# Convolution Neural Networks - CNN
# Part II

Deep Neural Network

Session 21

Pramod Sharma
pramod.sharma@prasami.com

---

## Agenda

2

- Introduction
- Classical Networks
- Network in Network
- Inception Network
- Transfer Learning
- Object Detection

pra-sami

## Classic Networks

3

LeNet-5

AlexNet

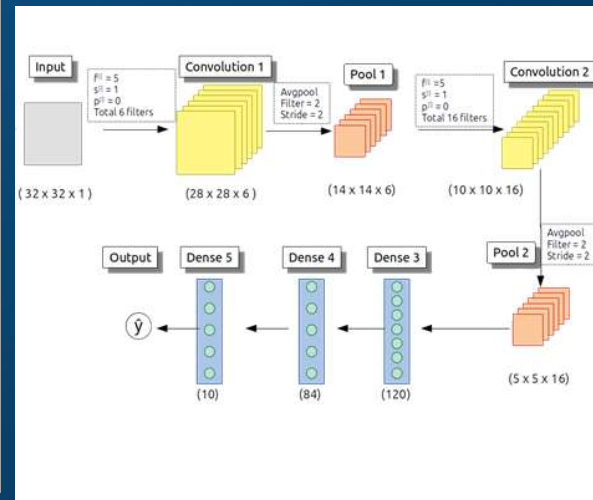VGG

pra-sami

## Classic Networks

4

ResNet

DenseNet

Unet

pra-sami

## LeNet - 5

5

- ❑ LeCun et. Al., 1998 – Gradient based learning applied to document recognition

- ❑ A number of Conv and Pool layers stacked together

- ❑ Followed by dense layers

- ❑ Softmax activation to predict probabilities

- ❑ Original LeNet -5 had 32 x 32 x 1 images and was used for handwriting dataset

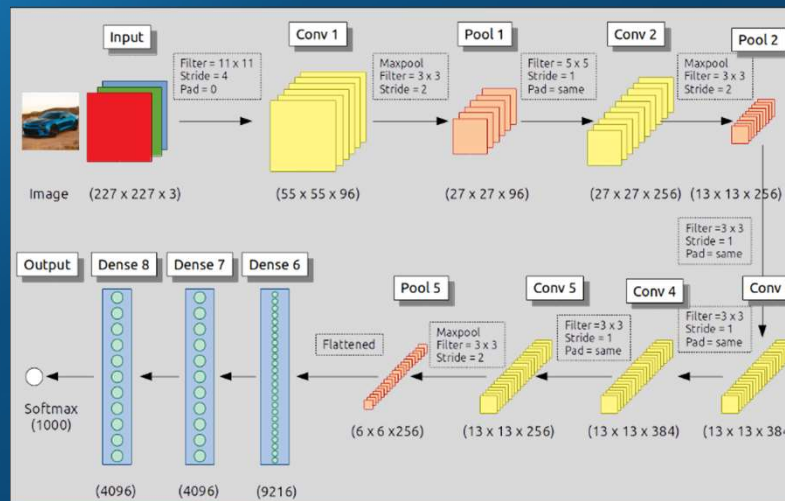- ❑ Had Average Pooling and used Tanh activation



12/3/2024

pra-sami

## AlexNet

6



- ❑ Alex net was considered very deep back then
  - ❖ It used ReLU
- ❑ First one to use 'Local Response Norm' and prove that it's not a good idea
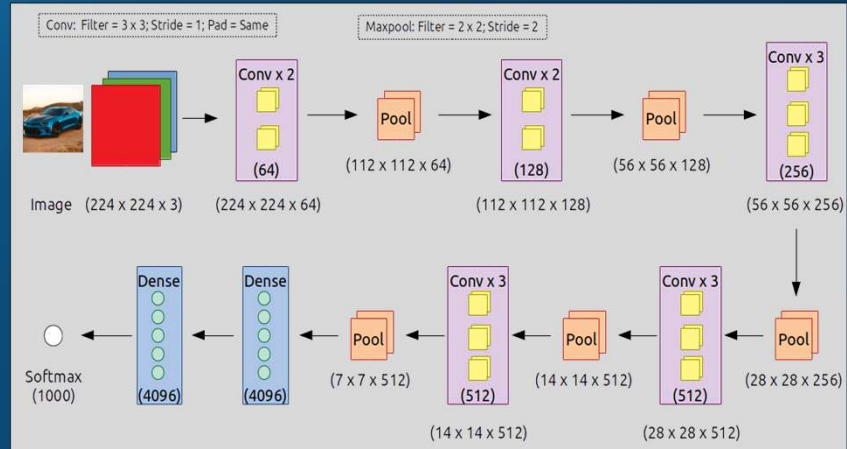
12/3/2024

pra-sami

# VGG-16

7

- ❑ Standardized the parameters
- ❑ It had 16 layers with weights
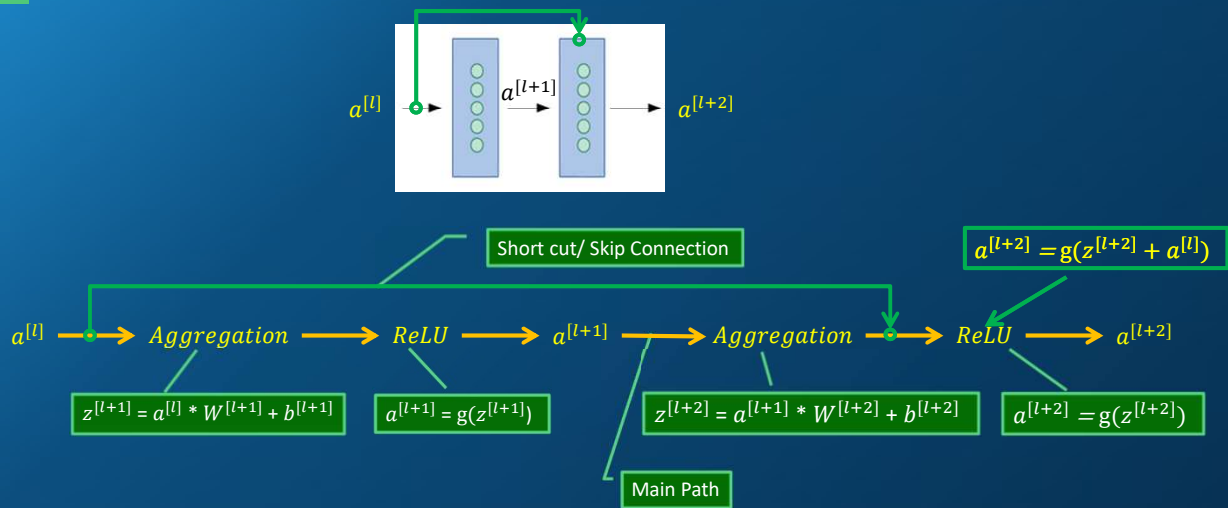- ❑ Uniformity made it very attractive for researchers

pra-sami

---

8

Those were Classical Networks

pra-sami

## Residual Block



$a^{[l]}$    $a^{[l+1]}$    $a^{[l+2]}$

Short cut/ Skip Connection

$$a^{[l+2]} = g(z^{[l+2]} + a^{[l]})$$

$a^{[l]} \longrightarrow Aggregation \longrightarrow ReLU \longrightarrow a^{[l+1]} \longrightarrow Aggregation \longrightarrow ReLU \longrightarrow a^{[l+2]}$

$$z^{[l+1]} = a^{[l]} * W^{[l+1]} + b^{[l+1]}$$

$$a^{[l+1]} = g(z^{[l+1]})$$

$$z^{[l+2]} = a^{[l+1]} * W^{[l+2]} + b^{[l+2]}$$
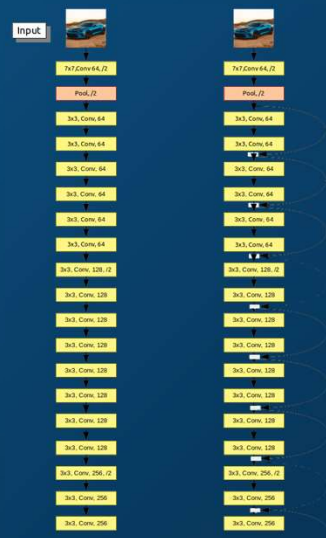
$$a^{[l+2]} = g(z^{[l+2]})$$

Main Path

12/3/2024

pra-sami

## ResNet

- ❑ Deeper networks had vanishing gradient problems

- ❑ Most networks resulted in higher errors and lesser accuracy as the depth increased

- ❑ ReLU activations solved it to some extent

- ❑ As networks became deeper (more layers), it lead to higher classification error

- ❑ It was not due to over-fitting as, as training errors were higher too!

- ❑ Expectation was that network with more layers should be as good if not better!

- ❑ Deeper networks are not good handling identity function (Output same as input)

- ❑ ResNet Architecture addressed it



12/3/2024

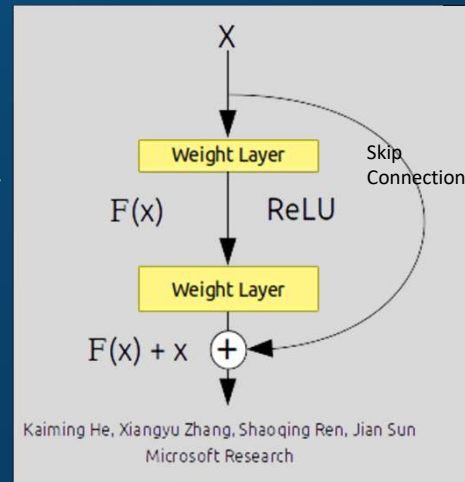pra-sami

5

## ResNet – Building block

11

- ❑ For normal convolutions:
  - ❖ F(a) = F(a) + a
- ❑ In case of Pooling
  - ❖ F(a) = F(a) + a . Ws
  - ❖ Where Ws is matrix of <previous layer size> x <size of layer L+2>



X

Weight Layer

F(x) ReLU

Weight Layer

F(x) + x

Skip Connection

Kaiming He, Xiangyu Zhang, Shaoqing Ren, Jian Sun
Microsoft Research

pra-sami

---

## ResNet – Building block

12

- ❑ if F(x) becomes zero, it is at least x
  - ❖ Relies on making identity function explicit
  - ❖ Simply, Input 'x' is processed by two conv. layers as earlier
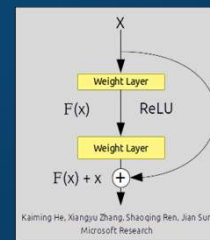  - ❖ Then 'x' is added to the output before applying ReLU

- ❑ Thus it is catering to both.
  - ❖ Old abstracts are retained and additional abstracts if any are added!

- ❑ Early layers are trying to learn some low level features such as edges, corners etc,
  - ❖ Later layers are focusing on high level abstractions such as wheels, wind shield, etc…
  - ❖ Subsequent layers may degrade or obfuscate these reliable signals
  - ❖ ResNet architecture gives the network a more explicit codes the output of the block defaulting to its input x, if F(x) is zero

- ❑ In short, don't forget what you have already learnt, at least….

pra-sami

## 1 x 1 Convolution – Network in Network



Not so obvious in a single layer…
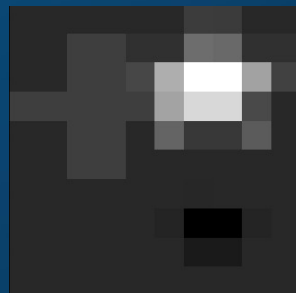
Lin et al., 2013 Network in Network

12/3/2024

pra-sami

## 1 x 1 Convolution – multiple layers


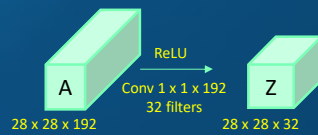
Nonlinearity is introduced over multiple layers…
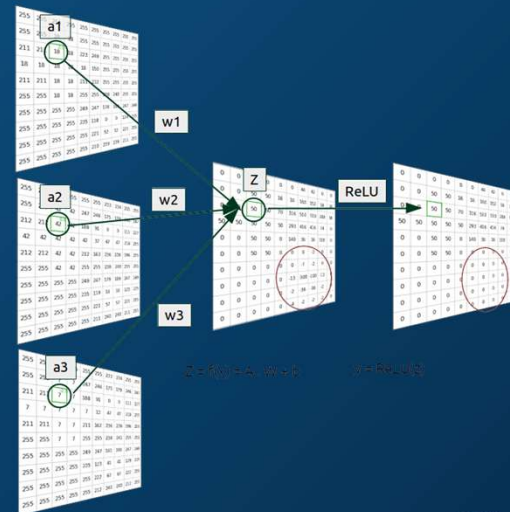
ReLU

12/3/2024

## Network in Network

15

- ❏ Another advantage is that it can be used to reduce dimensions

- ❏ Thus allowing us to shrink or expand or keep the averages of the channels,

- ❏ Of course, it permits us to add non-linearity

A → ReLU, Conv 1 x 1 x 192, 32 filters → Z

28 x 28 x 192          28 x 28 x 32

Network in Network

---

## Inception Network - Acknowledgements

16

- ❏ Takes inspiration from movie "Inception"… "We need to go deeper"

Going deeper with convolutions

| Christian Szegedy | Wei Liu | Yangqing Jia |
| Google Inc. | University of North Carolina, Chapel Hill | Google Inc. |

| Pierre Sermanet | Scott Reed | Dragomir Anguelov | Dumitru Erhan |
| Google Inc. | University of Michigan | Google Inc. | Google Inc. |

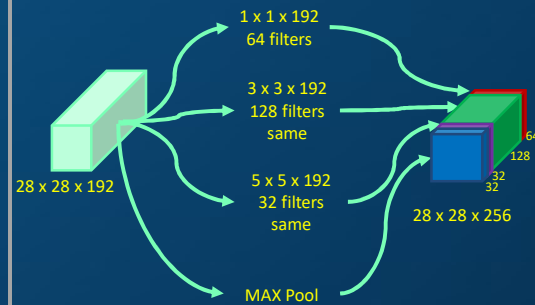| Vincent Vanhoucke | Andrew Rabinovich |
| Google Inc. | Google Inc |

# Inception Network – Building Block

17

- We are always faces with challenge of selecting the filters, pooling and their respective sizes

- Engineers though of a solution of adding all together and let the network decide what works best

- Enter combination of filters

- It has problem of computational cost

- Note that you have to use Padding with stride of one in the MaxPool layer to match the dimensions

28 x 28 x 192

1 x 1 x 192
64 filters

3 x 3 x 192
128 filters
same

5 x 5 x 192
32 filters
same

MAX Pool

64
128
32
32

28 x 28 x 256

12/3/2024

pra-sami

---

# Inception Network – Computational Cost

18

- Let's take one filter as an example

28 x 28 x 192

5 x 5 x 192
32 filters
same

28 x 28 x 32

- Overall computations:
  - ❖ 5 x 5 x 192 x 28 x 28 x 32 = 120,422,400
  - ❖ Say = 120 million

- A very computationally heavy operation

12/3/2024

pra-sami

9
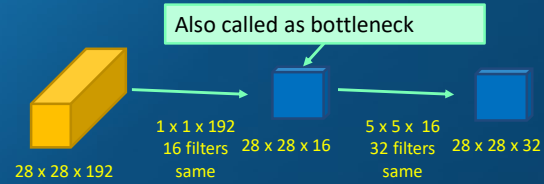
## Inception Network – Computational Cost

19

- ❏ Alternatively,

Also called as bottleneck

28 x 28 x 192

1 x 1 x 192
16 filters    28 x 28 x 16
same

5 x 5 x 16
32 filters    28 x 28 x 32
same

- ❏ Overall computations
  = {(1 x 1 x 192) x (28 x 28 x 16)} + {(5 x 5 x 16) x (28 x 28 x 32)}    = 2,408,448 + 10,035,200   = 12,443,648  Say = 12 million

- ❏ Reduced by 10 times!

- ❏ Caution: the size of bottleneck layer to be chosen carefully too much shrinking may harm the performance

- ❏ Also Helping us in reducing the number of channels!

12/3/2024

pra-sami

## Inception Module

20

DepthConcat

depth = 64+128+32+32 = 256

28 x 28 x 64       28 x 28 x 128       28 x 28 x 32

Hence add a 1 x 1 conv layer

Conv 1x1+1(S)     Conv 3x3+1(S)     Conv 5x5+1(S)     Conv 1x1+1(S)

28 x 28 x 32

# of channels is still 192

28 x 28 x 96       28 x 28 x 16

Conv 1x1+1(S)     Conv 1x1+1(S)     MaxPool 3x3+1(S)

28 x 28 x 192

Same padding

DepthConcat     28 x 28 x 192

12/3/2024

pra-sami

10

## Complete Network - GoogLeNet



Side branches to keep tab on over fitting

## Transfer Learning

❑ May take days or even weeks to train on very large datasets.

❑ In AI and ML world, its customary to publish one's work in open source
  ❖ Open source large datasets, pre-trained models and weights available
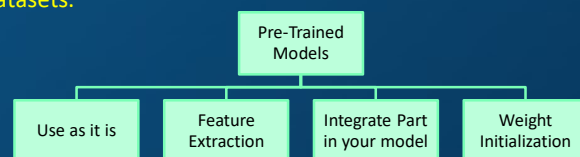
❑ Especially helpful in cases where we have limited pictures

❑ The models are complex and have multiple classes
  ❖ Image net ➔ 1000 classes (ImageNet Large Scale Visual Recognition Challenge, or ILSVRC or ImageNet)
  ❖ A range of high-performing models available

❑ Use top performing model directly, or integrated into a new model

❑ Of course with some modifications to last few layers
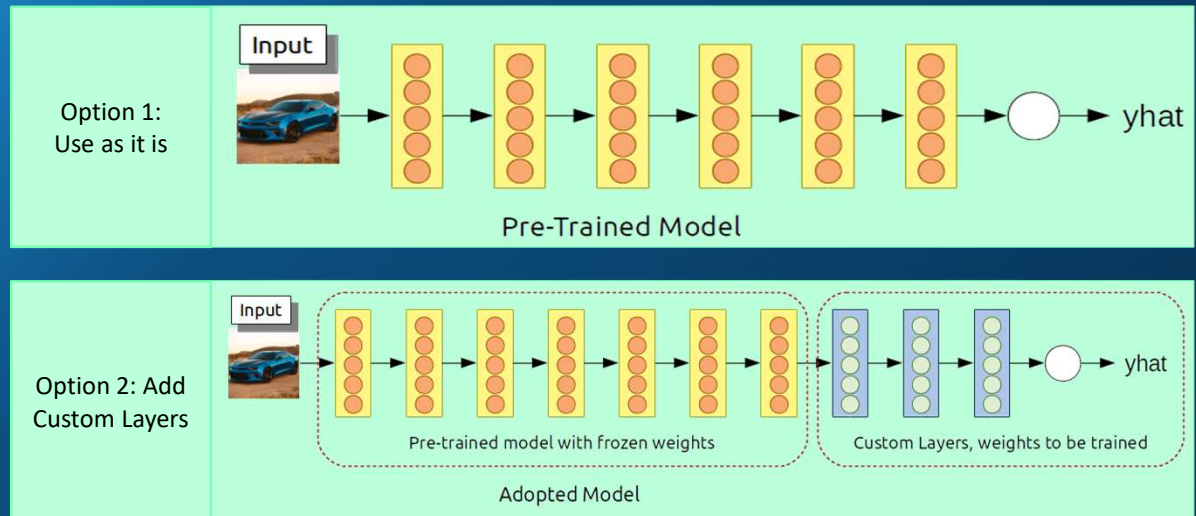
❑ Most pre-trained models APIs are available



Pre-Trained Models

| Use as it is | Feature Extraction | Integrate Part in your model | Weight Initialization |

12/3/2024

11

## Transfer Learning Options

23

**Option 1: Use as it is**

Input

Pre-Trained Model

yhat

**Option 2: Add Custom Layers**

Input

Pre-trained model with frozen weights

Custom Layers, weights to be trained

yhat

Adopted Model

pra-sami

---

## Transfer Learning Option : 3

24

Input

Pre-trained model with frozen weights

Save Features

**Extract Features**

Saved Features

Custom Layers, weights to be trained

yhat

**Train smaller model on saved features**

❑ Feel free to experiment by training frozen layers as well!

❑ If you have more data more layers could be used.

❑ It there is lots and lots of data, use this model to initialize and train all the weights

❑ These models are so well trained, it advantage to use existing weights!!

pra-sami

## Landmark Detection

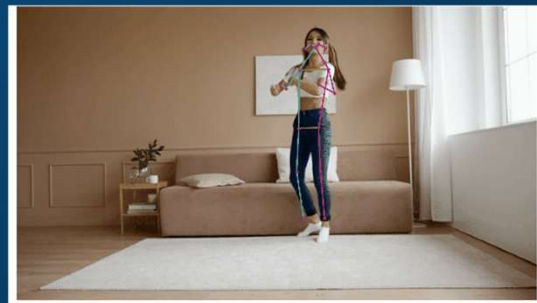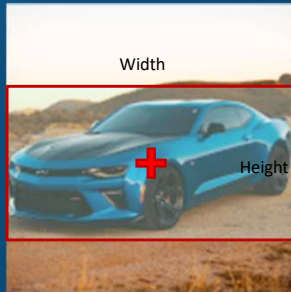## Gait Detection

## Object Localization

27



Classification:
Identify object
It's a Car!

- Pedestrian
- Car
- Truck
- Bike
- others

Classification with localization:
Identify object and mark its location

- Class of object
- Location of bounding box (mid point, height, width)
- $\hat{y}$ will be a vector

Detection:
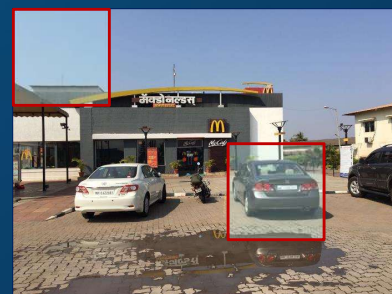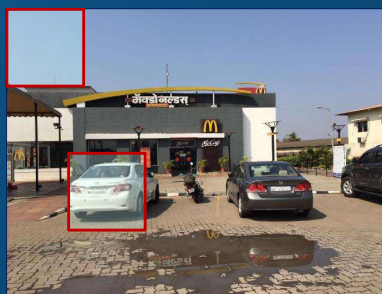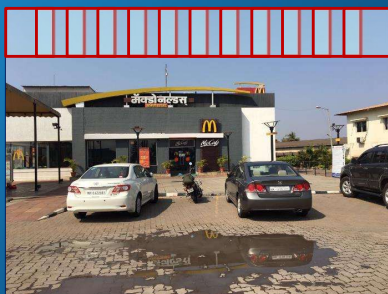Identify multiple object in the image

- Classes of all object
- Location of bounding box (mid point, height, width) of all objects
- $\hat{y}$ will be a vector

12/3/2024

pra-sami

## Sliding Window Detection

28



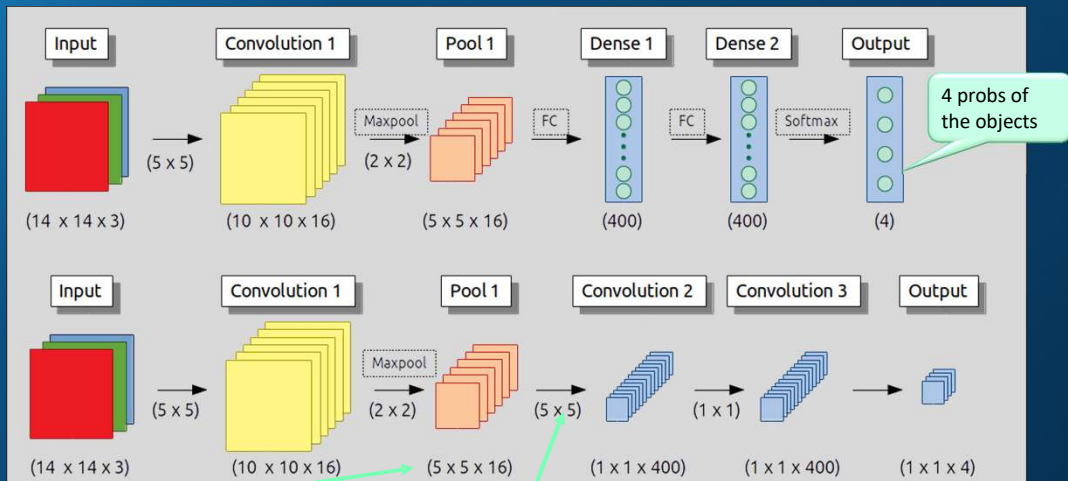❑ Analyzing for all these windows is resource consuming….

❑ We can convert logic to some what similar to convolutional networks and achieve better efficiencies.

12/3/2024

pra-sami

## Sliding Window Convolution way…

29



**Traditional ConvNet**

| Input | Convolution 1 | Pool 1 | Dense 1 | Dense 2 | Output |
|---|---|---|---|---|---|
| (14 x 14 x 3) | (5 x 5) → (10 x 10 x 16) | Maxpool (2 x 2) → (5 x 5 x 16) | FC → (400) | FC → (400) | Softmax → (4) |

4 probs of the objects

**Sliding window ConvNet**

| Input | Convolution 1 | Pool 1 | Convolution 2 | Convolution 3 | Output |
|---|---|---|---|---|---|
| (14 x 14 x 3) | (5 x 5) → (10 x 10 x 16) | Maxpool (2 x 2) → (5 x 5 x 16) | (5 x 5) → (1 x 1 x 400) | (1 x 1) → (1 x 1 x 400) | → (1 x 1 x 4) |

- ❑ Each 5 x5 x 16 layer is applied 5 x 5 x 16 filter and some activation to get 1 x 1 x 400 nodes
- ❑ Mathematically its same as fully connected layer!!

12/3/2024

pra-sami

## Convolution Implementation of Object Detection

30



- ❑ The computations are shared across the windows
- ❑ Results of each of region (1 x 1) are available using the convolution
- ❑ For bigger image size, output also increases
- ❑ This is telling us if in respective region, target object is present or not!

pra-sami

15

# Convolution instead of Sliding Window.

**31**



| Input | | | | | Output |
|---|---|---|---|---|---|
| (28 x 28 x 3) | (24 x 24 x 16) | (12 x 12 x 16) | (8 x 8 x 400) | (8 x 8 x 400) | (8 x 8 x 4) |

- Hence, by moving 14x14 region over the entire image we would know location of the region with maximum probability of containing a car.
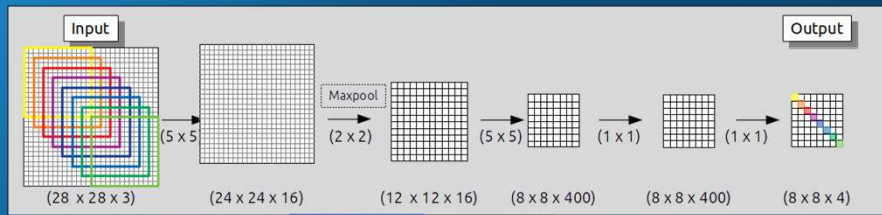
- Issue remains that size of bounding box ( region) is predefined

- Chances are that it is not very accurate.

12/3/2024

pra-sami

---

# Intersection over Union - IoU

**32**

- $IoU = \dfrac{Area\ of\ Ground\ Truth\ Box}{Area\ of\ Predicted\ Box}$



Ground Truth bounding box

Predicted

Actual

- IoU > 0.5 Acceptable

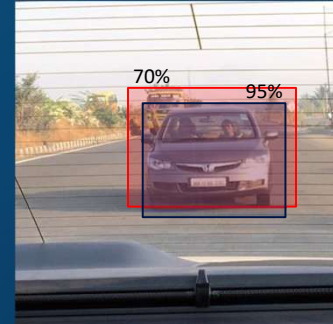- IoU = 1.0 Perfect

- IoU > 0.6 for little stringent requirements

12/3/2024

pra-sami

16

## Non Max Suppression



- Lets assume we are interested in only one object per image ➜ [$p_c$, x, y, h, w]
- First step will be to discard all detection below a certain threshold (e.g. $p_c \leq 0.70$)
- The output bounding boxes will have some overlap
- Retain one with highest probability
- If you are trying to identify multiple objects, say Cars, Pedestrians, Motorcycles output vector will have more dimensions
  - $p_c, C_1, C_2, C_3, x_1, y_1, h_1, w_1, \ldots$

12/3/2024

pra-sami

## Anchor Boxes

- Any anchor can be defined with
  - Presence : in any object is present in the anchor
  - Box location: mid point ( x, y ), height and width of the box
  - Class: What class is present- Car/person/motorcycle
- Fully defined anchor for three class
  - $p_c, b_x, b_y, b_h, b_w, c_1, c_2, c_3 \implies 8$ values

$$\hat{y} = \left\{ \begin{array}{c} Presence \\ Box\ location \\ Class \end{array} \right\} = \left\{ \begin{array}{c} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{array} \right\}$$



Car

Person

Person

12/3/2024

pra-sami

**35**

# YOLO – You Only Look Once - Training and Data Preparation

- ❑ Assume our image is divided in 3 x 3 grid
  - ❖ Real implementation : 16 x 16 or 19 x 19
- ❑ Assume we have only two anchor box per cell
  - ❖ i.e. not more than two items in a cell
- ❑ Thus ŷ will be 3 x 3 x 16

$$
\hat{y} = \begin{Bmatrix} p_{c1} \\ b_{x1} \\ b_{y1} \\ b_{h1} \\ b_{w1} \\ c_{11} \\ c_{21} \\ c_{31} \\ p_{c2} \\ b_{x2} \\ b_{y2} \\ b_{h2} \\ b_{w2} \\ c_{12} \\ c_{22} \\ c_{32} \end{Bmatrix} =
$$

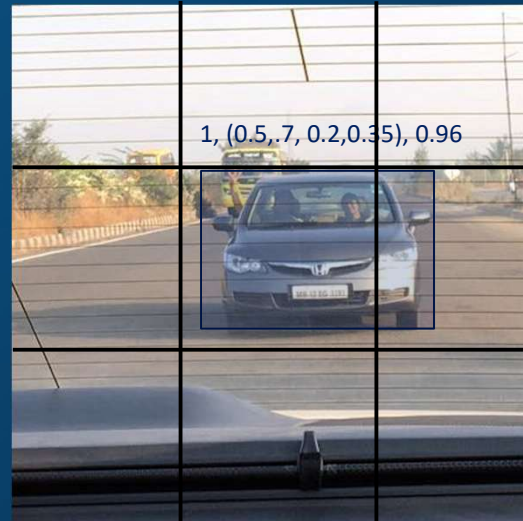| | | | | |
|---|---|---|---|---|
| 0 | 0 | 0 | | 0 |
| – | – | – | | – |
| – | – | – | | – |
| – | – | – | | – |
| – | – | – | | – |
| – | – | – | | – |
| – | – | – | | – |
| – | – | – | | – |
| 0 | 0 | 0 | ……… | 1 |
| – | – | – | | 0.5 |
| – | – | – | | 0.7 |
| – | – | – | | 0.2 |
| – | – | – | | 0.35 |
| – | – | – | | 1 |
| – | – | – | | 0 |
| – | – | – | | 0 |

1, (0.5,.7, 0.2,0.35), 0.96

12/3/2024

pra-sami

---

**36**

# YOLO – You Only Look Once - Predictions

- ❑ Thus ŷ will be 3 x 3 x 16

$$
\hat{y} = \begin{Bmatrix} p_{c1} \\ b_{x1} \\ b_{y1} \\ b_{h1} \\ b_{w1} \\ c_{11} \\ c_{21} \\ c_{31} \\ p_{c2} \\ b_{x2} \\ b_{y2} \\ b_{h2} \\ b_{w2} \\ c_{12} \\ c_{22} \\ c_{32} \end{Bmatrix} =
$$

| | | | | |
|---|---|---|---|---|
| 0 | 0 | 0 | | 0 |
| – | – | – | | – |
| – | – | – | | – |
| – | – | – | | – |
| – | – | – | | – |
| – | – | – | | – |
| – | – | – | | – |
| – | – | – | | – |
| 0 | 0 | 0 | ………… | 1 |
| – | – | – | | 0.5 |
| – | – | – | | 0.7 |
| – | – | – | | 0.2 |
| – | – | – | | 0.35 |
| – | – | – | | 1 |
| – | – | – | | 0 |
| – | – | – | | 0 |

12/3/2024

pra-sami

## YOLO – You Only Look Once - Predictions

37

- ❏ Get bounding boxes for each of the cells...

- ❏ Bounding boxes may overflow
  - ❖ We have not given any grid locations

- ❏ Except for those in red every one else would have low probability

- ❏ Keep Red ones and remove others.



12/3/2024

pra-sami

## YOLO5 – Most Stable Version (2023)

38

- ❏ Resize the input image to 488 x488

- ❏ Run a single convolutional network on the image

- ❏ Thresholds the resulting detections by the model's confidence

- ❏ Final output is the 7 × 7 ×30 tensor of predictions

- ❏ Leaky ReLU as activation in all the Layers (except last)

- ❏ Linear activation function for final layer

- ❏ Sum of Squares Error (SSE) as optimizing function

- ❏ Batch size of 64, Momentum of 0.9 and Decay of 0.0005

- ❏ Dropout (rate = .5) is used after the first connected layer

- ❏ Data Augmentation is used (random scaling, translation, exposure, saturation)

12/3/2024

pra-sami

## YOLO V1

39

pra-sami

---

## Now Darknet -53

40

❑ Starting YOLO version 3.0 started using Darknet-53

❖ Other networks can also be used

❑ 
$$loss_1 = -\sum_{i=0}^{S^2}\sum_{j=0}^{B} W_{ij}^{obj}[\hat{C}_i^j \log(C_i^j)+(1-\hat{C}_i^j)\log(1-C_i^j)] -$$
$$\lambda_{noobj}\sum_{i=0}^{S^2}\sum_{j=0}^{B}(1-W_{ij}^{obj})[\hat{C}_i^j \log(C_i^j)+(1-\hat{C}_i^j)\log(1-C_i^j)]$$

$$loss_2 = -\sum_{i}^{s^2}\sum_{j}^{B} W_{ij}^{obj}\sum_{c=1}^{C}[\hat{p}_i^j(c)\log(p_i^j(c))-(1-\hat{p}_i^j(c))\log(1-p_i^j(c))]$$

$$loss_3 = 1-IOU+\frac{\rho^2(b,b^{gt})}{c^2}+\frac{16}{\pi^4}\frac{\left(\arctan\frac{w^{gt}}{h^{gt}}-\arctan\frac{w}{h}\right)^4}{1-IOU+\frac{4}{\pi^2}\left(\arctan\frac{w^{gt}}{h^{gt}}-\arctan\frac{w}{h}\right)^2}$$

|    | Type | Filters | Size | Output |
|----|------|---------|------|--------|
|    | Convolutional | 32 | 3 × 3 | 256 × 256 |
|    | Convolutional | 64 | 3 × 3 / 2 | 128 × 128 |
|    | Convolutional | 32 | 1 × 1 | |
| 1× | Convolutional | 64 | 3 × 3 | |
|    | Residual | | | 128 × 128 |
|    | Convolutional | 128 | 3 × 3 / 2 | 64 × 64 |
|    | Convolutional | 64 | 1 × 1 | |
| 2× | Convolutional | 128 | 3 × 3 | |
|    | Residual | | | 64 × 64 |
|    | Convolutional | 256 | 3 × 3 / 2 | 32 × 32 |
|    | Convolutional | 128 | 1 × 1 | |
| 8× | Convolutional | 256 | 3 × 3 | |
|    | Residual | | | 32 × 32 |
|    | Convolutional | 512 | 3 × 3 / 2 | 16 × 16 |
|    | Convolutional | 256 | 1 × 1 | |
| 8× | Convolutional | 512 | 3 × 3 | |
|    | Residual | | | 16 × 16 |
|    | Convolutional | 1024 | 3 × 3 / 2 | 8 × 8 |
|    | Convolutional | 512 | 1 × 1 | |
| 4× | Convolutional | 1024 | 3 × 3 | |
|    | Residual | | | 8 × 8 |
|    | Avgpool | | Global | |
|    | Connected | | 1000 | |
|    | Softmax | | | |

Table 1. **Darknet-53.**

pra-sami

20

# R-CNN

41

❑ RCNN has nothing to do with RNN (Recurrent neural networks).

❑ R-CNN is short for "Region-based Convolutional Neural Networks."
   ❖ Takes in input image
   ❖ Extracts around 2000 bottom-up region proposals
   ❖ Computes features for each proposal using a large convolutional neural network (CNN)
   ❖ Classifies each region using class-specific linear SVMs

❑ This network was slow, hence
   ❖ Spate of other proposals are going on
   ❖ Fast RCNN
      ➢ Convolutional implementation of sliding window
   ❖ Faster R-CNN
      ➢ Use Convolutional Network to propose regions

pra-sami

---

42

Dense Net

pra-sami

43

12/3/2024

pra-sami

---

44

## A 5-layer dense block with a growth rate of k = 4.



DenseNets **simplify the connectivity pattern**

Concatenate all previous layer

But now the depths are exploding…

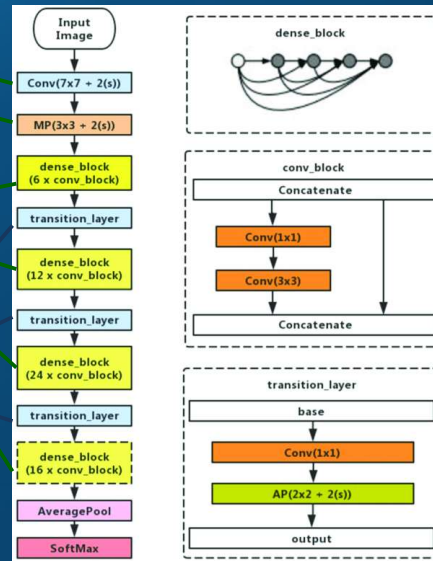Each layer is increasing the depth.

12/3/2024

pra-sami

## DenseNet 121 Architecture

45



Conv layer

MaxPool layer

Dense Blocks

Transition Layer

---

## DenseNet Architectures

46

| Layers | Output Size | DenseNet-121 | | DenseNet-169 | | DenseNet-201 | | DenseNet-264 | |
|---|---|---|---|---|---|---|---|---|---|
| Convolution | 112 × 112 | 7 × 7 conv, stride 2 | | | | | | | |
| Pooling | 56 × 56 | 3 × 3 max pool, stride 2 | | | | | | | |
| Dense Block (1) | 56 × 56 | $\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix}$ | × 6 | $\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix}$ | × 6 | $\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix}$ | × 6 | $\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix}$ | × 6 |
| Transition Layer (1) | 56 × 56 | 1 × 1 conv | | | | | | | |
| | 28 × 28 | 2 × 2 average pool, stride 2 | | | | | | | |
| Dense Block (2) | 28 × 28 | $\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix}$ | × 12 | $\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix}$ | × 12 | $\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix}$ | × 12 | $\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix}$ | × 12 |
| Transition Layer (2) | 28 × 28 | 1 × 1 conv | | | | | | | |
| | 14 × 14 | 2 × 2 average pool, stride 2 | | | | | | | |
| Dense Block (3) | 14 × 14 | $\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix}$ | × 24 | $\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix}$ | × 32 | $\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix}$ | × 48 | $\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix}$ | × 64 |
| Transition Layer (3) | 14 × 14 | 1 × 1 conv | | | | | | | |
| | 7 × 7 | 2 × 2 average pool, stride 2 | | | | | | | |
| Dense Block (4) | 7 × 7 | $\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix}$ | × 16 | $\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix}$ | × 32 | $\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix}$ | × 32 | $\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix}$ | × 48 |
| Classification Layer | 1 × 1 | 7 × 7 global average pool | | | | | | | |
| | | 1000D fully-connected, softmax | | | | | | | |

## Why Change?

47

❑ DenseNets require fewer parameters than an equivalent traditional CNN

❑ Some variations of ResNets have proven that many layers are barely contributing and can be dropped

❑ Inception Nets have proven that it's a good idea to concatenate layers

❑ Vanishing Gradients were always problems
  ❖ In DenseNets each layer has direct access to the gradients from the loss function and the original input image

❑ Traditional feed-forward neural networks connect the output of the layer to the next layer using:
  ❖ Activations ($a^{[l]}$) = g ($a^{[l-1]} * W^{[l]} + b^{[l]}$)

❑ ResNet modified them a bit:
  ❖ Activations ($a^{[l]}$) = g ($a^{[l-1]} * W^{[l]} + b^{[l]} + a^{[l-2]}$)

pra-sami

## DenseNets

48

❑ DenseNets : do not sum the output feature maps of the layer with the incoming feature maps but concatenate them:
  ❖ Activations ($a^{[l]}$) = g ( [$a^{[0]}$ , $a^{[1]}$ , $a^{[2]}$ , …, $a^{[l-2]}$ , $a^{[l-1]} * W^{[l]}$ ] + $b^{[l]}$)

❑ But Activations between various layers would have different shape
  ❖ To solve, DenseNets divide them in blocks
  ❖ Shape remain same in one DenseBlock

❑ Transition Layers: Layers in-between Dense Layers changing dimensions from one block to another block:
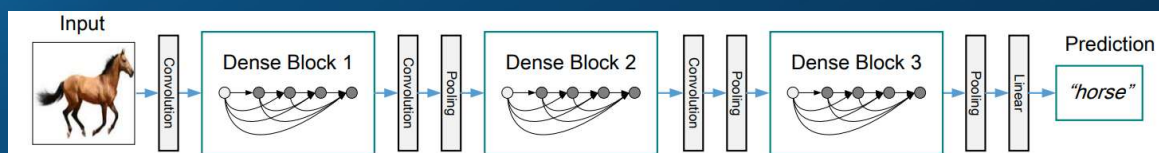  ❖ Apply 1 x 1, pooling, BatchNorm etc.

pra-sami

## DenseNets

❑ Every layer has access to its preceding feature maps
  ❖ i.e. to the collective knowledge
  ❖ Each layer is then adding a new information

❑ DenseNet layers are very narrow (e.g., 12 filters per layer)
  ❖ Adding only a small set of feature-maps to the "collective knowledge" of the network
  ❖ Keep the remaining feature-maps unchanged
  ❖ The final classifier makes a decision based on all feature-maps in the network

## Type of DenseNets

❑ DenseNets-B
  ❖ Regular DenseNets that take advantage of 1x1 convolution to reduce the feature maps size
  ❖ Then apply the 3x3 convolution
  ❖ B stands for bottleneck

❑ DenseNets-BC
  ❖ Another little incremental step to DenseNets-B, to reduce the number of output feature maps
  ❖ The compression factor (theta) determines the reduction.
  ❖ Instead of having m feature maps at a certain layer, we will have theta*m.
  ❖ Theta is in the range [0–1].
  ❖ DenseNets will remain the same when theta=1, and will be DenseNets-B otherwise.

## Reflect…

51

- ❑ Which of the following is true about AlexNet?
  - ❖ a) It uses 15 layers including fully connected layers
  - ❖ b) It introduced the concept of Residual Learning
  - ❖ c) It was the first CNN to win the ImageNet Large Scale Visual Recognition Challenge (ILSVRC)
  - ❖ d) It uses a 5x5 kernel in the first convolutional layer

- ❑ Answer: c) It was the first CNN to win the ImageNet Large Scale Visual Recognition Challenge (ILSVRC)

- ❑ What is the key innovation introduced by ResNet?
  - ❖ a) Use of deeper convolution layers
  - ❖ b) Use of 1x1 convolution kernels
  - ❖ c) Introduction of skip connections (residual connections)
  - ❖ d) Global average pooling for dimensionality reduction

- ❑ Answer: c) Introduction of skip connections (residual connections)

- ❑ Which of the following is true about ImageNet?
  - ❖ a) It is a dataset consisting of 10 million images
  - ❖ b) It contains over 22,000 object categories
  - ❖ c) It focuses on medical image segmentation
  - ❖ d) It contains only grayscale images

- ❑ Answer: b) It contains over 22,000 object categories

- ❑ What is the primary characteristic of VGGNet architecture?
  - ❖ a) It uses a large number of filters in each layer
  - ❖ b) It uses very small 3x3 filters in convolutional layers
  - ❖ c) It introduced skip connections
  - ❖ d) It employs global average pooling instead of fully connected layers

- ❑ Answer: b) It uses very small 3x3 filters in convolutional layers

pra-sami

## Reflect…

52

- ❑ What was the main innovation introduced by Google's Inception Net?
  - ❖ a) Introduction of the "bottleneck" layers
  - ❖ b) Use of parallel filters of different sizes in the same layer (Inception module)
  - ❖ c) Use of large convolution filters for all layers
  - ❖ d) Introduction of Dense blocks

- ❑ Answer: b) Use of parallel filters of different sizes in the same layer (Inception module)

- ❑ What is the key innovation of Faster R-CNN over Fast R-CNN?
  - ❖ a) It uses an RPN (Region Proposal Network) for faster region proposals
  - ❖ b) It replaces convolution layers with fully connected layers
  - ❖ c) It combines object detection and segmentation in one model
  - ❖ d) It removes the need for bounding box regression

- ❑ Answer: a) It uses an RPN (Region Proposal Network) for faster region proposals

- ❑ How does YOLO differ from traditional object detection models?
  - ❖ a) YOLO performs object detection by scanning the image in patches
  - ❖ b) YOLO predicts both class probabilities and bounding boxes in a single pass
  - ❖ c) YOLO uses a sliding window technique for localization
  - ❖ d) YOLO uses fully connected layers for region proposal

- ❑ Answer: b) YOLO predicts both class probabilities and bounding boxes in a single pass

- ❑ What is the primary characteristic of DenseNet?
  - ❖ a) It uses dilated convolutions to increase the receptive field
  - ❖ b) It uses skip connections from every layer to every other layer
  - ❖ c) It stacks convolutional layers without any pooling layers
  - ❖ d) It uses separable convolutions to reduce computational cost

- ❑ Answer: b) It uses skip connections from every layer to every other layer

pra-sami

## Reflect…

53

- Why does ResNet's performance degrade when the depth of the network increases, without residual connections?
  - a) The network begins to overfit due to an excessive number of parameters
  - b) The gradient vanishes as it backpropagates through the layers, making training ineffective
  - c) It reduces computational complexity too much, leading to poor feature extraction
  - d) It uses too many skip connections, leading to exploding gradients

- Answer: b) The gradient vanishes as it backpropagates through the layers, making training ineffective

- In DenseNet, how does feature reuse occur across layers?
  - a) Each layer receives the feature maps of all preceding layers as input
  - b) Feature maps from selected layers are concatenated to form the final feature vector
  - c) The output of each layer is summed with the output of the previous layer
  - d) DenseNet shares weights between alternate layers to reduce the number of parameters

- Answer: a) Each layer receives the feature maps of all preceding layers as input

- In Faster R-CNN, what is the role of the Region Proposal Network (RPN)?
  - a) To classify the entire image and then crop regions of interest
  - b) To predict regions that are most likely to contain objects, which are then classified by the detection network
  - c) To directly classify each pixel of the image into object categories
  - d) To generate bounding boxes based on edge detection algorithms

- Answer: b) To predict regions that are most likely to contain objects, which are then classified by the detection network

- Which domain is U-Net primarily designed for?
  - a) Object detection
  - b) Natural language processing
  - c) Image segmentation, especially in biomedical images
  - d) Image classification

- Answer: c) Image segmentation, especially in biomedical images

12/3/2024

pra-sami

54



12/3/2024

pra-sami

EXTRA MATERIAL

**pra-sâmi**

---

## Tips

| Data vs. Feature Engineering | Benchmark Performance |
|---|---|
| ❑ Depending upon size of data, you may need to do feature engineering | ❑ For benchmarking ➜ Ensamble |
| | ❖ Create multiple model ( 3 to 25 models) |
| | ❖ Train them independently |
| ❑ More data, lesser feature engineering | ❖ Average out the results (ŷ) |
| | ❑ Rarely used in production due to cost considerations |
| | ❑ Multi-crop at the test time |

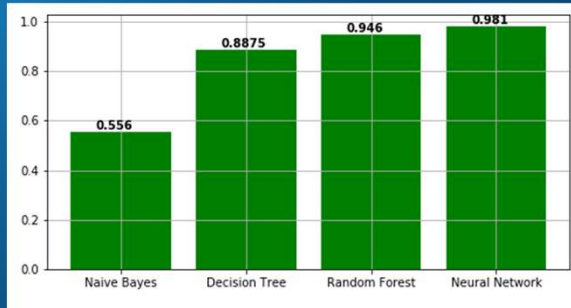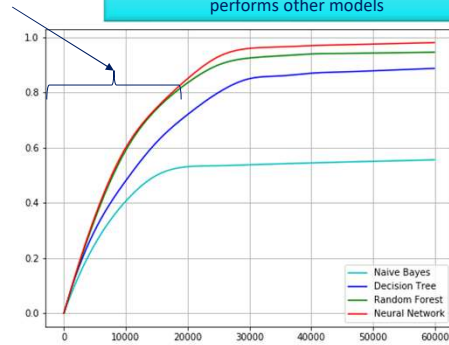**pra-sâmi**

# Relative performance of models

57

Small amount of data performance are comparable

As data size grows Neural networks out performs other models





12/3/2024

pra-sami