

March 21 (Tue).

Midterm is scheduled on
the 23rd (Thu). 1 hr Exam.

16:30 - 17:30

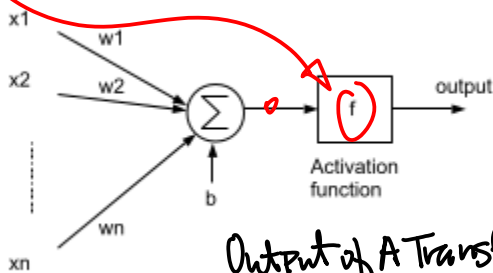
Then 15 minutes for prep &
Uploading the file.

Example: Softmax Activation
Function. MNIST CNN

Note:

$$\sigma(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}} \text{ for } i = 1, \dots, K \text{ and } \mathbf{z} = (z_1, \dots, z_K) \in \mathbb{R}^K. \quad \dots (1)$$

1° $f(\mathbf{z})$ Notation



Output of A Transfer
function

2° Index $f(z_i) = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}} \quad \dots (2)$

Total No. of Output Neurons = K
for Handwritten Digits Recognition
K=10;

where

$$\sum_{j=1}^K e^{z_j} = e^{z_1} + e^{z_2} + \dots + e^{z_K} \geq e^{z_i} \quad \dots (3)$$

Property 1.

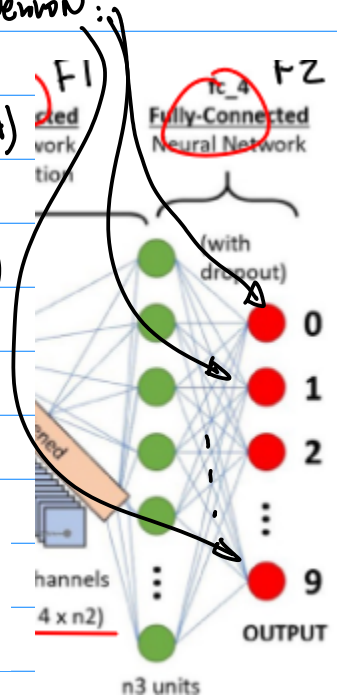
Output from Each Neuron:

$$f(z_1) = \frac{e^{z_1}}{\sum_{j=1}^{10} e^{z_j}} \quad \dots (4)$$

for Digit "0" (1st Output)

And

Dimension $\mathbf{z} \in \mathbb{R}^K$
Let $K=1$



$$f(z_2) = \frac{e^{z_2}}{\sum_{j=1}^{10} e^{z_j}} \text{ for Digit "1" (2nd output)} \quad \dots (5)$$

$$f(z_0) = \frac{e^{z_0}}{\sum_{j=1}^{10} e^{z_j}} \text{ for Digit "9" (10th o/p)}$$

if Add Eqn (1) + Eqn (2) + ...

$$f(z_1) + f(z_2) + \dots + f(z_0)$$

$$= \frac{e^{z_1}}{\sum_{j=1}^{10} e^{z_j}} + \frac{e^{z_2}}{\sum_{j=1}^{10} e^{z_j}} + \dots + \frac{e^{z_{10}}}{\sum_{j=1}^{10} e^{z_j}}$$

$$= \frac{\sum_{j=1}^{10} e^{z_j}}{\sum_{j=1}^{10} e^{z_j}} = 1. \quad \dots (6)$$

The softmax function takes as input a vector z of K real numbers, and normalizes it into a probability distribution consisting of K probabilities proportional to the exponentials of

$$f(z_i) = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

See Eqn (b).

Now, move to the 2nd Half of DCNN.

2022F-108a-Yolo-architecture-loss-function-2022-10-10.pdf

Typical Classification/Recognition Results are in Bounding Boxes

<https://arxiv.org/pdf/1506.02640v5.pdf>

S Divvala, Ross Girshick, Ali Farhadi
Allen Institute for AI, Facebook AI Research <http://pjreddie.com/yolo/>

Joseph Redmon*, Santosh Divvala[†], Ross Girshick*, Ali Farhadi[†]
University of Washington*, Allen Institute for AI[†], Facebook AI Research*
<http://pjreddie.com/yolo/>

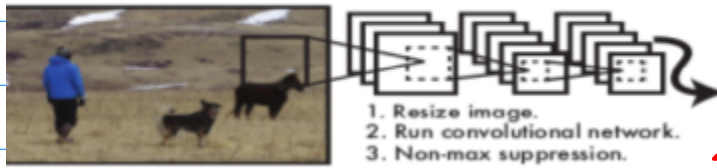


Figure 1: The YOLO Detection System. (1) resizes the input image to 448×448 , (2) runs a convolutional network, and (3) thresholds the resulting feature map.

Then, we would like to Achieve Semantic Segmentation.

Pixel by pixel

Based Segmentation/

Detection/Recognition



PART II (After the midterm)

April 4 (Tue)

Roadmap: Yolo (You-Only-Look-ONCE)

Semantic Segmentation.

Project Assignment to Implementation
Due 2 1/2 weeks. April 20th.
(Thursday)

CMPE258

Spring 2023

Homework (In-Class Presentation) Requirements Due April 6 (Thu).

1° One Paragraph Description (Abstract)
of the proposed Semester-Long
Project.

2° Title

Team members : First Name,
Last Name,
Major

Team Coordinator.
Contact E-mail.

3° Abstract Part.

Objective(s) : a) What is the
proposed work;

b) What is the coding / training /
Testing Task involved in
the project ?

c) Anticipated Result ?
And deliverable ?

d) Tools, platform, programming
Language Version, T.F.,

Pytorch, ChatGPT etc.

Also, Define Python Packages,

OpenCV.

Example: On Yolo.

Ref:

2022F-106-README-Tiny-Yolo4-GP...

Note 1: Readme for Yolo github
Installation & Testing.

Title: README Tiny Yolo v4 GPU Ubuntu

Document Number: 105-1b

CTI One Corporation

Table 1a. Document History

2022-10-6	Establish this document, document archive: (base) harry@workstation:~/yolo-2022-10-19\$	YY
-----------	--	----

1. Setup YOLO v4 environment

1.1. Clone the GitHub folder;

\$ git clone https://github.com/pythonlessons/TensorFlow-2.x-YOLOv3.git

1.2. Create YAML file for building the YOLO v4 Anaconda environment;

Create TensorFlow-2.x-YOLOv3/conda-gpu.yml as the following;

=====

name: yolo4-gpu

Ref: Introduction

2022F-108a-Yolo-architecture-loss-function-2022-10-10.pdf

Base-Line Ref for Yolo Technique

2022S-112-yolo-paper.pdf

You Only Look Once:
Unified, Real-Time Object Detection

Joseph Redmon*, Santosh Divvala*[†], Ross Girshick[¶], Ali Farhadi*[†]

University of Washington*, Allen Institute for AI[†], Facebook AI Research[¶]

<http://pjreddie.com/yolo/>

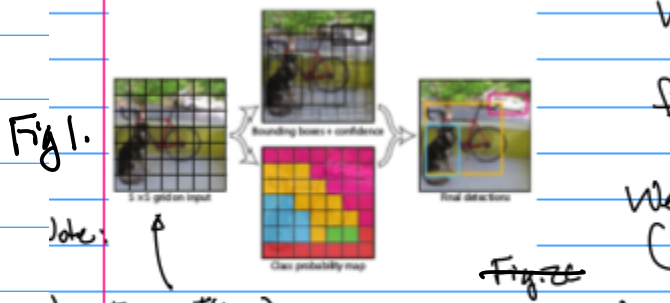
Lecture Notes: Base Line Ref / Requirements.

2022F-101-cmpe 258-note-2022-11-1.pdf

Example: Notations for Yolo.
Ref, pp 36.

CMPE258
Oct. 13, 22

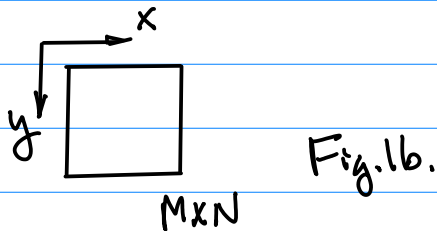
2. Bounding Boxes $B_{ij}(x, y) \dots (z)$
 i for x -direction, j for y -direction



1. Image $I(x, y)$ is divided into $S \times S$ grids.
Denote it as $G_{p,q}(x, y), \dots (i)$
where $p, q = 0, 1, 2, \dots$ indicate the

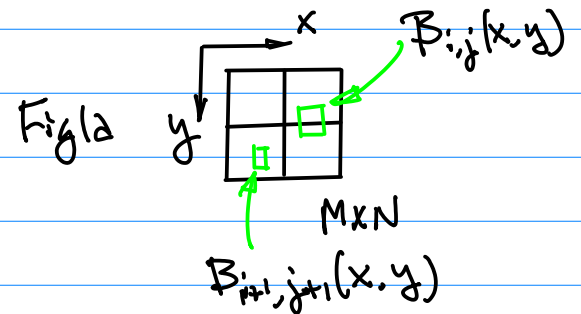
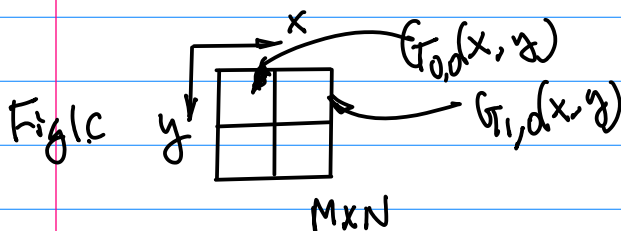
1. Image $I(x, y)$. with Resolution

$M \times N$
No. of Col. No. of Rows.



Divide $I(x, y)$ into $S \times S$ grids.
Each grid is Denoted as

$G(x, y) \dots (i)$
 p, q matches to x
Where $p, q = 0, 1, 2, \dots, S-1$



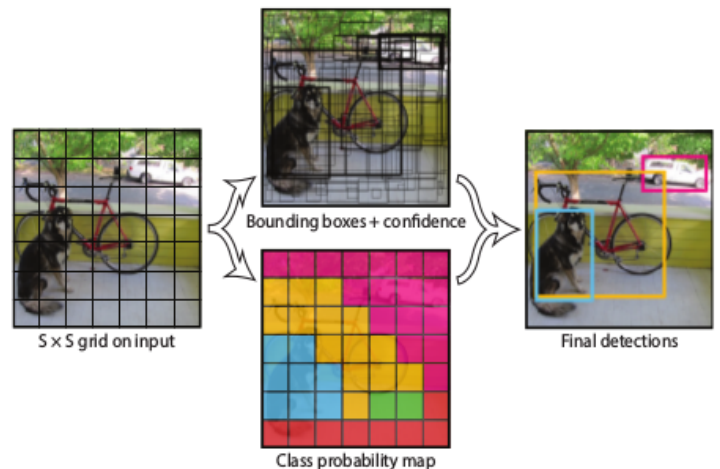
3. Five Parameters to define each Bounding Box.

(x, y) : Location of the top Left Corner of $B_{ij}(x, y)$

w, h : Width and Height of $B_{ij}(x, y)$

f : Confidence level, Probability distribution to Describe the likelihood of the B.B. (B^2) belongs to a certain Class of objects.

$(x, y, w, h, f) \dots (3)$



Note 1. α . $G_{ij}(x,y)$ Grid.

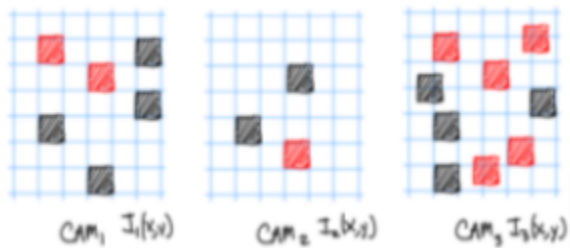
Figure 2: The Model. Our system models detection as a regression problem. It divides the image into an $S \times S$ grid and for each grid cell predicts B bounding boxes, confidence for those boxes, and C class probabilities. These predictions are encoded as an $S \times S \times (B * 5 + C)$ tensor.

$b, B_{ij}(x,y)$ Class probability.
 $C. (x,y,w,h,f)$ Probability
Confidence.

Aprih (Th).

Example: Discussion on Notation/Formulation.

Ref. [2022F-101-cmpe258-note-2022-11-1.pdf](#) PP38



Camera 1: $I_1(x,y)$ Camera 2: $I_2(x,y)$ Camera 3: $I_3(x,y)$

$$R = RI_1 + RI_2 + RI_3 \dots (1)$$

R : Red Squares, Persons.
 B : (Black) Vehicles. for "Union"

Intersection. " \cap "

Consider the probability of the event " R " (meaning Person(s) being captured on any one of these images).

$$\text{Prob}(R) = \text{Prob}(RI_1) + \text{Prob}(RI_2) + \text{Prob}(RI_3) \dots (2)$$

$$I_1 \cap I_2 \cap I_3 = \phi \text{ (Empty set)}$$

Consider Each Individual Camera:

$$\text{Prob}(RI_1) = \text{Prob}(R|I_1) \text{Prob}(I_1) \dots (3a)$$

Similarly,

$$\text{Prob}(RI_2) = \text{Prob}(R|I_2) \text{Prob}(I_2) \dots (3b)$$

$$\text{Prob}(RI_3) = \text{Prob}(R|I_3) \text{Prob}(I_3) \dots (3c)$$

Rewrite Egn (2):

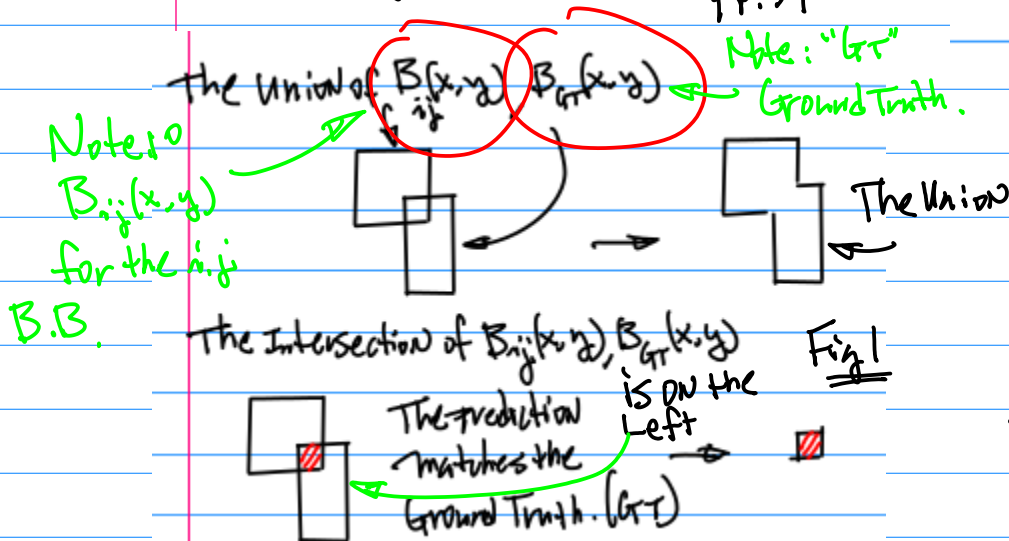
$$\begin{aligned} \text{Prob}(R) &= \text{Prob}(R|I_1) \text{Prob}(I_1) \\ &+ \text{Prob}(R|I_2) \text{Prob}(I_2) \\ &+ \text{Prob}(R|I_3) \text{Prob}(I_3) \\ &= \sum_{i=1}^3 \text{Prob}(R|I_i) \text{Prob}(I_i) \dots (4) \end{aligned}$$

Ref: 20225-112-yolo-paper.pdf

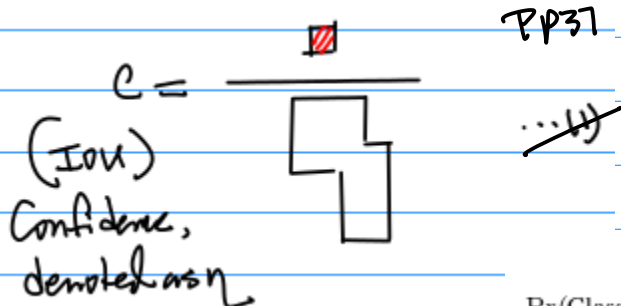
$$\Pr(\text{Class}_i | \text{Object}) * \Pr(\text{Object}) * \text{IOU}_{\text{pred}}^{\text{truth}} = \Pr(\text{Class}_i) * \text{IOU}_{\text{pred}}^{\text{truth}} \quad (1)$$

Note 1. \uparrow Conditional Probability. \uparrow $\Pr(C_i)$
 $\Pr(C_i | \text{Obj}) \Pr(\text{Obj})$

1. IOU (Intersection of Union)
Index for the purpose of Comparing
2 Bounding Boxes at a time
PP.37



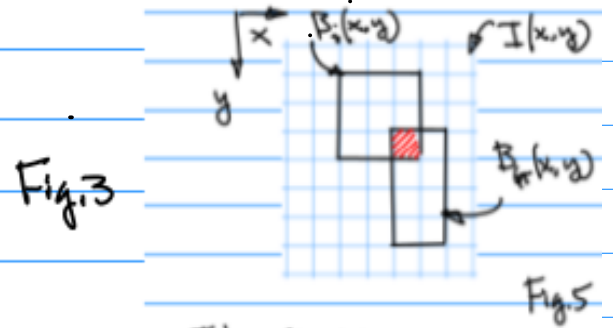
$$\text{IOU} = \frac{\text{Intersection of } B_{ij}(x,y) \text{ and } B_{gr}(x,y)}{\text{The Union of } B_{ij}(x,y) \text{ and } B_{gr}(x,y)} \dots (5)$$



Hand Calculation of IOU.

Example: PP36

5. IOU (Intersection of Union).
Example: Illustration of IOU



Sol. First, find the Number of
pixels of the $B_{ij}(x,y) \cap B_{gr}(x,y)$

$$B_{ij}(x,y) \cap B_{gr}(x,y) = 1$$

then, Find the Union

$$B_{ij}(x,y) \cup B_{gr}(x,y) = N[B_{ij}(x,y)] + N[B_{gr}(x,y)] - N[B_{ij}(x,y) \cap B_{gr}(x,y)]$$

$$= (3 \times 3) + (2 \times 4) - 1 = 9 + 8 - 1 = 17 - 1 = 16$$

$$\therefore \text{IOU} = \frac{1}{16}$$

Now, from Eqn (1) of the Ref. (Research Paper)

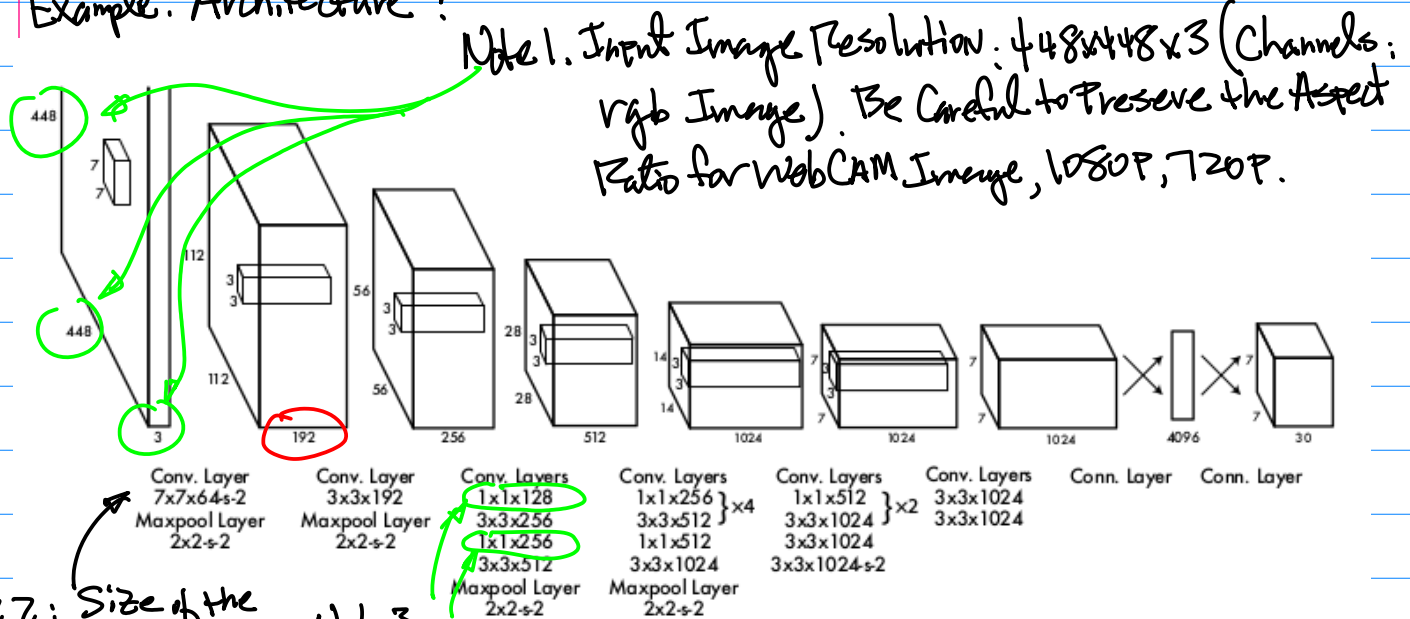
Note

$$\Pr(\text{Class}_i | \text{Object}) * \Pr(\text{Object}) * \text{IOU}_{\text{pred}}^{\text{truth}} = \Pr(\text{Class}_i) * \text{IOU}_{\text{pred}}^{\text{truth}} \quad (1)$$

IOU^{truth} ← "G.T." B.B.
IOU^{pred} ← Predicted B.B.

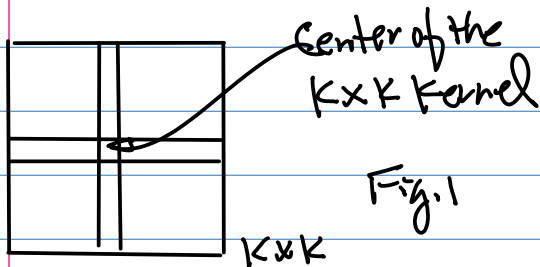
YOLO

Example: Architecture:



Note 2: Size of the Kernel: 7×7 , 64 of them.

Note 3: 1×1 Convolution is utilized here



For $K \times K$ 2D Convolution, the output of the convolution. "Spatial Information", Neighbouring pixels under the $K \times K$ kernel are counted for (for feature extraction @ the center of the kernel)

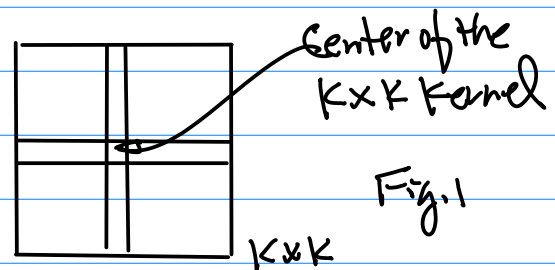
April 1 (Tue).

Presentation (Brief) on Each Team Project.

Example: Continuation on With Formulation.

1×1 Convolution.

Note 1. Background on $K \times K$ Convolution.



Output: 1 pixel
Input: $K \times K$ pixels.

Captures All Neighbouring pixels at a time, And Produce one pixel Output.

Note 2. For each convolution kernel, the convolution conducted will result in one output feature layer

As we continue the Convolution Process, the Number of Output Feature Layers will grow Significantly, Therefore, there's a need to Reduce the Number of Layers without missing crucial features.

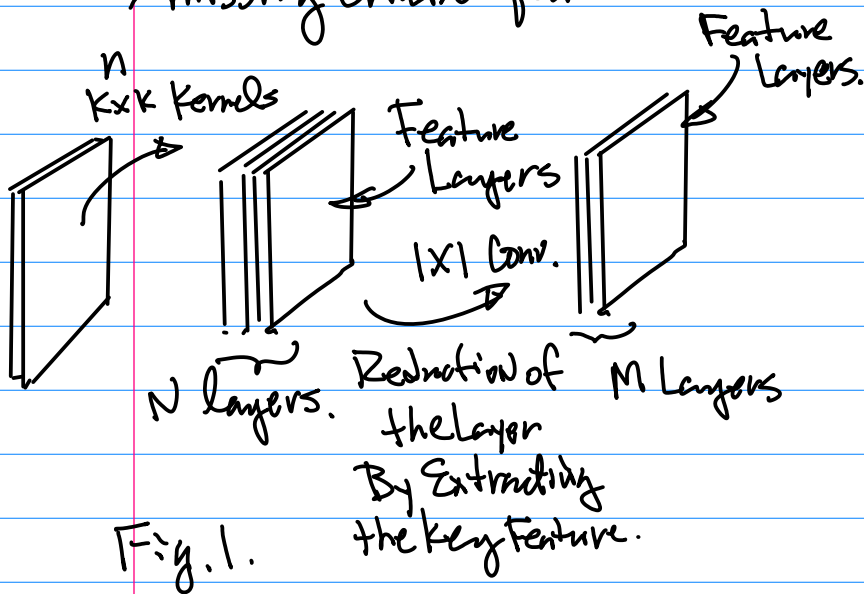


Fig. 1.

To Be Able to Extract/Preserve the Key features to Achieve Reduction of Layers. We are using the following Technique.

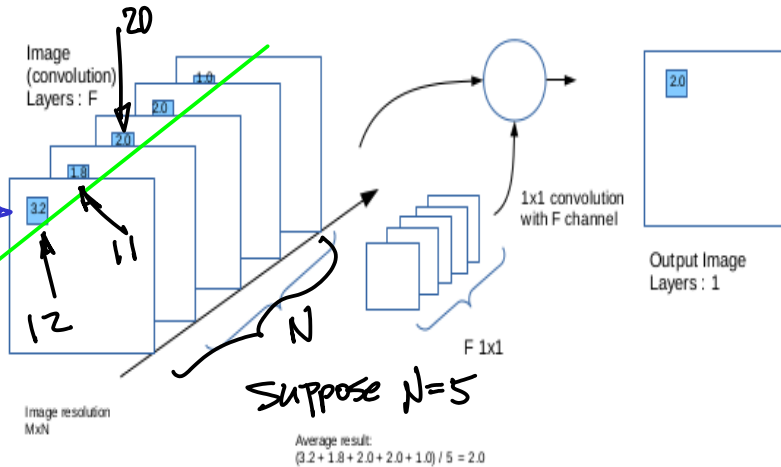
1x1 Convolution for Dimension Reduction and Pooling

The 1x1 convolution enables dimension reduction by reducing the number of channels in convolution layers

1. Suppose the input layers is $C \times H \times W$, where C is its channels. The 1x1 convolution generates one average result in shape $H \times W$. The 1x1 (filter) is a vector of length C .
2. Now if you have F 1x1 filters, you get F layers of output, the output shape is $F \times H \times W$. For input layer $C \times H \times W$ with F 1x1 convolution (with channel is C), you will get $F \times H \times W$ layers.

Note 1.

One pixel across the entire stack of the feature layer.



Harry Li, Ph.D.

Reduction Requirement: Combine $N=5$ layers into 1 layer.

To preserve the feature in this process.

Question: What is the technique to combine them (pixels at different layers) with equal contribution from each layer?

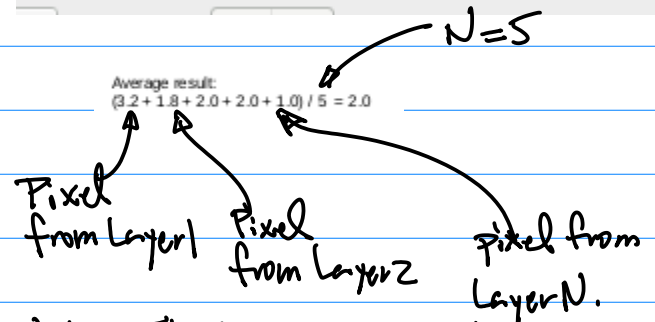
$$\frac{1}{N} [I_1(x_1, y_1) + I_2(x_1, y_1) + \dots + I_5(x_1, y_1)] \dots (1)$$

$$\frac{1}{N} I_1(x_1, y_1) + \frac{1}{N} I_2(x_1, y_1) + \dots + \frac{1}{N} I_5(x_1, y_1)$$

April 12 (Th).

Example: 1x1 Convolution.

2022F-108a-Yolo-architecture-loss-function-2022-10-10.pdf



Note: The Average operation treats the feature Equally from Each Layer.

More General Case:

$$\alpha_1 I_1(x_1, y_1) + \alpha_2 I_2(x_1, y_1) + \dots + \alpha_N I_N(x_1, y_1)$$

where $\alpha_1 + \alpha_2 + \dots + \alpha_N = 1 \dots (2)$

Note: Project Assignment (10 pts)
on Object Recognition.
Due April 30th (Sunday, 11:59pm).

CMPE258
YOLACT Semantic Segmentation and Comparison with YOLO
HL

PART I YOLO (5 points)

test run YOLO v4 based on the readme document given in the class github
<https://github.com/hualili/opencv/blob/master/deep-learning-2022s/2022F-106-README-YY-HL-v2-2022-10-20.pdf>

smart phone to take record a 15 - 30 second video clip for YOLO v4 object

Note: To Config Anaconda Environment,
Use the Pre-Created Configuration
files here, from the github

2022S-104b-conda-gpu.yml

2022S-104c-conda-cpu.yml

Then, Create Conda environment
Config. file

\$ conda env create -f conda-gpu.yml

Next, Activate the Conda Env.

\$ conda activate yolo4-gpu

Then, perform the following task

1.5. Download Tiny YOLO v4 model files;

\$ wget -P model_data
<https://github.com/AlexeyAB/darknet/>

1.6. Modify the configuration file, TensorFlow-2.x-YOLOv3/yolov3/configs.py;

13 YOLO_TYPE = "yolov4" # yolov4 or yolov3

37 TRAIN_YOLO_TINY = True

Example: Yolo Architecture.

Note: 1° Understand the Composition of the Yolo Architecture.

To Be Able to Describe the function (s) of Each Block;

2° Analyze the Parameter(s) of Each Block.

Make the following modification to Run Yolo for the Video file Input. Need to modify the python Code,

1. Input image size: 448x448x3; 2. Resolution reduction for feature extraction/abstraction Pooling and convolution with stride = 2;

Base Line Yolo Architecture

Design guideline: The block of convolutional layers to extract image features, the fully connected layers to predict the output probabilities and the locations (coordinates).

Diagram illustrating a 3D CNN architecture for video classification. The input is a 3D volume (448x448x3). The architecture consists of several layers:

- Conv. Layer** (7x7x64+2)
- Maxpool Layer** (2x2+2)
- Conv. Layer** (3x3x192)
- Maxpool Layer** (2x2+2)
- Conv. Layer** (1x1x128)
- Conv. Layer** (3x3x256)
- Conv. Layer** (1x1x256)
- Conv. Layer** (3x3x512)
- Conv. Layer** (3x3x1024)
- Conv. Layer** (1x1x256) x4
- Conv. Layer** (3x3x512) x4
- Conv. Layer** (1x1x512)
- Conv. Layer** (3x3x1024) x2
- Conv. Layer** (3x3x1024)
- Conv. Layer** (3x3x1024+2)
- Conv. Layer** (3x3x1024)
- Conv. Layer** (3x3x1024)

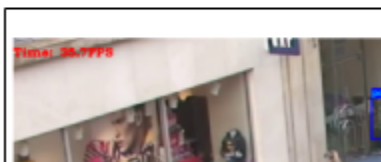
The final output is a 3D volume (4096x7x30).

For Feature Extraction
Convolutional Operations.
Like/Similar to the Creation of
high Dimensional feature Vectors.

$$\mathbb{X}(x_1, x_2, \dots, x_n)$$

Note: Image/Video Input should be square with preservation of the Original Aspect Ratio.

1. Input image size: 448x448x3; 2. Resolution reduction for feature extraction/abstraction Pooling and convolution with stride = 2;



Be Sure to Record Video Clip(s) Yourself for YoLo Testing, you need to use these Videos in the future.

Stride = 1 Convolution.

"default" version of Convolution.

Shift/move ~~the mask~~ from Top left
corner Left to Right ONE pixel
at a time, And Top to Bottom ONE
Row at a time.

Stride = 2.

Shift/move ^{the} mask from Top left
corner Left to Right 2 pixels
at a time, And Top to Bottom 2
Row at a time.

Example/Exercise Off-Line.

Sub-Sampling + Convolution \rightarrow Reduction
of feature layers.