

Object localization

- Loss:

$$L = L_{\text{reg}} + L_{\text{cls}}$$

↗ ↘
 regression classification
 e.g. cross entropy e.g. L_2
 ↓
 $L_i(y^{(i)}, \hat{y}^{(i)}) = \begin{cases} (y_i^{(i)} - \hat{y}_i^{(i)})^2 & \text{if } y_i^{(i)} \neq 0 \\ \sum_j (y_j^{(i)} - \hat{y}_j^{(i)})^2 & \text{otherwise} \end{cases}$
 ↓
 w.r.t cross entropy for class labels

1st element in vector $y \rightarrow$ Pc

- Suitable for one object in an image. What should we do if there are multiple objects?

Sliding window detection

- Slide a window through the image. At each location:
 - Crop sub window
 - Classify each sub window (object/background)
- To deal with different object scales, crop different sub windows with different sizes and aspect ratios at each location
- This may be slow because there are many sub windows to classify



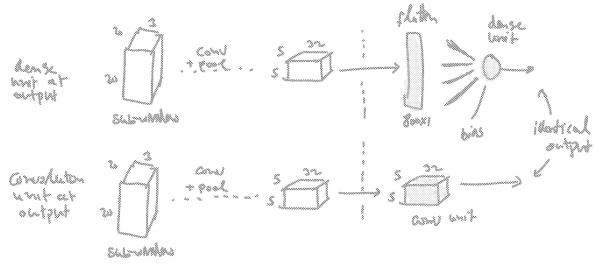
(1000 × 1000 image → 10,000 sub-windows
 with 10 sub-windows to classify for each
 at each location image

Sliding window detection

- Problems:
 - Running dense layers for classifying each sub window is inefficient
 - In addition, dense layers need a fixed input size and so the image input size needs to be fixed
- Fully convolutional implementation:
 - A fully convolutional implementation of the dense layers is more efficient: no need to process each sub window separately, and convolutions can be implemented efficiently (e.g. as in mobileNet)
 - A fully convolutional implementation can work on variable size images

Sliding window detection

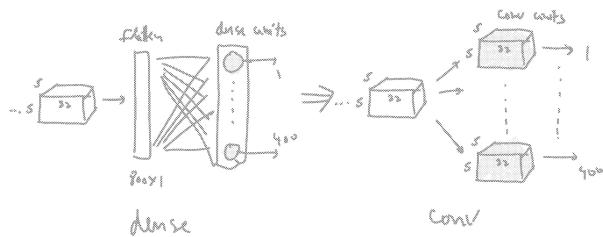
- A single fully connected (Dense) unit may be replaced by a single convolutional unit that has the exact dimensions of the input tensor (without flattening it)



Sliding window detection

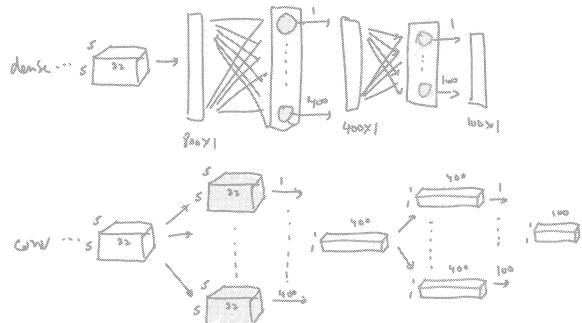
- Consider replacing 400 fully connected (FC) units

- The output of 400 FC (dense) units is a 400×1 vector
- We can replace 400 FC units with $5 \times 5 \times 32$ convolution filters and so their output will be a 400×1 vector



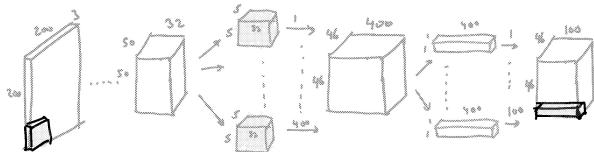
Sliding window detection

- Replacing a subsequent dense layer with 100 units:



Sliding window detection

- So far we performed a single sub window classification
- For a single sub window there is no reduction in computations
- The advantage of convolutions is when sharing computations between overlapping windows



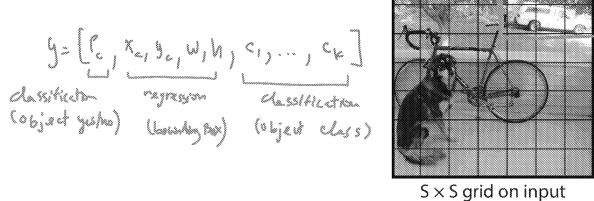
- Each column in the output provides probability for 100 classes in a corresponding image region (sub window)

Sliding window detection

- Summary:
 - Use convolution + pooling to extract features and reduce spatial resolution
 - Each location in a reduced resolution layer corresponds to a larger patch in the original image (patches overlap)
 - Use a dense layer or 1×1 convolution to classify or regress at each location in reduced resolution layer (classify class and regress bounding box)
- Problem:
 - Detection using a sliding window does not produce accurate detection because:
 1. It is based on classification in a pooled (lower resolution) layer
 2. The region covered by filters is of fixed size and so we detect fixed size bounding boxes
 - To solve this we can, in addition to classifying each location, regress bounding boxes

Grid cell detection

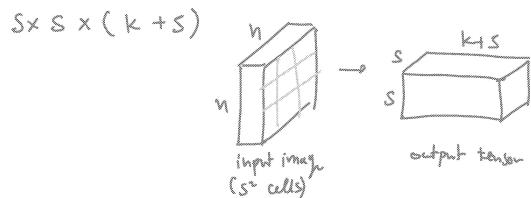
- Instead of scanning the image, divide the image into grid cells and attempt to detect an object in grid cell
- Cells do not necessarily match objects and so use regression to find a bounding box for each detected object
- Similar to object localization (single object) except that now at each grid cell detect and localize objects
- At each grid cell train and predict using:



Refer image

Grid cell detection

- For an $S \times S$ grid and K classes the output is a tensor of size:



- At each grid cell with an object find a precise bounding box instead of using a fixed size window
- Only one object is detected at each grid cell

Refer image

Grid cell detection

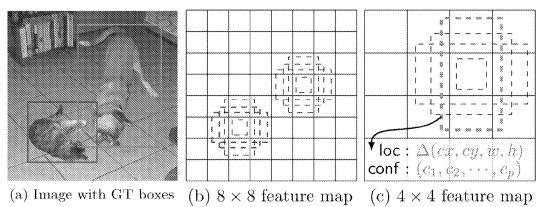
- Problems:

- What should be the grid cell size?
 - A too small cell will fail in classifying big objects
 - A too large cell will fail in classifying small objects
- What should be the aspect ratio of candidate cells?
 - Different objects have different aspect ratios of boxes that should be used for their detection
- How to detect more than one object in each grid cell?

Grid cell detection

- Anchor boxes:

- To solve three problems we use multiple anchor boxes of each grid location
- Anchor boxes have different sizes and different aspect ratios (selected according to what we see in the training data)



Grid cell detection

- Multiple anchor boxes:
 - To allow detecting two different (overlapping) objects at each location (e.g. objects with different aspect ratios) double the target factor y
- Most grid cells will have zero or one anchor boxes

$$y = [l_c^1, x_c^1, y_c^1, w_c^1, h_c^1; c_1^1, \dots, c_k^1; l_c^2, x_c^2, y_c^2, w_c^2, h_c^2; c_1^2, \dots, c_k^2]$$

$p_c^i = 0$ and $p_c^j = 0$ if object does not exist

- Anchor box assignment during training:
 - During training, the anchor box assigned to an object is the one having the highest IOU with a ground truth box



Refer image

Single shot detectors

- Single shot versus double shot methods:
 - Single shot detectors perform object detection in a single network that looks at the image once
 - Two shot detectors look at the image twice: one to propose object regions (region proposals) and a second to refine and classify regions
 - Two shot methods have a higher computational cost (lower frame rate) but may be more accurate
- Single shot methods:
 - Consider predefined anchor boxes at grid cells
 - Refine and classify region proposals
- Two shot methods:
 - Consider predefined anchor boxes at grid cells
 - Classify anchor boxes as objects or not to get region proposals
 - Refine and classify region proposals (higher quality classifier)
- Single shot approaches:
 - YOLO (You Only Look Once)
 - MultiBox detection
 - Single Shot MultiBox Detector (SSD)

YOLO (You Only Look Once)

- A fast single-shot detector
- Algorithm:
 - Divide image into ($s \times s$) grid cells (e.g. $s = 19$). Each grid cell is responsible for detecting objects with centers in it.
 - Each grid cell predicts B boundary boxes and confidence scores (confidence that the box contains the object and is accurate).
- Confidence:
 - Confidence = Probability(object) + IOU(truth, prediction)
 - Confidence = 0 if no object in cell.
- Training:
 - Use a vector of B bounding boxes per cell.
- Inference:
 - Predict using dense layers or 1x1 convolutions.
 - Apply non-maximum suppression (NMS).

$$y = [l_c^1, b_x^1, b_y^1, w_c^1, h_c^1; c_1^1, \dots, c_k^1; l_c^2, b_x^2, b_y^2, w_c^2, h_c^2; c_1^2, \dots, c_k^2]$$

$p_c^j = 0$ if boundary box does not contain object

YOLO (You Only Look Once)

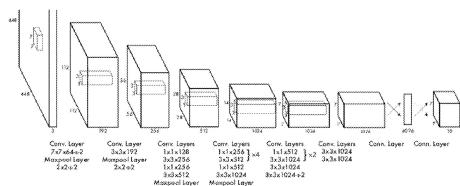


Figure 3: The Architecture. Our detection network has 24 convolutional layers followed by 2 fully connected layers. Alternating 1×1 convolutional layers reduce the features space from preceding layers. We pretrain the convolutional layers on the ImageNet classification task at half the resolution (224×224 input image) and then double the resolution for detection.

- Use Leaky ReLU activations inside the network

YOLO (You Only Look Once)

Yolo loss function:

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

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

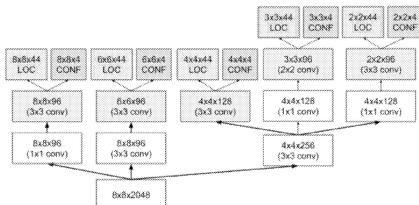
$$+ \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \underbrace{(C_i - \hat{C}_i)^2}_{\text{object classification}}$$

$$+ \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{noobj}} \underbrace{(C_i - \hat{C}_i)^2}_{\text{no-object classification}}$$

$$+ \lambda_{\text{noobj}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{noobj}} \underbrace{(p_i(c) - \hat{p}_i(c))^2}_{\text{box classification}}$$

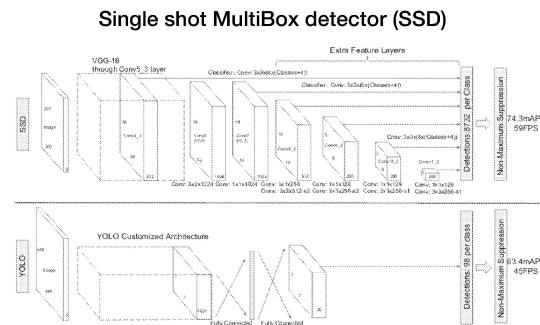
$$+ \sum_{i=0}^{S^2} \mathbb{1}_i^{\text{obj}} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2$$

MultiBox detection

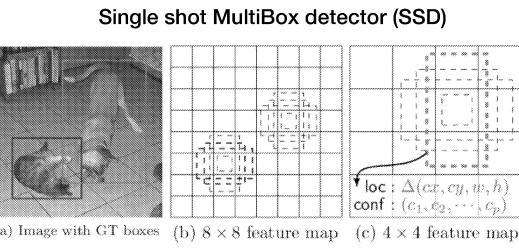


Architecture of multi-scale convolutional prediction of the location and confidences of multibox

- Use prior for anchor boxes determined per class
 - Inception-like detection at different scales
 - Loss
 - Confidence_loss = class categorical cross entropy
 - Location loss = bounding box IoU
 - MultiBox loss = Confidence loss + alpha * Location loss



- SSD = YOLO + MultiBox
- Use feature maps from intermediate convolution layers



- Use different priors for different convolution layers
- Apply non-maximum suppression (NMS)
- Use hard negatives (FP predicted by the model) to train instead of all negatives. This helps in balancing negative and positive examples.

• Loss:

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

$$L_{conf} = \frac{1}{n_{positives}} \left(\sum_{positives} \text{cross-entropy loss} + \sum_{hard\ negatives} \text{cross-entropy loss} \right)$$