

April 5, 2023

DSML: Computer Vision.

Object localization and Detection - 1

Class starts
@ 9:05 pm.



What normal people see
when they walk on street



What Computer Vision
folks see



WHO WOULD WIN?



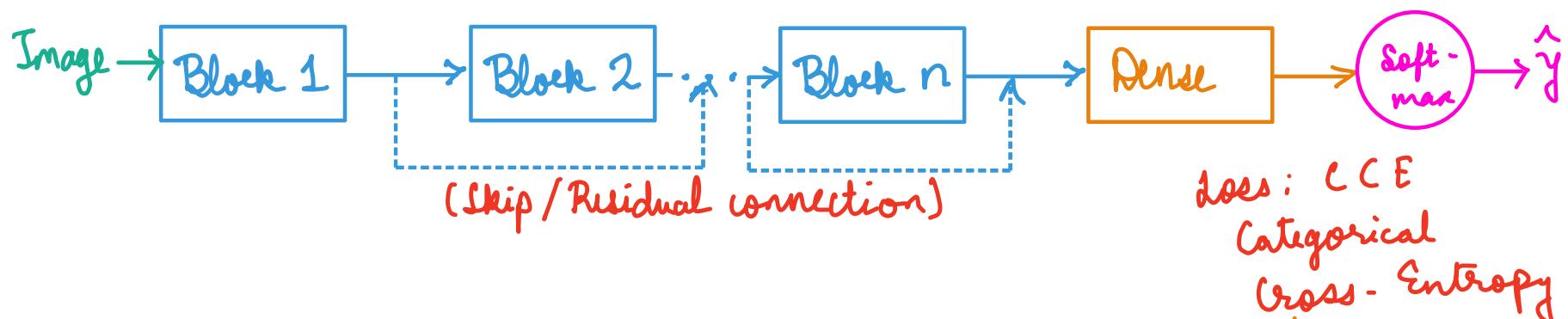
STATE OF THE ART
NEURAL NETWORK



ONE NOISY BOI

Recap:

- * Basic CNN architecture for Image classification:



- * Popular architectures and their performance:

Model Name	Number of params	Top 1 Acc	Top 5 Acc
EfficientnetB0	5.3M	77.3	93.5
MobileNet	2.3M	71.0	90.5
ResNet50	25.6M	83.2	96.5
Inception	22.9M	79.0	94.5
VGG16	138M	74.4	91.9
AlexNet	62M	63.3	84.6

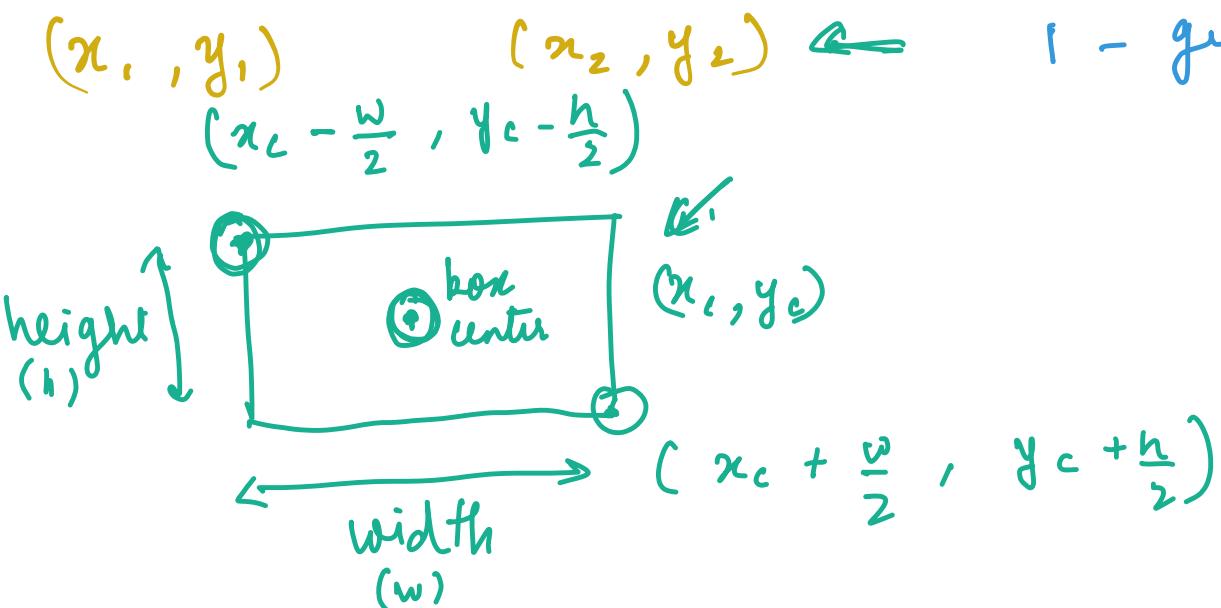
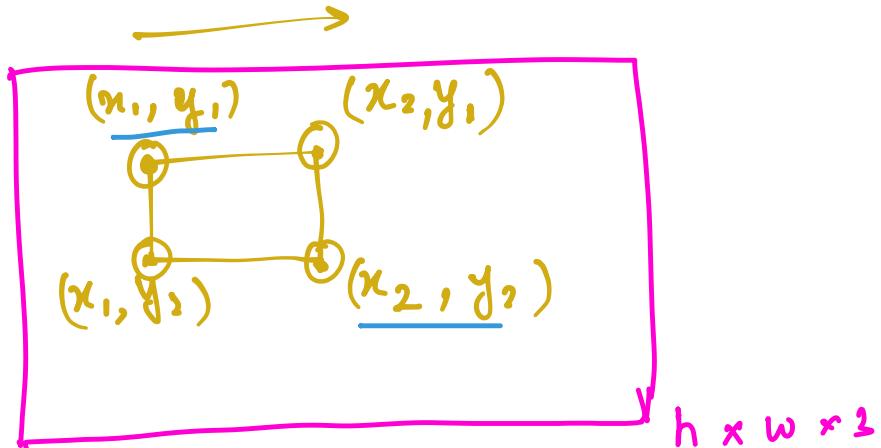
Novelty
 Automated scaling strategy.
 → Factorizing convolutions
 Residual / skip connections.
 Inception module, 1×1 conv
 3×3 kernels, blocks.
 First method to use GPUs,
 SGD.

Recap:

* Major CNN related concepts:

- (a) Transfer learning: Reuse an already trained CNN for image classification.
 - Helpful when we have a small dataset.
 - As long as images are similar to Imagenet, better than training from scratch.
- (b) Image similarity: Reuse of an already trained CNN for unsupervised similarity search.
 - Select a convolutional representation. This is referred to as "activations" or "embeddings"
 - Store these image representations in a database optimized for KNN.
 - Similar images have "activations" which are close in the Euclidean sense.

Discussion : Bounding boxes



for our dataset:

Annotations: contain 5 numbers.

—
—
—
—
—
↓ $\underbrace{2}_{x_c, y_c}$ $\underbrace{3}_{w, h}$ $\underbrace{4}_{x_c, y_c}$ $\underbrace{5}_{w, h}$

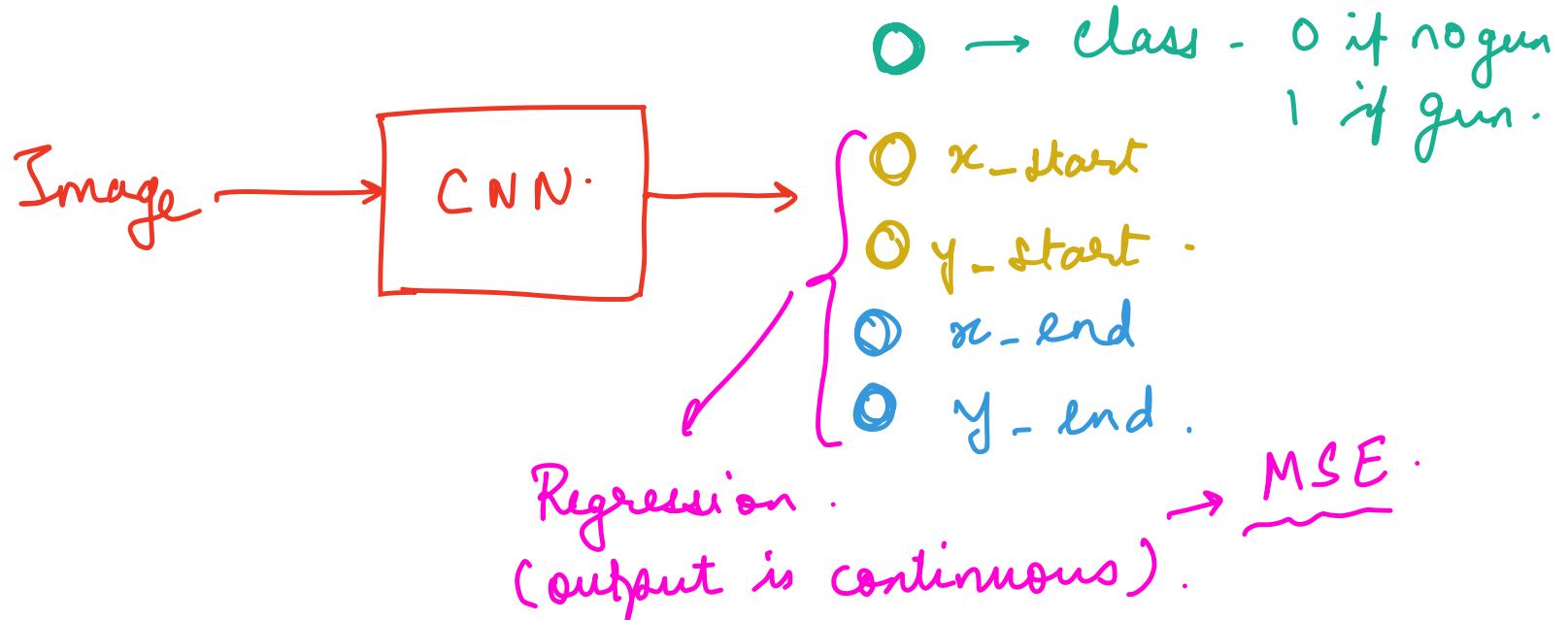
0 - no gun

or

1 - gun

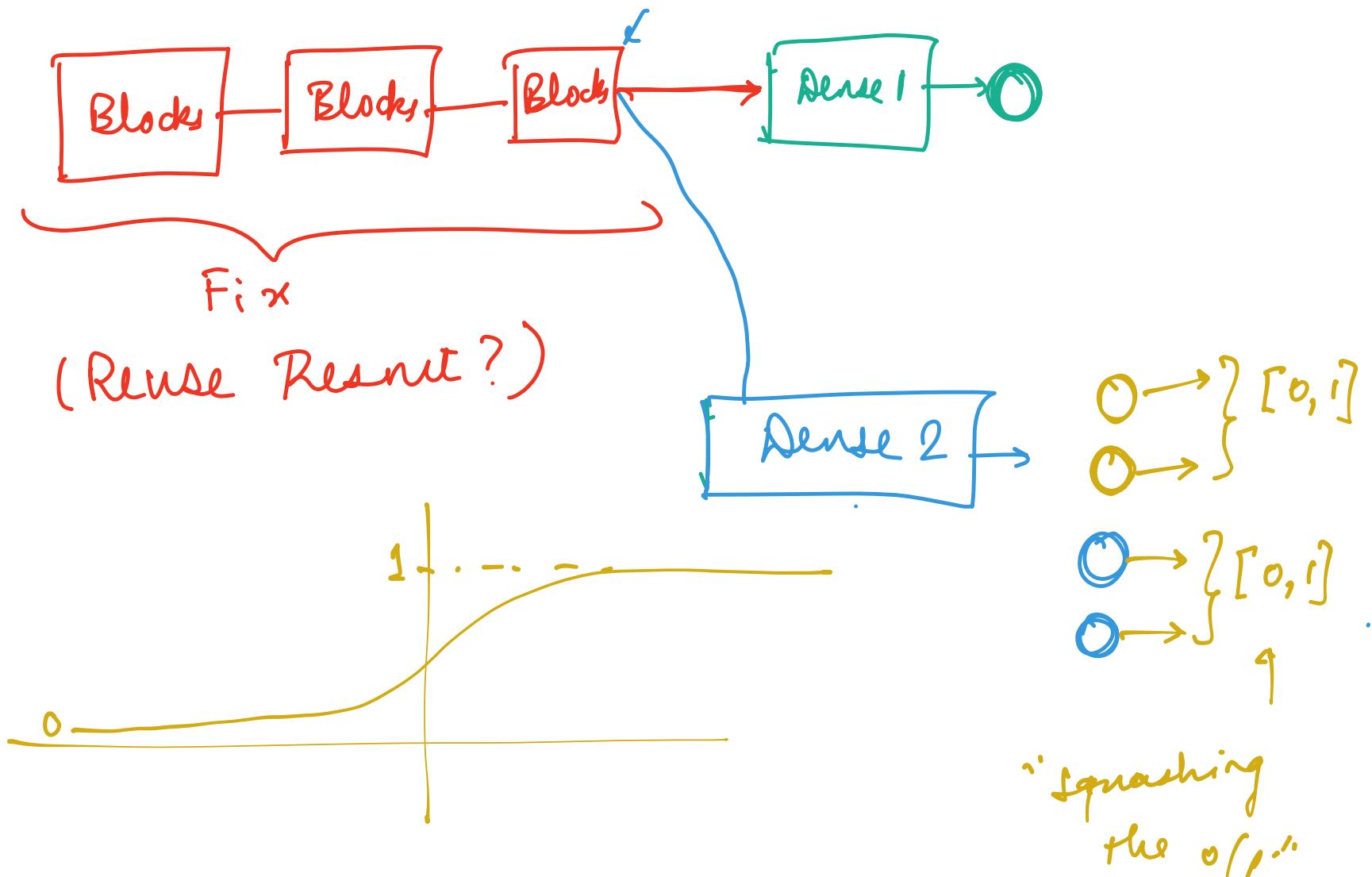
} normalized.

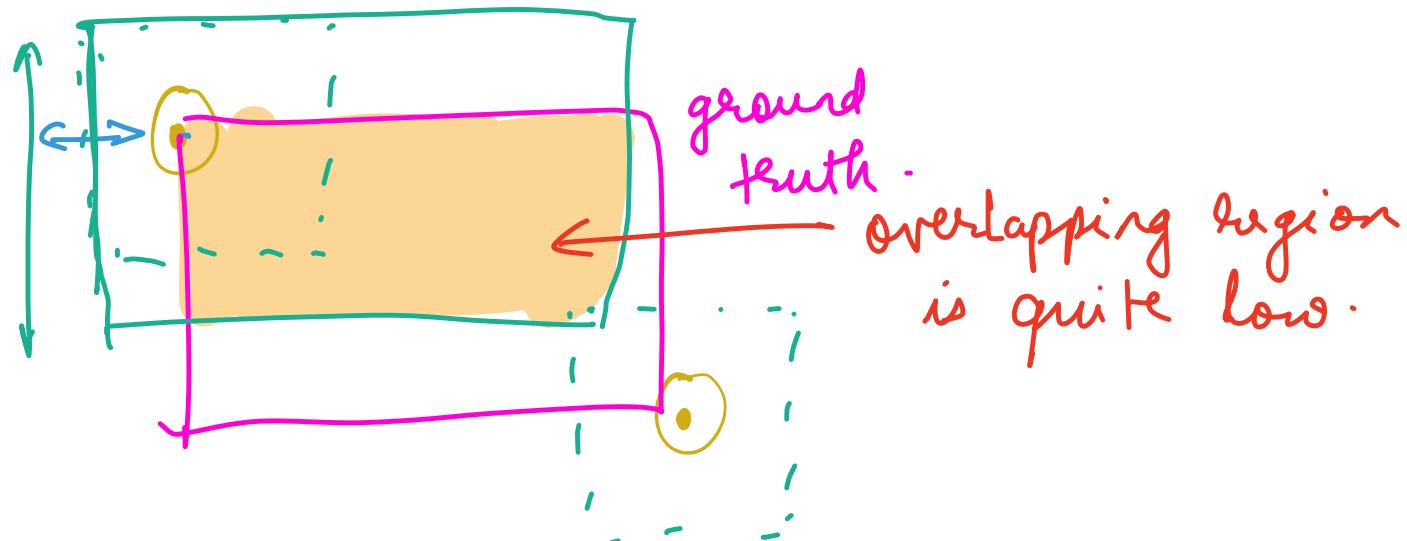
I/O relationship



Issue \rightarrow only 1 gun per image !!

Architecture for the CNN.





Interpreting  error.

~~$0.05 \times \text{height}$~~

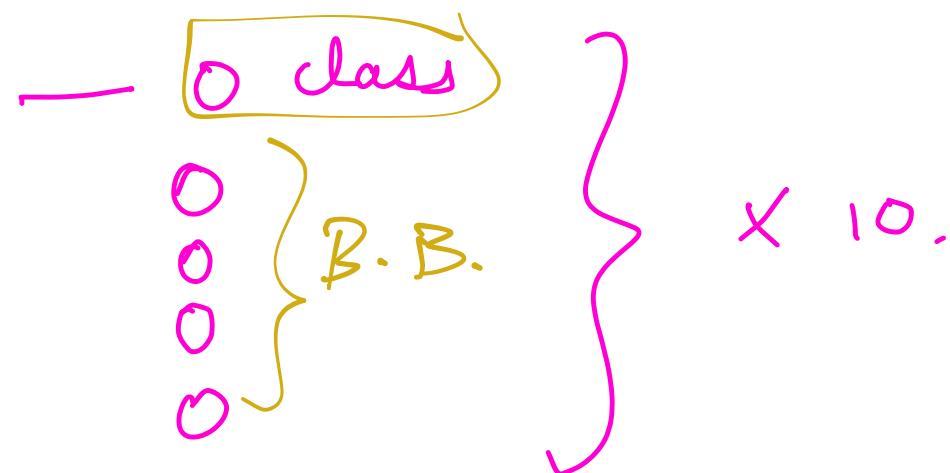
~~$0.05 \times \text{width}$~~

Discussion : Naive approach :

- ① Our model gives only 1 gun per image.
→ more suggestions.
- ② Our Model's performance is quite bad;
only 12% overlap.

How to fix issue 1 ?

* At max, we can have 10 guns.



Con: Not scalable
at all !!.

50 outputs in
o/p layer.

Suggestions :

✓ 1] Sliding window over image

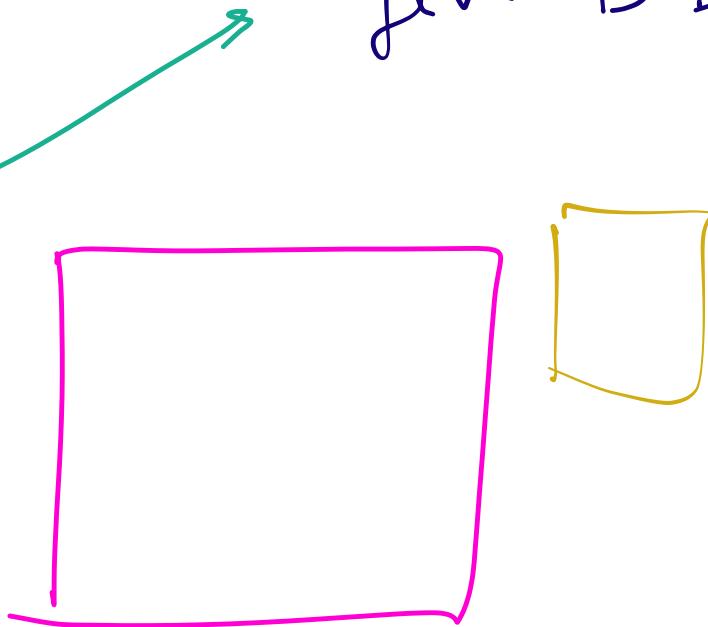


Predict if gun $\xrightarrow{\text{No}} 0.$

↓ yes

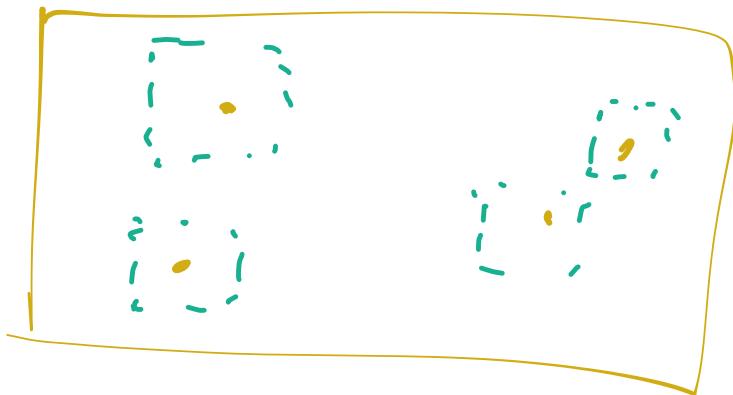
give B.B co-ordinates.

2]



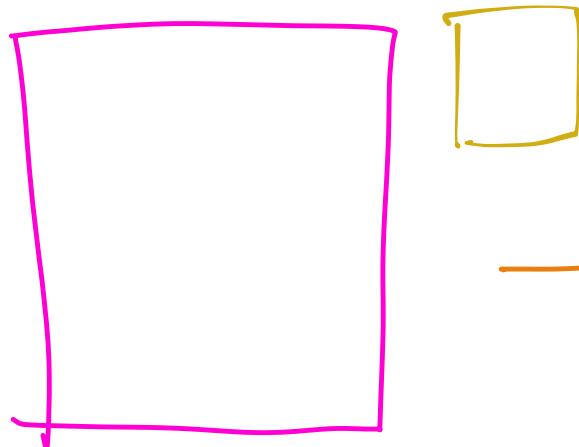
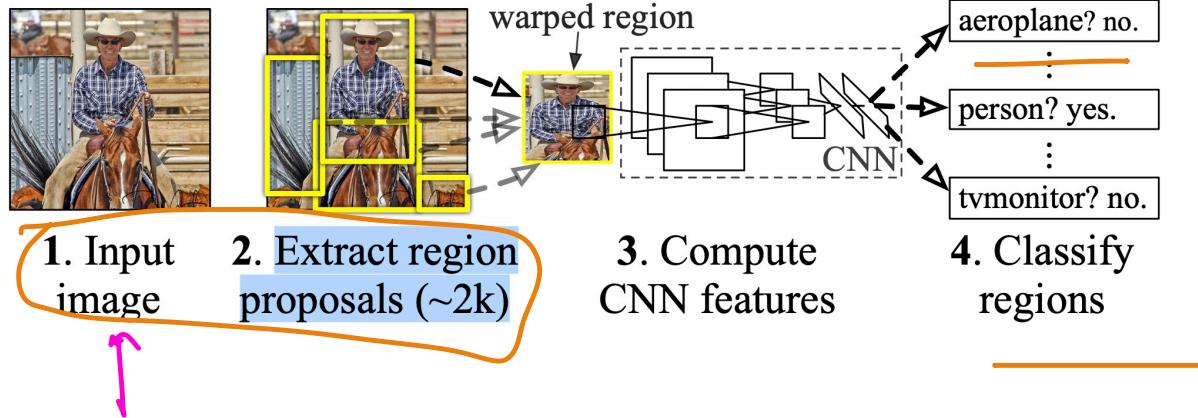
3] Use Grad CAM
to find
regions of
interest &
then give B.B.

g) Predict multiple centers and fit
B.B. around it.



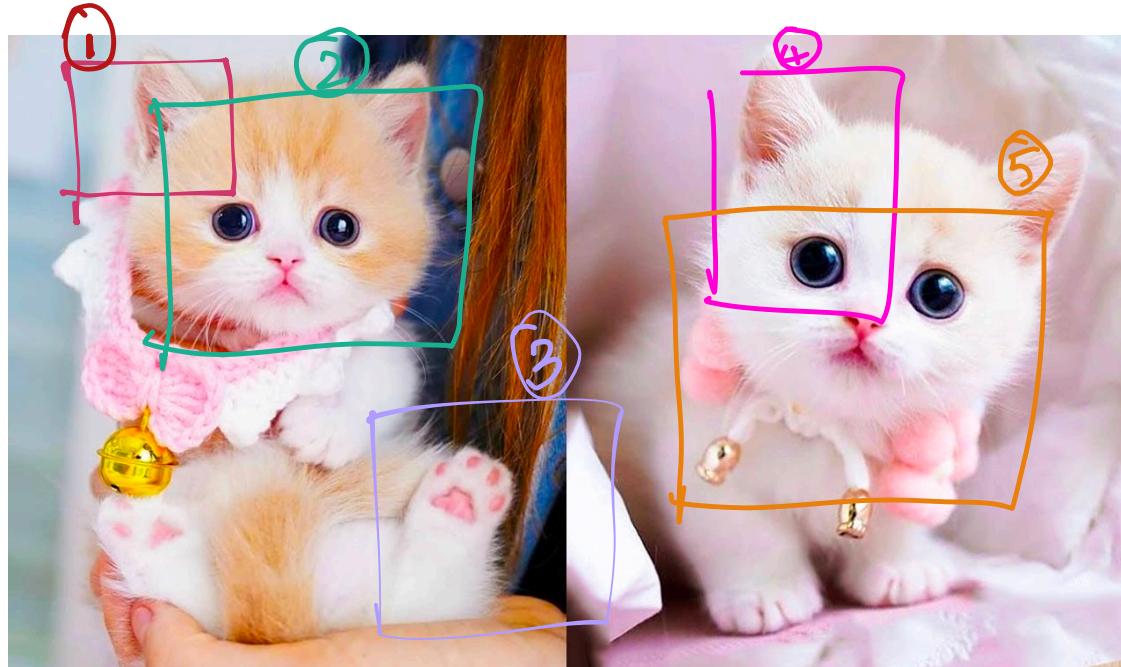
Solution 1 :

R-CNN: Regions with CNN features



Extract patches.

Pass each patch through an image classifier → get the prob.



Keep if confidence > 70%.

① - 10%

② - 90% +

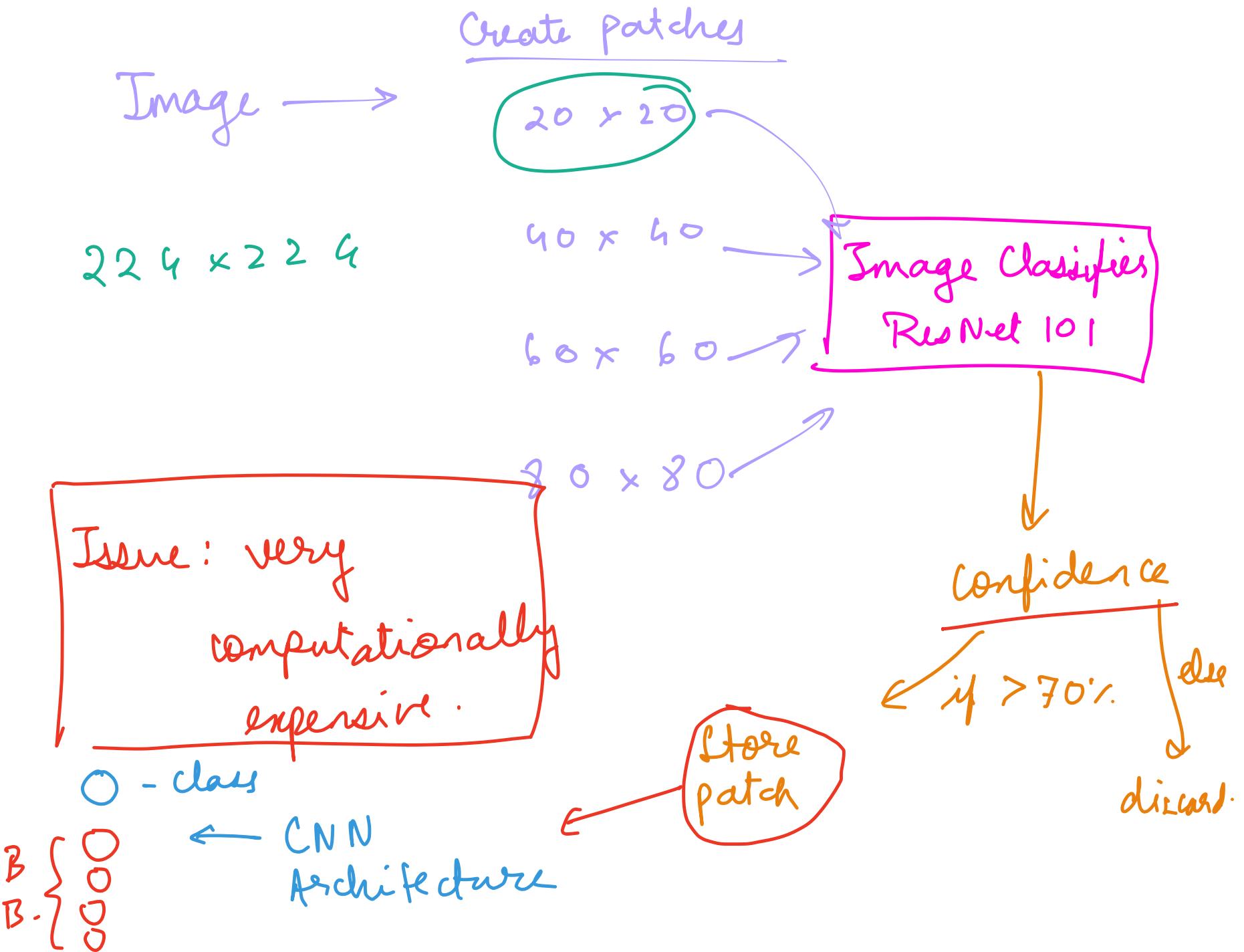
③ - 15%

④ - 45%

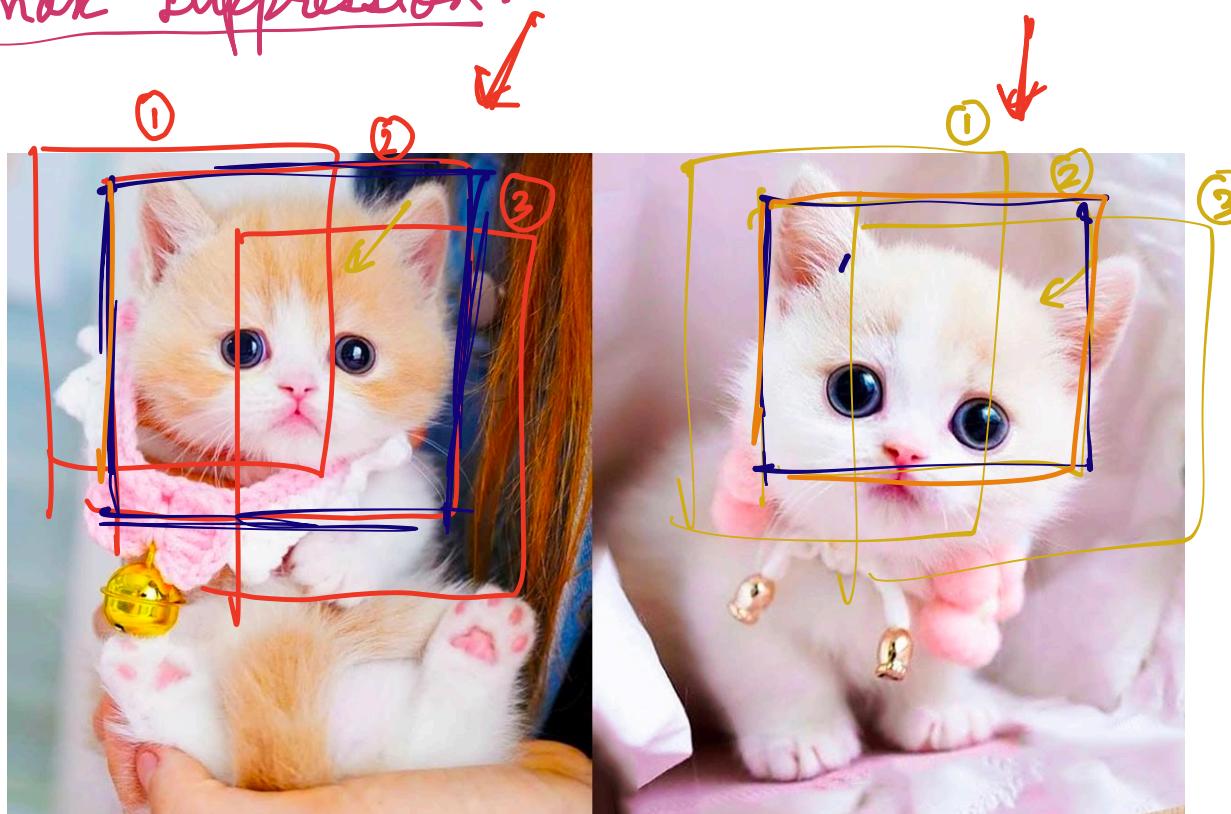
⑤ - 90%

②

⑤



Non max suppression:



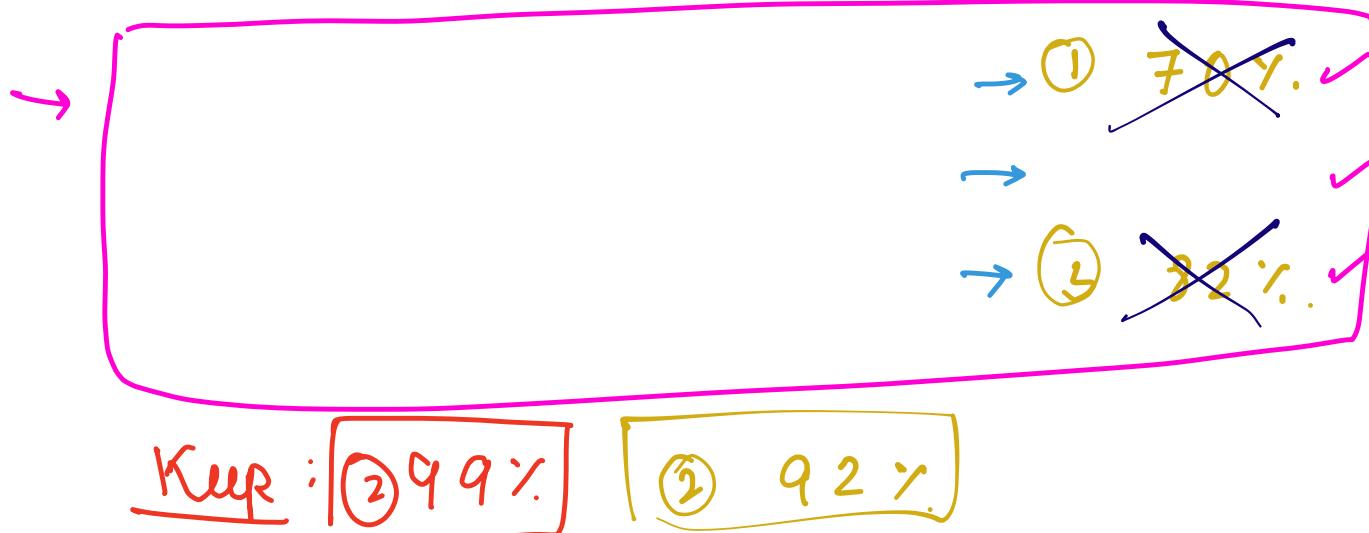
Step 1:

Select
the
max.

Step 2:

Eliminate
windows
which have
overlap.

IOU threshold.
50%.



Step 3: keep
repeating
1 & 2 till
array is empty

$3 \times 3 \times 3$ → Old way of doing convolutions.

3×3 → Depthwise convolution.

No FCN

If we want to make a CNN without a **fully connected layer** at the end, for the task of classification, What is the replacement for the 'Dense' Layers used in traditional CNNs?

HINTS



Complete Solution

You will get full points if and only if you give CORRECT ANSWER in first attempt. All later attempts will get you ZERO score.



Avg-Pooling Layers



Adding Dropout Layers

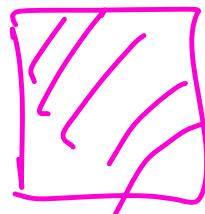


1x1 convolutions Layers



Max-Pooling Layers

correct answer .



Question wrong .

Data augmentation effectiveness



Suppose you wish to train a neural network to locate lions anywhere in the images, and you use a training dataset that has images similar to the ones shown above. In this case, if we apply the data augmentation techniques, it will be _____ as there is _____ in the training data.

HINTS



Hint 1

Complete Solution

You will get full points if and only if you give CORRECT ANSWER in first attempt. All later attempts will get you ZERO score.

- effective, size bias
- effective, position bias
- ineffective, position bias
- ineffective, angle bias