: Limitations

- Detects only a small number of objects
- Cannot detect objects that are small or close due to strong spatial constraints
- High localization error
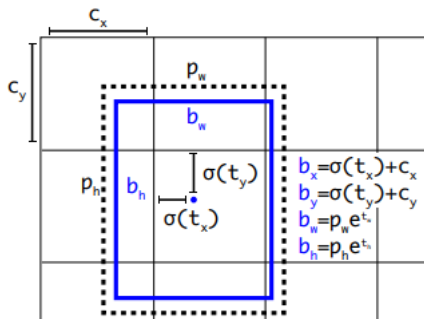- Relatively low recall

# YOLO v2: Anchor Boxes



$$b_x = \sigma(t_x) + c_x$$
$$b_y = \sigma(t_y) + c_y$$
$$b_w = p_w e^{t_w}$$
$$b_h = p_h e^{t_h}$$

# YOLO v2: Anchor Boxes



- Predicts 5 coordinates per anchor box: $t_x, t_y, t_h, t_w, t_o$

# YOLO v2: Anchor Boxes



- Predicts 5 coordinates per anchor box: $t_x, t_y, t_h, t_w, t_o$
- If cell is offset from top left corner of image by $(c_x, c_y)$ and bounding box has width and height $p_w, p_h$, then predictions correspond to:
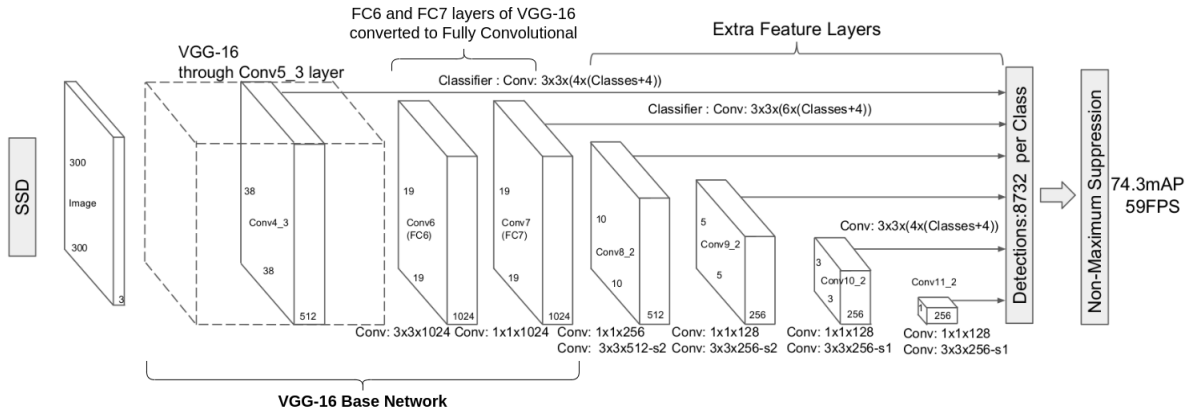
$$b_x = \sigma(t_x) + c_x$$
$$b_y = \sigma(t_y) + c_y$$
$$b_w = p_w e^{t_w}$$
$$b_h = p_h e^{t_h}$$
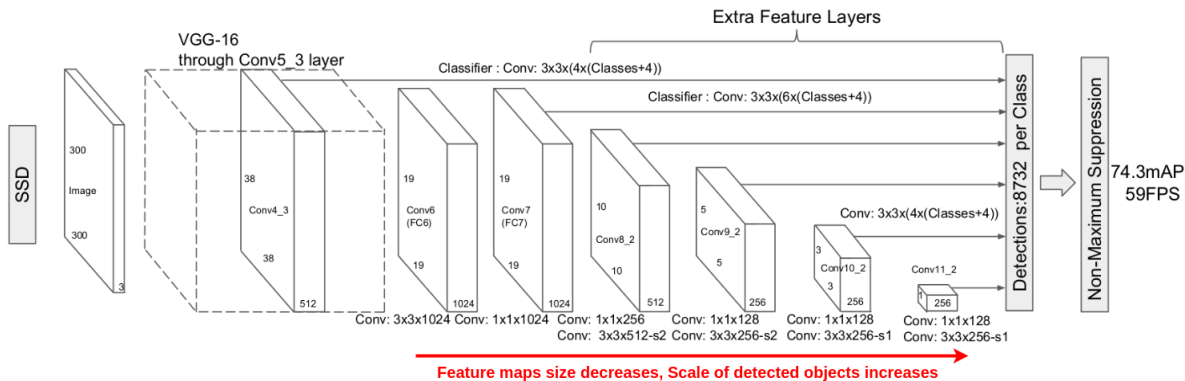$$Pr(\text{object}) * IOU(b, \text{object}) = \sigma(t_o)$$

# Single Shot MultiBox Detector (SSD)[2]

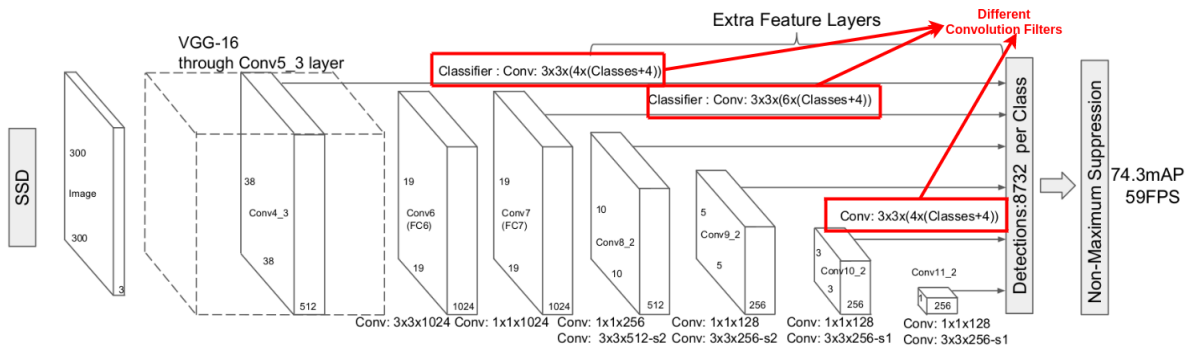[2]Liu et al, SSD: Single Shot MultiBox Detector, ECCV 2016

# SSD: Model

## Multi-scale Feature maps for Detection



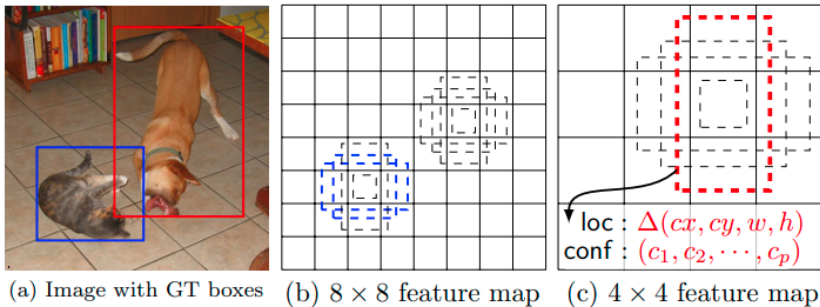Feature maps size decreases, Scale of detected objects increases

# SSD: Model

## Exclusive convolutional predictors for each feature map



A feature layer of size $m \times n$ with $c$ channels gives $m \times n$ locations (grid cells); bounding box offsets output values relative to grid cell location (like in Faster R-CNN)

# SSD: Anchor Boxes and Aspect Ratios



(a) Image with GT boxes    (b) $8 \times 8$ feature map    (c) $4 \times 4$ feature map

- For each of $k$ default boxes with different aspect ratios, SSD predicts $c$ class-specific scores and 4 anchor box offsets
- For an $m \times n$ feature map, there are $(c + 4)kmn$ outputs

# SSD: Loss Function

Weighted sum of localization loss (`loc`) and confidence loss (`conf`):

$$L(x, c, l, g) = \frac{1}{N}(L_{conf}(x, c) + \alpha L_{loc}(x, l, g))$$

where:

- $x_{ij}^p$ is {1,0} if $i^{\text{th}}$ default box matches $j^{\text{th}}$ ground truth box of category $p$
- $c =$ class probabilities
- $l, g =$ predicted and ground truth box parameters

# SSD: Loss Function

- **Localization loss** $L_{loc}$ given by:

$$L_{loc}(x, l, g) = \sum_{i \in Pos}^{N} \sum_{m \in \{cx, cy, w, h\}} x_{ij}^{k} \text{smooth}_{\text{L1}}(l_i^m - \hat{g}_j^m)$$

$$\hat{g}_j^{cx} = (g_j^{cx} - d_i^{cx})/d_i^w \qquad \hat{g}_j^{cy} = (g_j^{cy} - d_i^{cy})/d_i^h$$

$$\hat{g}_j^w = \log\left(\frac{g_j^w}{d_i^w}\right) \qquad \hat{g}_j^h = \log\left(\frac{g_j^h}{d_i^h}\right)$$

Regress offsets for center $(cx, cy)$ of anchor box $(d)$ and its width $(w)$ and height $(h)$

# SSD: Loss Function

- **Localization loss** $L_{loc}$ given by:

$$L_{loc}(x, l, g) = \sum_{i \in Pos}^{N} \sum_{m \in \{cx, cy, w, h\}} x_{ij}^k \text{smooth}_{L1}(l_i^m - \hat{g}_j^m)$$

$$\hat{g}_j^{cx} = (g_j^{cx} - d_i^{cx})/d_i^w \qquad \hat{g}_j^{cy} = (g_j^{cy} - d_i^{cy})/d_i^h$$

$$\hat{g}_j^w = \log\left(\frac{g_j^w}{d_i^w}\right) \qquad \hat{g}_j^h = \log\left(\frac{g_j^h}{d_i^h}\right)$$

Regress offsets for center $(cx, cy)$ of anchor box $(d)$ and its width $(w)$ and height $(h)$

- **Confidence loss** $L_{conf}$ is softmax loss of class confidences (probabilities):

$$L_{conf}(x, c) = -\sum_{i \in Pos}^{N} x_{ij}^p log(\hat{c}_i^p) - \sum_{i \in Neg} log(\hat{c}_i^0) \quad \text{where} \quad \hat{c}_i^p = \frac{\exp(c_i^p)}{\sum_p \exp(c_i^p)}$$

# SSD: Practical Implementation

**Hard Negative Mining:**

- Similar to Faster R-CNN, most anchor boxes are negative
- To counter it, select negative examples that have highest confidence loss such that ratio between negatives and positives is $\sim 3 : 1$

# SSD: Practical Implementation

**Hard Negative Mining:**

- Similar to Faster R-CNN, most anchor boxes are negative
- To counter it, select negative examples that have highest confidence loss such that ratio between negatives and positives is $\sim 3 : 1$

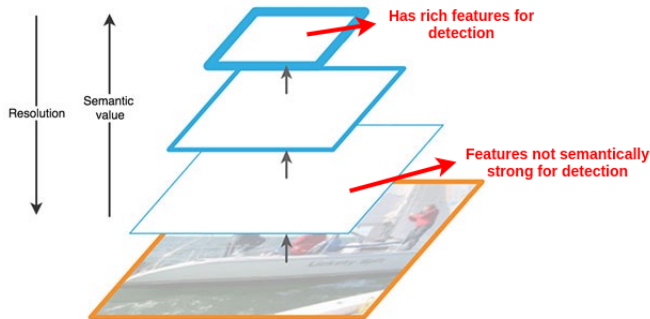**Data Augmentation:** Training samples are obtained as below:

- Use original image
- Randomly sample a patch, such that minimum IoU with objects is in $\{0.1, 0.3, 0.5, 0.7, 0.9\}$

# SSD: Performance

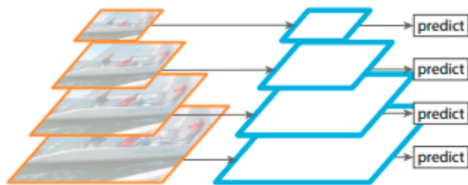| Method | mAP | FPS | batch size | # Boxes | Input resolution |
|---|---|---|---|---|---|
| Faster R-CNN (VGG16) | 73.2 | 7 | 1 | $\sim 6000$ | $\sim 1000 \times 600$ |
| Fast YOLO | 52.7 | 155 | 1 | 98 | $448 \times 448$ |
| YOLO (VGG16) | 66.4 | 21 | 1 | 98 | $448 \times 448$ |
| SSD300 | 74.3 | 46 | 1 | 8732 | $300 \times 300$ |
| SSD512 | 76.8 | 19 | 1 | 24564 | $512 \times 512$ |
| SSD300 | 74.3 | 59 | 8 | 8732 | $300 \times 300$ |
| SSD512 | 76.8 | 22 | 8 | 24564 | $512 \times 512$ |

# Feature Pyramid Network (FPN)[3]

- Feature maps from initial layers (which are high resolution) cannot be used for detection
- FPN provides a top-down pathway to construct higher resolution layers from a semantically rich layer
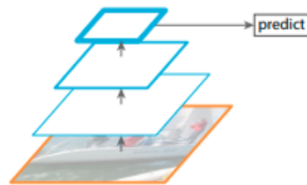


---

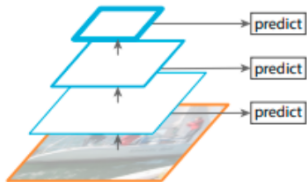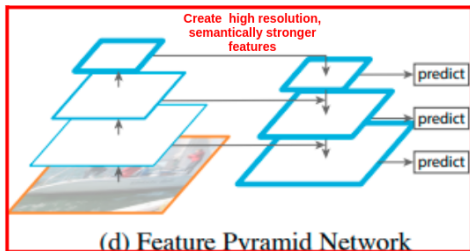[3]Lin et al, Feature Pyramid Networks for Object Detection, CVPR 2017

# Feature Pyramid Network
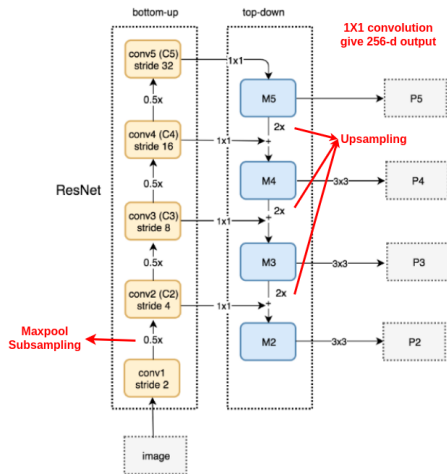


(a) Featurized image pyramid

(b) Single feature map

(c) Pyramidal feature hierarchy

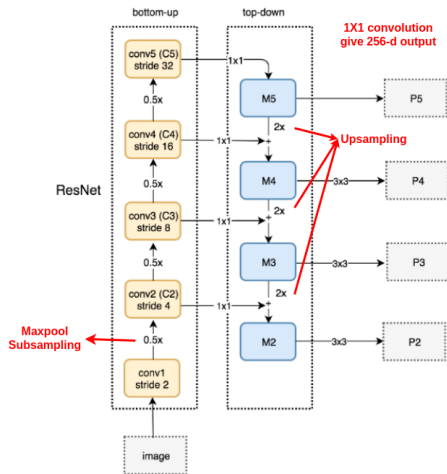Create high resolution, semantically stronger features

(d) Feature Pyramid Network

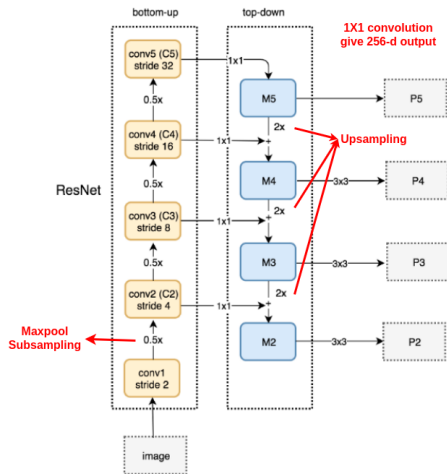# Feature Pyramid Network: Methodology



- All convolutional feature maps are treated with a $1 \times 1$ convolution with 256 channels

# Feature Pyramid Network: Methodology



- All convolutional feature maps are treated with a $1 \times 1$ convolution with 256 channels
- M5 is upsampled by a factor of 2 (Maxpool at $1/2 \times 1/2$)

# Feature Pyramid Network: Methodology



- All convolutional feature maps are treated with a $1 \times 1$ convolution with 256 channels
- M5 is upsampled by a factor of 2 (Maxpool at $1/2 \times 1/2$)
- M5 and C4 signals are element-wise added to give M4; similarly followed in the downward direction
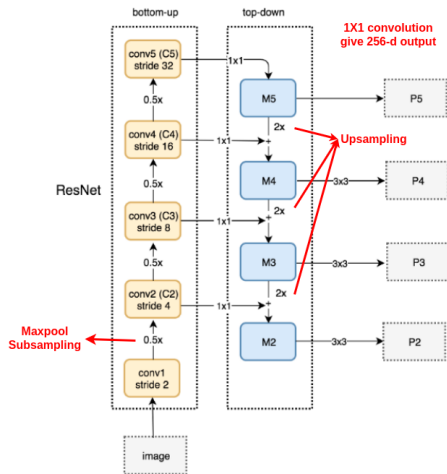
# Feature Pyramid Network: Methodology



- All convolutional feature maps are treated with a $1 \times 1$ convolution with 256 channels

- M5 is upsampled by a factor of 2 (Maxpool at $1/2 \times 1/2$)

- M5 and C4 signals are element-wise added to give M4; similarly followed in the downward direction

- $3 \times 3$ convolution is used to reduce aliasing effect of $M$s

# Feature Pyramid Network: Methodology
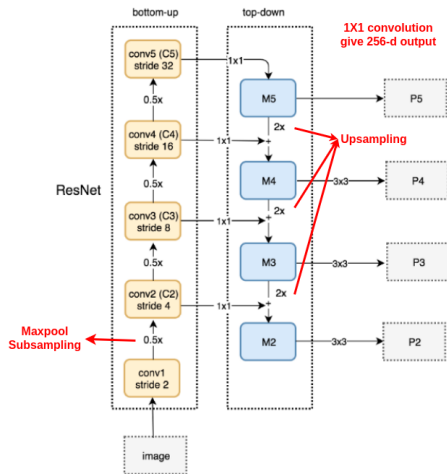


*Credit: Jonathan Hui, Medium.com*

- All convolutional feature maps are treated with a $1 \times 1$ convolution with 256 channels
- M5 is upsampled by a factor of 2 (Maxpool at $1/2 \times 1/2$)
- M5 and C4 signals are element-wise added to give M4; similarly followed in the downward direction
- $3 \times 3$ convolution is used to reduce aliasing effect of $M$s
- Finally, $P$s are individually fed into exclusive object detectors

# RetinaNet[4]

**Intuition:** Two-stage detectors are more accurate than one-stage detectors due to lesser class imbalance between background (negative) and object-containing (positive) proposals



**Two-stage Detectors**

Selective Search/Region Proposal Network → 1-2k regions

Lesser background to objects ratio

Very high background to objects ratio

Dense sampling based Sliding Window → 100k regions

**One-stage Detectors**

[4]Lin et al, Focal Loss for Dense Object Detection, ICCV 2017

# Class Imbalance in Object Detection: The Problems



Many negative examples, no useful signal

Few positive examples, rich information

# Class Imbalance in Object Detection: The Problems



Many negative examples, no useful signal

Few positive examples, rich information

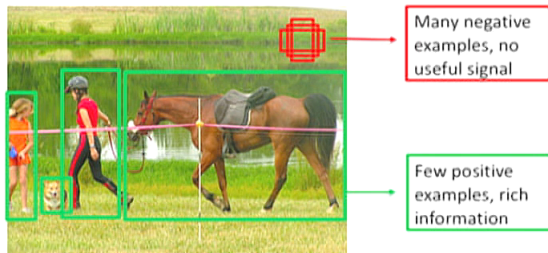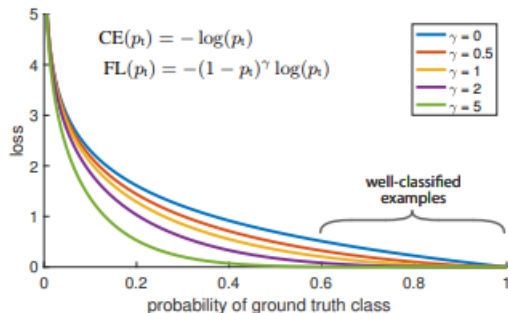- Training is inefficient as easy negatives contribute no useful signal

# Class Imbalance in Object Detection: The Problems



- Training is inefficient as easy negatives contribute no useful signal
- Loss due to easy negatives overwhelms loss due to positives and thereby, training process can lead to degenerate models; hard negative training alleviates it to some extent

*Credit: Sik-Ho Tsang, TowardsDataScience.com*

# CE Loss is Bad[5]



$$\text{CE}(p_t) = -\log(p_t)$$
$$\text{FL}(p_t) = -(1-p_t)^\gamma \log(p_t)$$

- $\gamma = 0$
- $\gamma = 0.5$
- $\gamma = 1$
- $\gamma = 2$
- $\gamma = 5$

well-classified examples

loss

probability of ground truth class

- CE Loss for binary classification is typically implemented as:

$$\text{CE}(p, y) = \begin{cases} -\log(p) & \text{if } y = 1 \\ -\log(1-p) & \text{otherwise.} \end{cases}$$

- Can be rewritten as
$$CE(p, y) = CE(p) = -\log(p_t)$$

---

[5]Lin et al, Focal Loss for Dense Object Detection, ICCV 2017

# CE Loss is Bad[5]



$$\text{CE}(p_t) = -\log(p_t)$$

$$\text{FL}(p_t) = -(1 - p_t)^\gamma \log(p_t)$$

legend:
- $\gamma = 0$
- $\gamma = 0.5$
- $\gamma = 1$
- $\gamma = 2$
- $\gamma = 5$

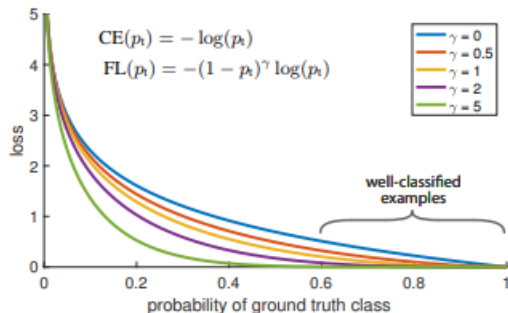well-classified examples

probability of ground truth class

- CE Loss for binary classification is typically implemented as:

$$\text{CE}(p, y) = \begin{cases} -\log(p) & \text{if } y = 1 \\ -\log(1 - p) & \text{otherwise.} \end{cases}$$

- Can be rewritten as $CE(p, y) = CE(p) = -\log(p_t)$

- Can be observed in graph ($\gamma = 0$ curve) that easily classifiable examples ($p >> .5$) incur loss of non-trivial magnitude
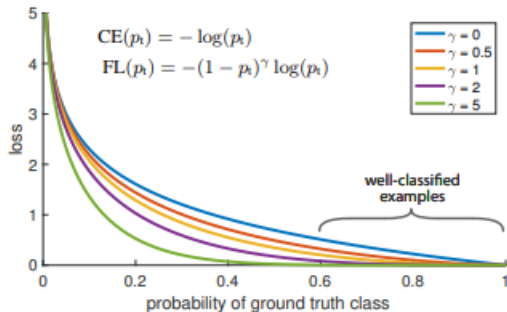
---

[5]Lin et al, Focal Loss for Dense Object Detection, ICCV 2017

# RetinaNet: Balanced Cross Entropy and Focal Loss



$$CE(p_t) = -\log(p_t)$$
$$FL(p_t) = -(1 - p_t)^\gamma \log(p_t)$$

- $\gamma = 0$
- $\gamma = 0.5$
- $\gamma = 1$
- $\gamma = 2$
- $\gamma = 5$

well-classified examples

loss

probability of ground truth class

**Balanced Cross Entropy:**

- Introduces a weighting factor $\alpha \in [0, 1]$ which can be inverse class frequency or a hyperparameter
- $\alpha$-balanced CE loss is given by:

$$CE(p_t) = -\alpha \log(p_t)$$

# RetinaNet: Balanced Cross Entropy and Focal Loss



$$\mathbf{CE}(p_t) = -\log(p_t)$$
$$\mathbf{FL}(p_t) = -(1 - p_t)^\gamma \log(p_t)$$

- $\gamma = 0$
- $\gamma = 0.5$
- $\gamma = 1$
- $\gamma = 2$
- $\gamma = 5$

well-classified examples

loss

probability of ground truth class

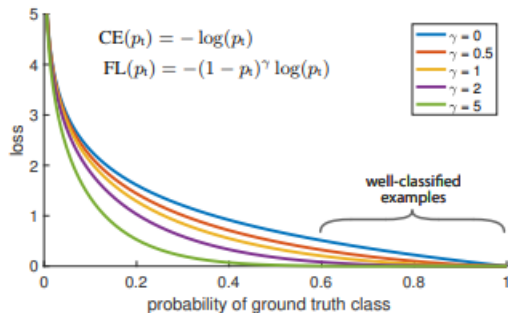**Balanced Cross Entropy:**

- Introduces a weighting factor $\alpha \in [0, 1]$ which can be inverse class frequency or a hyperparameter

- $\alpha$-balanced CE loss is given by:
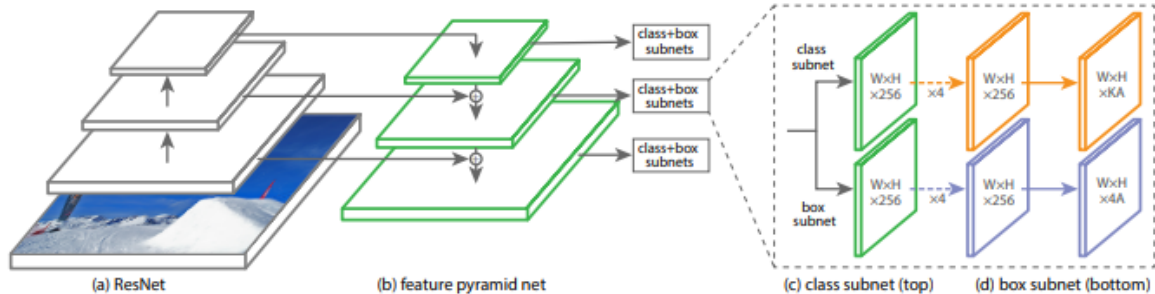
$$CE(p_t) = -\alpha \log(p_t)$$

**Focal Loss:**

- Adds a modulating factor $(1 - p_t)^\gamma$ to CE loss, with tunable focusing parameter $\gamma \geq 0$.

- Focal loss is given by:

$$FL(p_t) = -(1 - p_t)^\gamma \log(p_t)$$

**FPN + Focal Loss**



(a) ResNet  (b) feature pyramid net  (c) class subnet (top)  (d) box subnet (bottom)

# Detectron

- Framework developed by Facebook AI Research (FAIR) to implement state-of-the-art object detection algorithms

- Highly flexible providing support across various algorithms and backbone networks

- See Detectron and Detectron2 for details

# Homework

## Readings

- Object Detection for Dummies, Part 4

- YOLO Family: All you want to know

- Understanding SSD

- Understanding FPN

- Understanding RetinaNet

# Homework

## Exercises

- YOLO9000 and YOLOv3 were follow-ups of YOLOv2. What was different in these extensions? Find out! (*Hint:* see the YOLO Family: All you want to know link)

- Given two bounding boxes in an image: an upper-left box which is $2 \times 2$, and a lower-right box which is $2 \times 3$ and an overlapping region of $1 \times 1$, what is the IoU between the two boxes?

- Consider using YOLO object detector on a $19 \times 19$ grid, on a detection problem with 20 classes, and with 5 anchor boxes. During training, for each image, you will need to construct an output volume $y$ as the target value for the neural network; this corresponds to the last layer of the neural network. ($y$ may include background). What is the dimension of this output volume?

# References

Wei Liu et al. "Ssd: Single shot multibox detector". In: *European conference on computer vision*. Springer. 2016, pp. 21–37.

Joseph Redmon et al. "You only look once: Unified, real-time object detection". In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016, pp. 779–788.

Tsung-Yi Lin et al. "Feature pyramid networks for object detection". In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2017, pp. 2117–2125.

Tsung-Yi Lin et al. "Focal loss for dense object detection". In: *Proceedings of the IEEE international conference on computer vision*. 2017, pp. 2980–2988.

Joseph Redmon and Ali Farhadi. "YOLO9000: better, faster, stronger". In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2017, pp. 7263–7271.