(R-CNN, Fast RCNN, Fast RCNN)

Two shot

(Yolo & SSD)

Single shot

✳ R-CNN

i/p image

↓ selective search

extract region proposals (~2k)

↓ (warp warped image regions)

compute CNN feature

↓ (forward each region to CNN)

classify regions.

(classify regions with svm)

✳ How it uses (linear Regr of bbox)

→ multi task loss

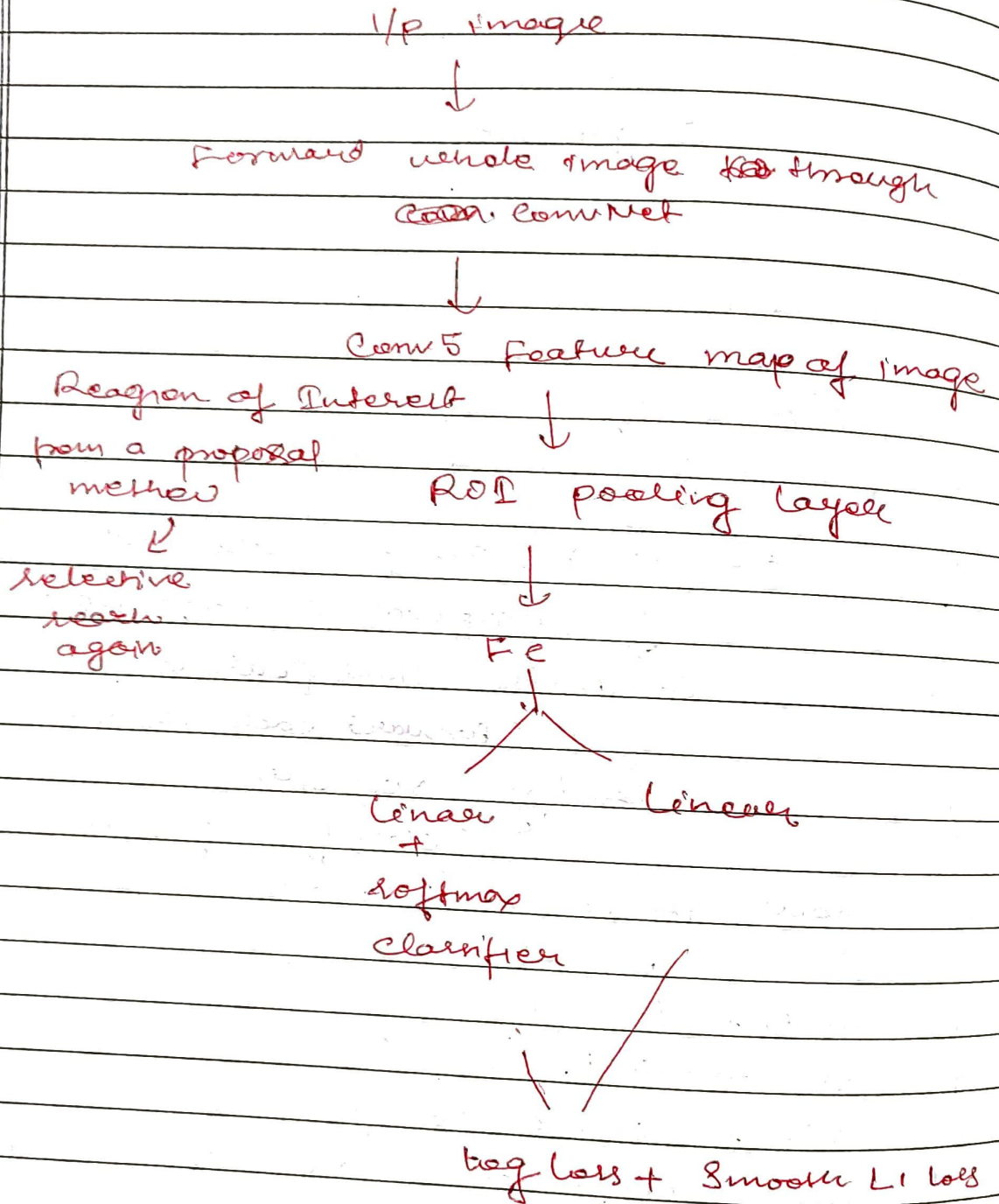→ Gradient descent

→ wrt fn to the loss fn

→ other performance matric as hyper param is tricky to choose
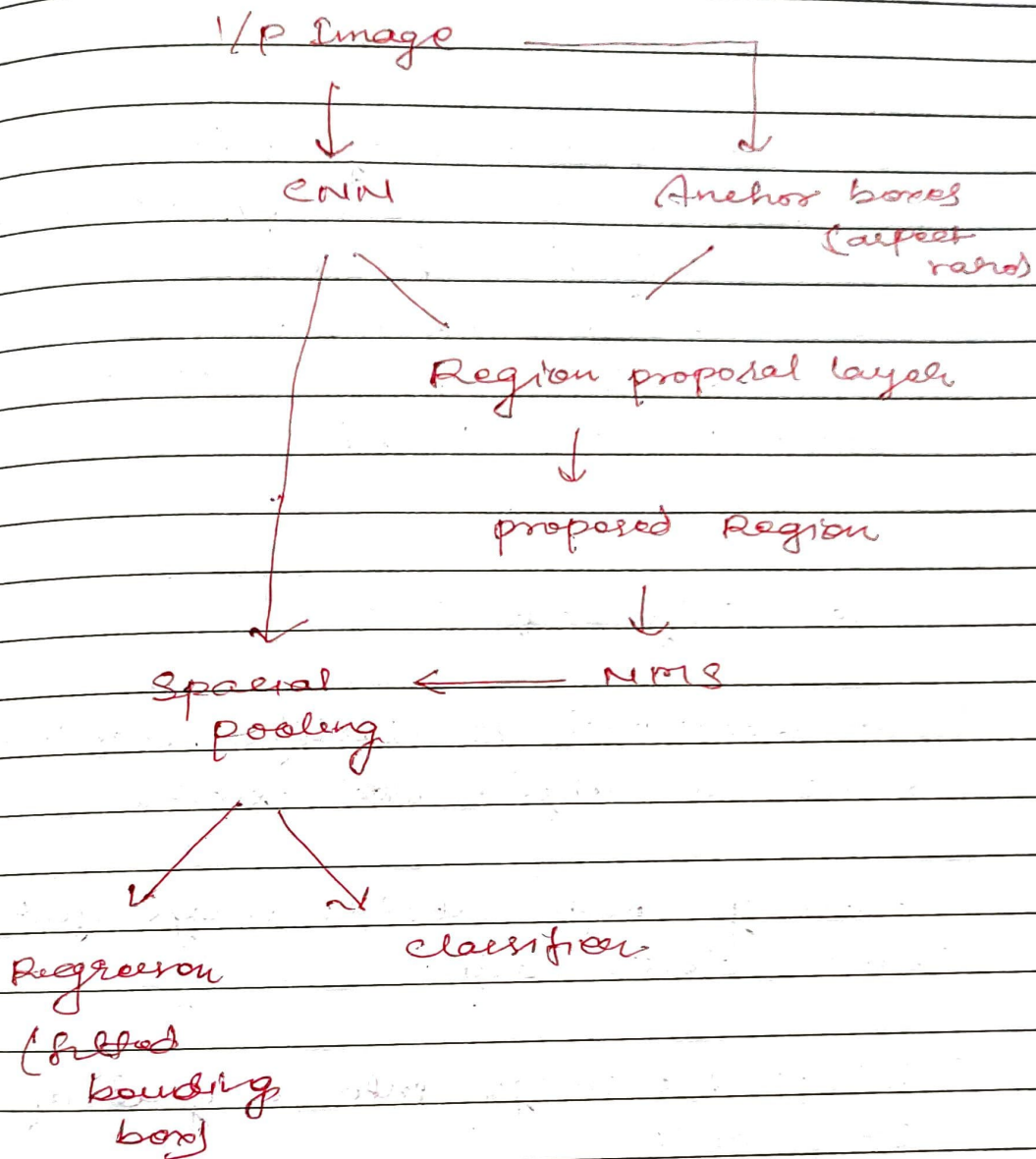
→ weighted sum of two losses.

→ Disadvantage

Separate CNN for each box (slow)

Fast R-CNN

I/P image

↓

Forward whole image through Conv. ConvNet

↓

Conv 5 Feature map of image

Region of Interest
from a proposal          ↓
method
                    ROI pooling layer
↓
Selective
search                      ↓
again
                            Fc
                           / \
            Linear              Linear
              +
           Softmax
           Classifier          /
                              / 
                             /
                    Reg loss + Smooth L1 loss

# Faster R-CNN

I/P Image

↓                              ↓

CNN                        Anchor boxes
                              (aspect
                                ratio)

Region proposal layer

↓

proposed Region

↓

Spacial   ←——   NMS
pooling

↓        ↓

Regresson        classifier
(shifted
bouding
box)

→ micro Region proposal N/w to predict
from features

Jointly train with 4 loses

→ RPN ⟳ classify object / not object
→ " regress box co-ordinatn               4
→ Final classification score              3
→ Final box co-ordin

Detection without Proposals

Yolo / SSD

→ Yolo

It divides the image into $S \times S$ grid and for each grid cell predict B bounding boxes, confidence for those boxes and C class probabilities.
These predictions are encoded as an $S \times S \times (B*5 + C)$ tensor.

→ You Only Look Once

→ Not traditional a classifier that is repurposed to be an object detector

→ Actually looks at the image just once but in clever way

→ Divide the image into a grid of say $13 \times 19$ cells

→ Each of these cells is responsible for predicting 5 bounding boxes

→ A bounding box describes the rectangle that encloses an object.

→ YOLO outputs a confidence score that tells how good the shape of the box is.

→ For each bounding box, the cell also predicts a class.

→ Yolo was trained on PASCAL VOC dataset of 20 different classes

→ The confidence score of bounding box and class prediction are combined into final score.
probability that his bounding box contains a specific object.

*
13×13 = 169 grid cell
169×5 = 845 bounding boxes
most of them have low confidence score
Threshold of 30% or more ⇒3
i/p image 416×416 resized
13×13×125 tensor describing the bounding boxes for grid cells

⭒

y

g

# Yolo Bounding box:

→ The i/p image is divided into $S \times S$ grid $(S = 7)$. If the center of an object falls into a grid cell, that grid cell is responsible for detecting that object.

→ Each grid cell predict B bounding boxes $(B = 2)$ and confidence scores for those boxes.

These confidence scores reflect how confident the model is that box contains an object i.e.

any objects in the box, $P(object)$

→ Each bounding box consists of 5 predictions $x, y, w, h,$ and confidence.

→ The $(x, y)$ coordinates represent the center of the box relative to the bounds of the grid cell.

→ The width w and height h are predicted relative to the bounds of the grid cell

→ The width w & height h are predicted relative to the whole image

→ The confidence represents the IOU b/w the predicted box and the any ground truth box.

→ By using SSD, we only need to take one single ~~sho~~ shot to detect multiple objects within the image.

→ Regional proposal network (RPN) based approaches such as R-CNN, Fast R-CNN series needs two shots, one for generating region proposals, one for detecting the object of each proposal

→ SSD is much faster

✳ ~~2 loss fundn~~ loss function has 2 ter. $L_{conf}$, $L_{loc}$

✳ A feature layer of size $m \times n$ (# of locations) with p channels

✳ for each location, we got k bounding boxes

✳ for each of the bounding box, we will compute c class scores and 4 offset relative to the original default bounding box shape.

✳ Thus, we got $(c+4) \, kmn \; o/p$