

Standard CV probs: object classification, object localization (bounding box), object detection (every object), semantic segmentation (scene understanding)

Object localization: $\mathcal{D} = \{(x_i, y_i, w_i, h_i)\}$ (bboxes)

not a loss func! →

Measuring localization accuracy → matter of convention

Typical: IoU (intersection over union). Correct if $IoU > 0.5$ (threshold).

O.L as regression: multi-task model, cross-entropy loss for class, (c_i) (two heads) regression loss for x_i, y_i, w_i, h_i

OR train classification THEN regression head

O.L using sliding window: classify every patch in image to form bounding box.
→ consider scaling w/ different aspect ratios, matching aspect ratio of object.

OverFeat! For every patch, output class prob AND bbox coords.

"Average" together boxes to get answer.

Save compute for sliding window thru "convolutional classification", Reuse calculations across windows.

The intuition here is that everything is rolled down to convolutions.

The FC layers → 1x1 convolutions!

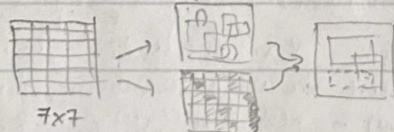
n_i objects, differs for each x_i .

Object Detection → $(c_{x_i}, c_{y_i}, x_{min}, y_{min}, w_i, h_i, \dots, c_{x_{n_i}}, c_{y_{n_i}}, x_{min}, y_{min}, \dots)$

Non-maximal suppression: "kill" off any detections that have other high-ranking detections of the same class nearby.

2 Ideas for Model: - output best class at each window > YOLO! (You only look 49 times!)
- train on fixed # of boxes

YOLO:



For each cell: $x, y, w, h \rightarrow c \rightarrow Iou \rightarrow B$ of these implements "conv. classification"

essentially an overfitted OverFeat.

If too many objects, skip them. R-CNN? A "smarter" sliding window.

Fast R-CNN: First forward image through conv.

Then process ROIs thru maxpool + FC.

Now: ONE convolution. Can be trained end-to-end.

Use same model to train ROIs using bbox coords to be optimized.

Semantic segmentation: label every single pixel with its class.

Main issue: expensive, head-on.

Idea: High-rec → low-rec → high-rec, much cheaper. "bottle-neck".

convolutions

Upsampling / Transpose convolution: incorporate fractional "stride" (ex. 1/2).

"Unpooling": remember the positions of max elements and place them back.

Problem: we may lose spatial information.

V-Net: implements skip connections.

concatenate down w/ up (initial) (processed)

theoretically contains spatial info

down-sample

up-sample

