

Deep Learning

Saideep

June 2019

1 Introduction

This is a tutorial on the basics of Deep Learning based on the Deep Learning Specialization on Coursera. Let us begin with a classic classification problem and keep building on it. Let us also get familiar with the notations that we are going to use throughout.

Classification Problem:

Take the example of image classification where you have to identify the presence of cat in an image. The dimensions of the image being considered are: 64x64x3 (3 channels for R, B, G).



Figure 1: Sample Image

$$n_x = 64 \times 64 \times 3 = 12288$$

(x,y) where $x \in R^{n_x}$ & $y \in \{0,1\}$

m training examples $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), (x^{(3)}, y^{(3)}) \dots (x^{(m)}, y^{(m)})\}$

$m_{train} = \# \text{ of training examples}$; $m_{test} = \# \text{ of test examples}$

$$X = \begin{bmatrix} \vdots & \vdots & \vdots & \dots & \vdots \\ x^{(1)} & x^{(1)} & x^{(1)} & \dots & x^{(m)} \\ \vdots & \vdots & \vdots & \dots & \vdots \end{bmatrix} \quad X.\text{shape} = (n_x, m)$$

$$Y = [y^{(1)}, y^{(2)}, y^{(3)} \dots y^{(m)}] \quad Y.\text{shape} = (1, m)$$

2 Logistic Regression

$$\hat{y} = P(y = 1/x); \quad 0 \leq \hat{y} \leq 1 \quad x \in R^{n_x}$$

parameters: $w \in R^{n_x}$, $b \in R$
 Output = $\sigma(w^T x + b)$

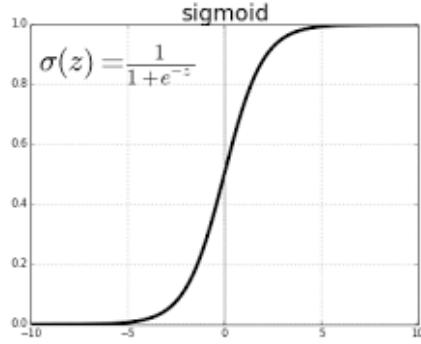


Figure 2: Sigmoid Function

2.1 Cost function:

$$L(\hat{y}, y) = 1/2(\hat{y} - y)^2 \Rightarrow \text{non-convex}$$

$$\hat{y} = \sigma(w^T x + b), \sigma(z) = 1/(1 + \exp(-z))$$

$$\text{Hence, } L(\hat{y}, y) = -(y \log(\hat{y}) + (1 - y) \log(1 - \hat{y}))$$

Cost function over entire training set:

$$\begin{aligned} J(w, b) &= 1/m (\sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})) \\ &= -1/m \sum_{i=1}^m (y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})) \end{aligned}$$

2.2 Gradient Descent:

Repeat {

$$w := w - \alpha dJ(w)/dw$$

}

Gradient Descent on one training example:

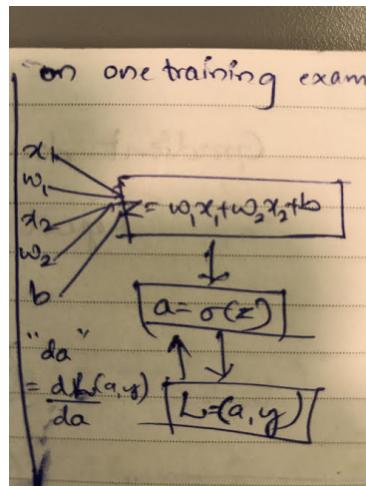


Figure 3:

$$\begin{aligned}
da &= -y/a + (1-y)/a \\
dz &= dL/dz = (dL/da) * (da/dz) = a - y \\
dL/dw_1 &= dw_1 = x_1 dz \\
dL/dw_2 &= dw_2 = x_2 dz \\
db &= dz \\
w_1 &:= w_1 - \alpha dw_1; w_2 := w_2 - \alpha dw_2; b := b - \alpha db
\end{aligned}$$

Gradient Descent on m training examples:

$$\begin{aligned}
J(w, b) &= 1/m \sum_{i=1}^m L(a^{(i)}, y) \\
a^{(i)} &= \widehat{y^{(i)}} = \sigma(z^{(i)}) = \sigma(w^T x^{(i)} + b)
\end{aligned}$$

```

for i = 1 to m do
    zi = wTx(i) + b
    a(i) = σ(z(i))
    J += -[y(i)log(y(i)) + (1 - y(i))log(1 - y(i))]
    dz(i) = a(i) - y(i)
    dw(1) += x1(i)dz(i)
    dw(2) += x2(i)dz(i)
    db+ = dz(i)
end for
J/ = m dw1/ = m dw2/ = m
dw1 = dJ/dw1
w1 := w1 - αdw1; w2 := w2 - αdw2; b := b - αdb
    ▷ n=2

```

2.3 Vectorizing Logistic Regression:

$$\begin{aligned}
z^{(1)} &= w^T x^{(1)} + b & z^{(2)} &= w^T x^{(2)} + b & z^{(3)} &= w^T x^{(3)} + b \\
a^{(1)} &= \sigma(z^{(1)}) & a^{(2)} &= \sigma(z^{(2)}) & a^{(3)} &= \sigma(z^{(3)}) \\
X &= \begin{bmatrix} \vdots & \vdots & \vdots & \dots & \vdots \\ x^{(1)} & x^{(1)} & x^{(1)} & \dots & x^{(m)} \\ \vdots & \vdots & \vdots & \dots & \vdots \end{bmatrix} & X.shape &= (n_x, m) \\
Z &= [z^{(1)}, z^{(2)}, z^{(3)} \dots z^{(m)}] & Z.shape &= (1, m) \\
&= w^T X + [b, b, b, \dots b]_{1xm} \\
&= w^T \begin{bmatrix} \vdots & \vdots & \vdots & \dots & \vdots \\ x^{(1)} & x^{(1)} & x^{(1)} & \dots & x^{(m)} \\ \vdots & \vdots & \vdots & \dots & \vdots \end{bmatrix} + b \\
&= \begin{bmatrix} \vdots & \vdots & \vdots & \dots & \vdots \\ w^T x^{(1)} + b & w^T x^{(1)} + b & w^T x^{(1)} + b & \dots & w^T x^{(m)} + b \\ \vdots & \vdots & \vdots & \dots & \vdots \end{bmatrix}
\end{aligned}$$

In python,

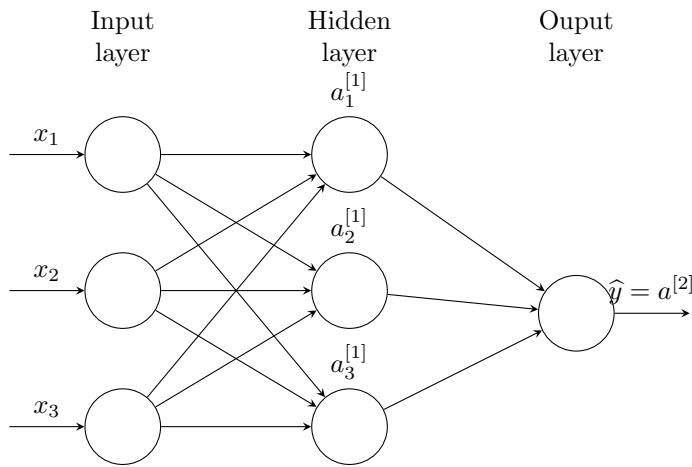
$$Z = np.dot(w.T, X) + b_{(1,1)}$$

$$\begin{aligned}
dz^{(1)} &= a^{(1)} - y^{(1)} & dz^{(2)} &= a^{(2)} - y^{(2)} \\
dZ &= [dz^{(1)}, dz^{(2)}, dz^{(3)} \dots dz^{(m)}]_{(1,m)} \\
A &= [a^{(1)}, a^{(2)}, a^{(3)} \dots a^{(m)}] & A.shape &= (1, m) \\
Y &= [y^{(1)}, y^{(2)}, y^{(3)} \dots y^{(m)}] & Y.shape &= (1, m)
\end{aligned}$$

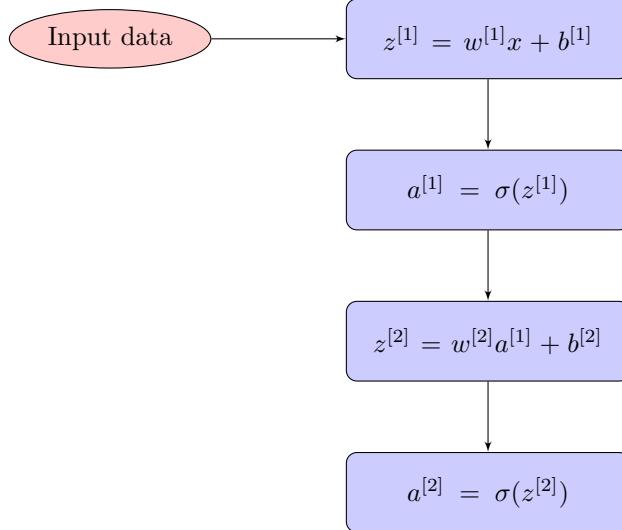
So,
 $Z = w^T X + b = np.dot(w.T, X) + b$
 $A = \sigma(Z)$
 $dZ = A - Y$
 $dw = (X dZ^T)/m$
 $db = np.sum(dZ)/m$
 $w := w - \alpha dw$ $b := b - \alpha db$

3 Neural Networks:

Using these ideas, let us now build a simple neural network and learn some notations as well.



superscript $[1],[2]$ to denote the layer of input & subscript $_{1,2}$ to denote the neuron number in the layer.
By convention we do not count the input layer for reporting the number of layers in the neural network.



$$L(a^{[2]}, y)$$

3.1 For one training example:

$$\begin{aligned} z^{[1]} &= w_1^{[1]T}x + b_1^{[1]}, \quad a_1^{[1]} = \sigma(z_1^{[1]}) \\ z^{[2]} &= w_2^{[1]T}x + b_2^{[1]}, \quad a_2^{[1]} = \sigma(z_2^{[1]}) \\ z^{[3]} &= w_3^{[1]T}x + b_3^{[1]}, \quad a_3^{[1]} = \sigma(z_3^{[1]}) \\ z^{[4]} &= w_4^{[1]T}x + b_4^{[1]}, \quad a_4^{[1]} = \sigma(z_4^{[1]}) \end{aligned}$$

$$\begin{aligned} &\begin{bmatrix} \dots & w_1^{[1]T} & \dots \\ \dots & w_2^{[1]T} & \dots \\ \dots & w_3^{[1]T} & \dots \\ \dots & w_4^{[1]T} & \dots \end{bmatrix}_{4 \times 3} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}_{3 \times 1} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix}_{4 \times 1} \\ &= \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}_{4 \times 1} = z^{[1]} \\ a^{[1]} &= \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \\ a_4^{[1]} \end{bmatrix}_{4 \times 1} = \sigma(z^{[1]}) \end{aligned}$$

Given input x,

$$\begin{aligned} Z^{[1]} &= W_{4 \times 3}^{[1]}x + b_{3 \times 1}^{[1]} \\ a^{[1]} &= \sigma(z^{[1]}) \\ Z^{[2]} &= W_{1 \times 4}^{[2]}a^{[1]} + b_{1 \times 1}^{[2]} \\ a^{[2]} &= \sigma(z^{[2]}) \end{aligned}$$

3.2 Vectorizing across multiple examples:

We denote $x^{(i)}$ for i^{th} training example,

3.2.1 Non-vectorized :

```
for i = 1 to m do
    z^{[1](i)} = w^{[1]}x^{(i)} + b^{[1]}
    a^{[1](i)} = σ(z^{[1](i)})
    z^{[2](i)} = w^{[2]}a^{[1](i)} + b^{[2]}
    a^{[2](i)} = σ(z^{[2](i)})
end for
```

3.2.2 Vectorized :

$$\begin{aligned} Z^{[1]} &= W^{[1]}X + b^{[1]} \\ A^{[1]} &= \sigma(Z^{[1]}) \\ Z^{[2]} &= W^{[2]}A^{[1]} + b^{[2]} \\ A^{[2]} &= \sigma(Z^{[2]}) \\ X &= \begin{bmatrix} \vdots & \vdots & \vdots & \dots & \vdots \\ x^{(1)} & x^{(1)} & x^{(1)} & \dots & x^{(m)} \\ \vdots & \vdots & \vdots & \dots & \vdots \end{bmatrix} \quad X.\text{shape} = (n_x, m) \end{aligned}$$

$$\begin{array}{c}
 \text{Training examples} \\
 \xrightarrow{\hspace{1cm}} \\
 Z^{[1]} = \begin{pmatrix} \vdots & \vdots & \cdots & \vdots \\ z^{1} & z^{[1](2)} & \cdots & z^{[1](m)} \\ \vdots & \vdots & \cdots & \vdots \end{pmatrix} \\
 A = [a^{1}, a^{[1](2)}, \dots, a^{[1](m)}] \\
 \downarrow \text{Hidden Units}
 \end{array}$$

3.3 Activation Functions :

tanh

$$g(z^{[1]}) = \tanh(z^{[1]}); \tanh(z) = (e^z - e^{-z})/(e^z + e^{-z})$$

ReLU : Rectified Linear Unit

$$g(z) = \max(0, z)$$

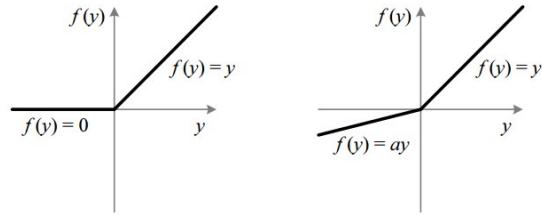


Figure 4: ReLU and Leaky ReLU

3.4 Back Propagation

3.4.1 Forward Propagation

$$Z^{[1]} = W^{[1]}X + b^{[1]}$$

$$A^{[1]} = \sigma(Z^{[1]})$$

$$Z^{[2]} = W^{[2]}A^{[1]} + b^{[2]}$$

$$A^{[2]} = \sigma(Z^{[2]})$$

3.4.2 Back Propagation

$$dZ^{[2]} = A^{[2]} - Y$$

$$dW^{[2]} = dZ^{[2]}A^{[1]T}/m$$

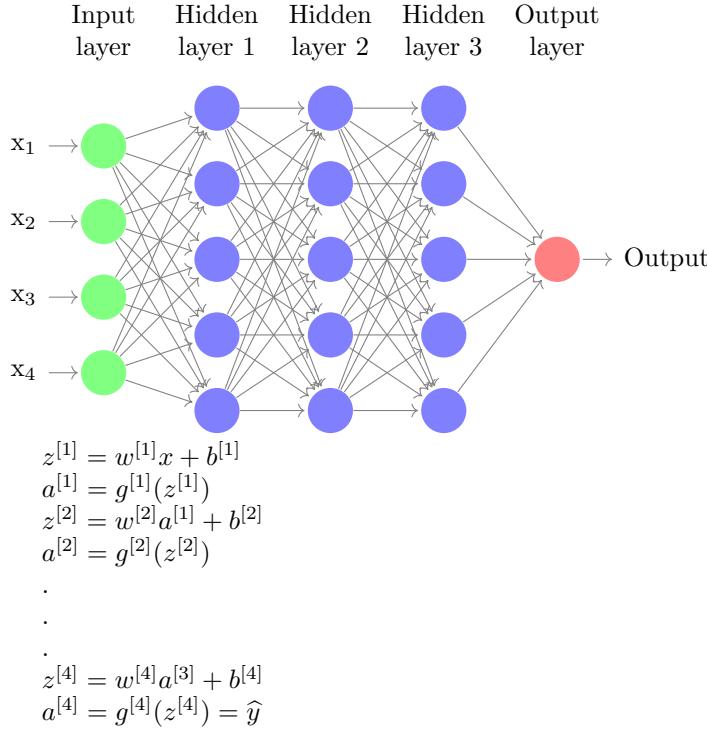
$$db^{[2]} = np.sum(dZ^{[2]}A^{[1]T}, axis=1, keepdims=True)/m$$

$$dZ^{[1]} = W^{[2]T}dZ^{[2]} \quad \overbrace{*}^{elem-wise-product} \quad g^{[1]}'(Z^{[1]})$$

$$dW^{[1]} = dZ^{[1]}X^T/m$$

$$db^{[1]} = np.sum(dZ^{[1]}, axis=1, keepdims=True)/m$$

4 Deep Neural Networks



To generalize,
 $z^{[l]} = w^{[l]}a^{[l-1]} + b^{[l]}$
 $a^{[l]} = g^{[l]}(z^{[l]})$

4.1 Forward Propagation for all training examples:

```

for i = 1 to 4 do
     $Z^{[i]} = W^{[i]}A^{[i-1]} + b^{[i]}$   $\triangleright X = A^{[0]}$ 
     $A^{[i]} = g^{[i]}(Z^{[i]})$   $\triangleright A^{[4]} = \hat{Y} = g^{[4]}(Z^{[4]})$ 
end for

```

Dimensionality check,

For single training example -

$$Z_{(n^{[l]}, 1)}^{[l]} = W_{(n^{[l]}, n^{[l-1]})}^{[l]} A_{(n^{[l-1]}, 1)}^{[l-1]} + b_{(n^{[l]}, 1)}^{[l]} \quad (W \text{ and } dW \text{ will have the same dimensions})$$

For all training examples -

$$A^{[l]} \& Z^{[l]} - (n^{[l]}, m)$$

$$W^{[l]} - (n^{[l]}, n^{[l+1]})$$

$$b^{[l]} - (n^{[l]}, 1) \quad (\text{python broadcasts hence it will add to each column when added to matrix})$$

4.2 Backward Propagation

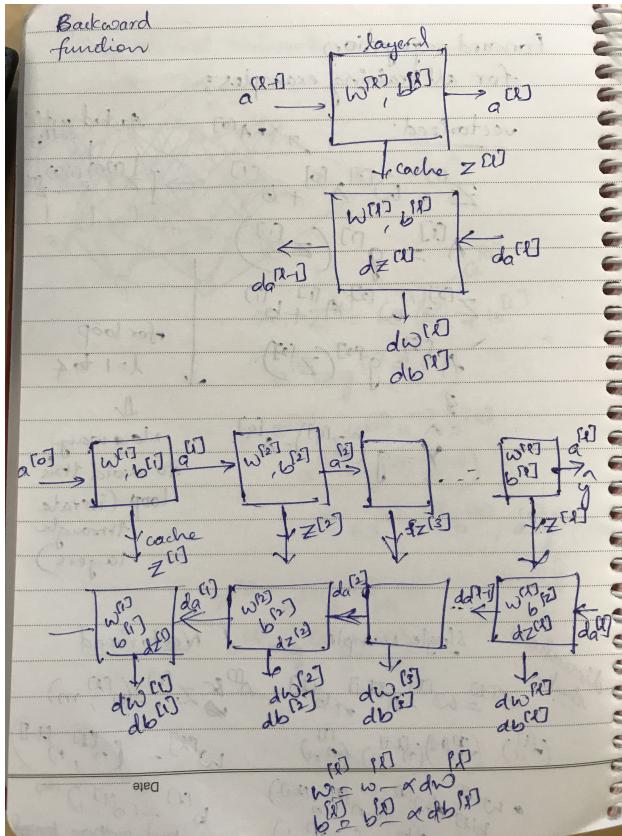


Figure 5:

4.3 Backward Propagation for layer l:

Input : $da^{[l]}$ Output : $da^{[l-1]}, dw^{[l]}, db^{[l]}$

$$dz^{[l]} = da^{[l-1]} * g^{[l]}'(z^{[l]})$$

$$dw^{[l]} = dz^{[l]} a^{[l-1]}$$

$$db^{[l]} = dz^{[l]}$$

$$da^{[l-1]} = W^{[l]T} dz^{[l]}$$

$$da^{[l]} = (-y)/a + (1 - y)/(1 - a)$$

4.4 Vectorized version for layer 1 (all training examples):

$$dZ^{[l]} = dA^{[l]} * g^{[l]}'(Z^{[l]})$$

$$dW^{[l]} = dZ^{[l]} A^{[l-1]T}$$

$$db^{[l]} = np.sum(dZ^{[l]}, axis=1, keepdims=True)/m$$

$$dA^{[l-1]} = W^{[l]T} dZ^{[l]}$$

$$dA^{[l]} = \left(\frac{-y^{(1)}}{a^{(1)}} + \frac{1-y^{(1)}}{1-a^{(1)}}, \frac{-y^{(2)}}{a^{(2)}} + \frac{1-y^{(2)}}{1-a^{(2)}}, \dots, \frac{-y^{(m)}}{a^{(m)}} + \frac{1-y^{(m)}}{1-a^{(m)}} \right)$$

5 Convolutional Neural Networks

5.1 Applications in Computer Vision:

- * Classification
- * Position Identification
- * Edge detection

Convolution operator:

$$\begin{array}{|c|c|c|} \hline 3 & 0 & 1 \\ \hline 1 & 5 & 8 \\ \hline 2 & 7 & 2 \\ \hline \end{array} \quad * \quad \begin{array}{|c|c|c|} \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline \end{array}$$

$$= 3*1 + 1*1 + 2*1 + 0*0 + 5*0 + 1*0 + 1*-1 + 8*-1 + 2*-1 = -5$$

Vertical Edge detection:

$$\begin{array}{|c|c|c|} \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline \end{array}$$

Horizontal Edge detection:

$$\begin{array}{|c|c|c|} \hline 1 & 1 & 1 \\ \hline 0 & 0 & 0 \\ \hline -1 & -1 & -1 \\ \hline \end{array}$$

5.2 CNNs:

In Deep Learning we try to learn the weights for the filter/kernel:

$$\begin{array}{|c|c|c|} \hline w_1 & w_2 & w_3 \\ \hline w_4 & w_5 & w_6 \\ \hline w_7 & w_8 & w_9 \\ \hline \end{array}$$

If we convolve a $n*n$ image with $f*f$ filter the resulting image has the dimensions $(n-f+1)*(n-f+1)$, i.e. (i) it shrinks the output as well as (ii) some pixels are used less than the others when we convolve (edges).

Padding

When we pad the image by p , image size would be $(n+2p)$ hence output of the convolved image is $(n+2p-f+1)$

Valid Convolution:

- * No padding
- $(nxn)*(fxf) \rightarrow (n-f+1)x(n-f+1)$

Same Convolution:

*Pad so that the output size is same as the input size.

$$\Rightarrow n+2p-f+1 = n$$

$$\Rightarrow p = (f-1)/2$$

f is almost always chosen to be odd probably because (i) asymmetry in padding when the filter size is even (ii) positioning of filter (it's nice to have a central pixel)

Strided Convolution:

Steps of convolution,

stride(s) = 2 means jump the window/filter by 2 steps

$$\left\lfloor \frac{n+2p-f}{s} + 1 \right\rfloor \times \left\lfloor \frac{n+2p-f}{s} + 1 \right\rfloor$$

Convolutions on RGB images

$$n \times n \times n_c * f \times f \times n_c \rightarrow (n-f+1) \times (n-f+1) \quad n_c = \text{depth}/\# \text{ of channels}$$

5.3 One CNN layer:

$d^{[l]}$ = filter size

$p^{[l]}$ = padding

$s^{[l]}$ = stride

$n_c^{[l]}$ = number of filters

Input : $n_H^{[l-1]} \times n_W^{[l-1]} \times n_c^{[l-1]}$

Output : $n_H^{[l]} \times n_W^{[l]} \times n_c^{[l]}$

$$n_{H/W}^{[l]} = \left\lfloor \frac{n_{H/W}^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1 \right\rfloor$$

Each filter is : $f^{[l]} \times f^{[l]} \times n_c^{[l-1]}$

Activations :

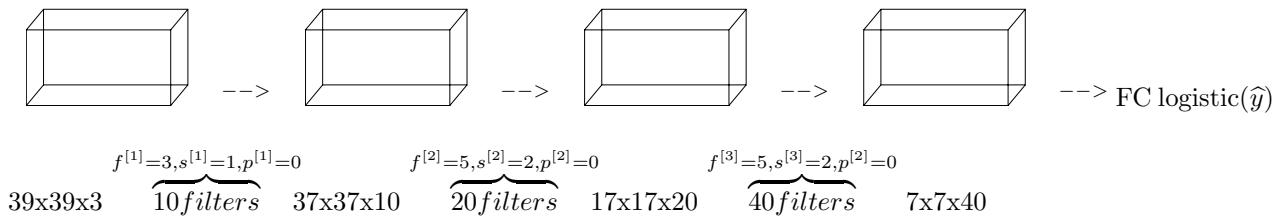
$a^{[l]} \rightarrow n_H^{[l]} \times n_W^{[l]} \times n_c^{[l]}$

$A^{[l]} \rightarrow m \times n_H^{[l]} \times n_W^{[l]} \times n_c^{[l]}$

Weights : $f^{[l]} \times f^{[l]} \times n_c^{[l-1]} \times n_c^{[l]} \rightarrow \# \text{ of filters in layer 1}$

bias : $n_c^{[l]} \rightarrow (1, 1, 1, n_c^{[l]})$

5.4 Typical CNN example



5.4.1 Types of layer in CNN

Convolution (CONV)

Pooling (POOL)

Fully Connected (FC)

Pooling Layer : Max Pooling

1	3	2	1
2	9	1	1
1	3	2	3
5	6	1	2

$\overbrace{\quad\quad\quad}^{f=2, s=2} \overbrace{\quad\quad\quad}^{-->}$

9	2
6	3

Case Studies:

Classic Networks:

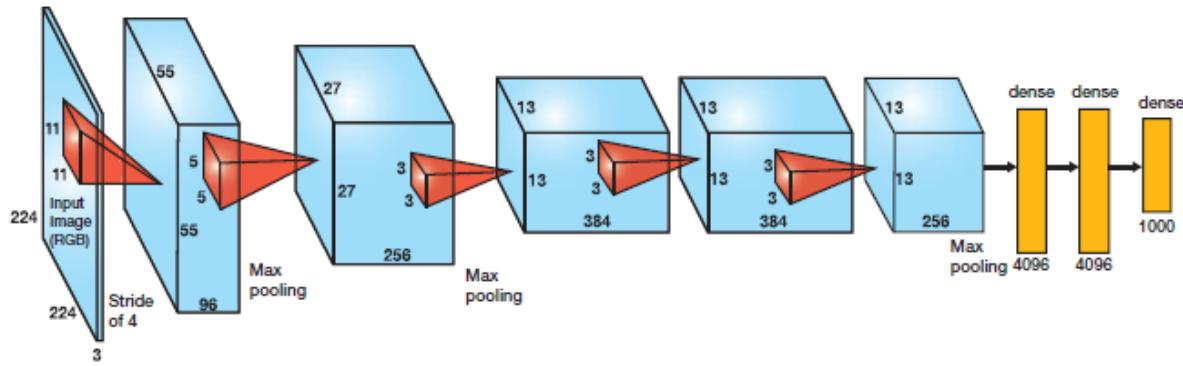
LeNet - 5

AlexNet

VGG

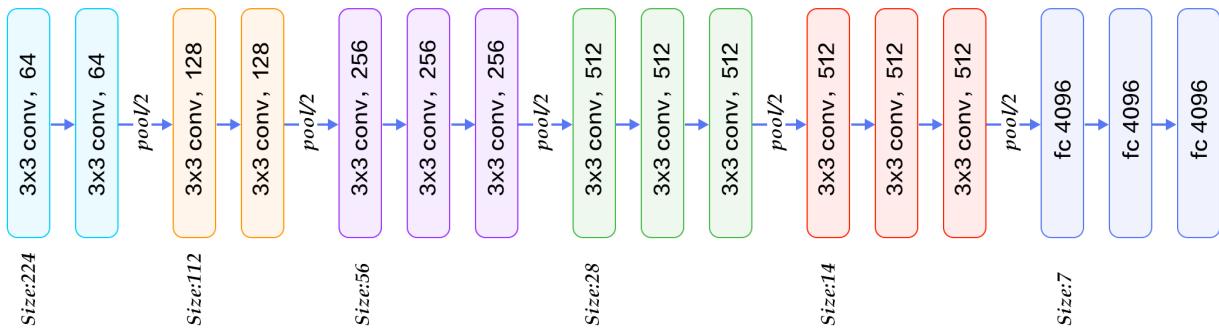
ResNet (152)

AlexNet



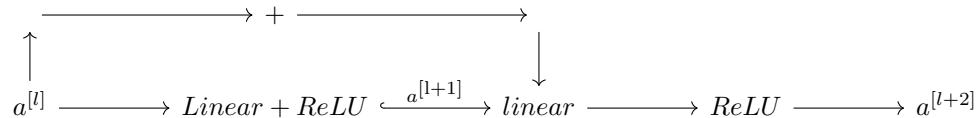
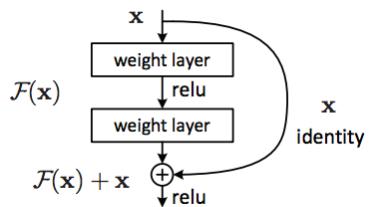
VGG

Here, CONV = 3x3 filter, s=1, same & max pool = 2x2, s=2



ResNet

Residual Networks



$$z^{[l+1]} = w^{[l+1]} a^{[l]} + b^{[l+1]} \quad a^{[l+1]} = g(z^{[l+1]})$$

$$z^{[l+2]} = w^{[l+2]} a^{[l+1]} + b^{[l+2]} \quad a^{[l+2]} = g(z^{[l+2]} + a^{[l]})$$

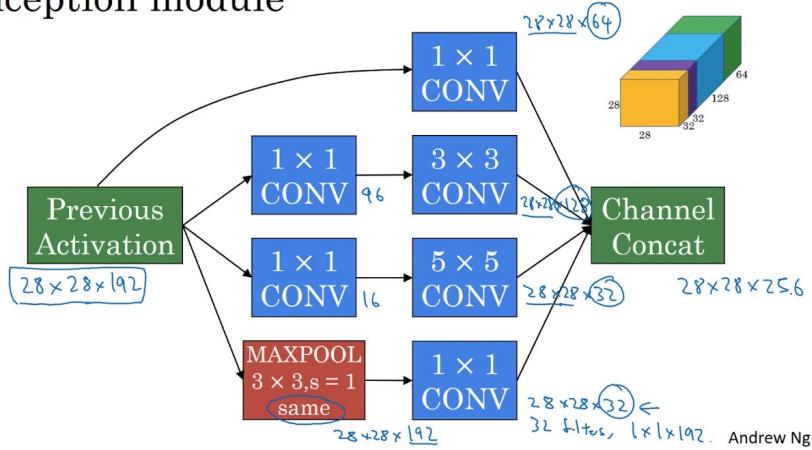
Why ResNets work?

As we go deeper & deeper it becomes difficult to learn parameters but it's easy to learn identity when residual block is added.

Inception network

Instead of choosing one filter size, use all sizes in the network.

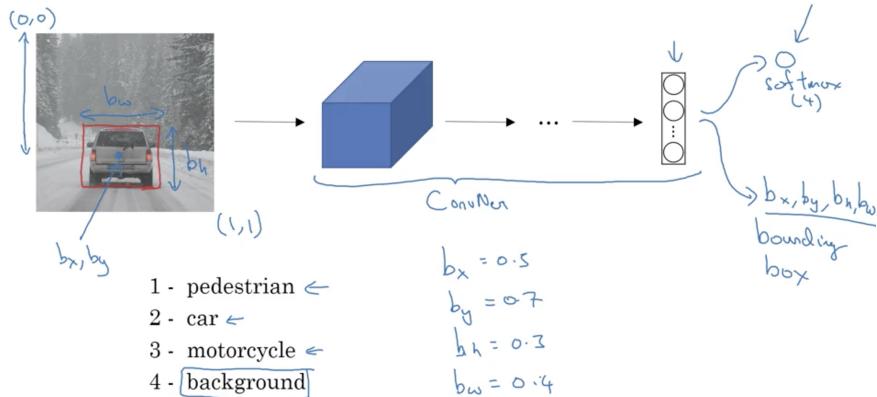
Inception module



6 Detection Algorithms:

6.1 Object Localization

In Object localization, instead we try to predict the x, y co-ordinates of the bounding box as well as the relative width and height of the box.



$$\text{Target label} = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

- $p_c > 0$ {background}
- 1 {car, pedestrian or motor-cycle}
- $c_1 >$ pedestrian/not
- $c_2 >$ car/not
- $c_3 >$ motor cycle/not

Loss function:

$$L(y, \hat{y}) = (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 + \dots + (\hat{y}_8 - y_8)^2$$

For landmark detection, one needs to annotate all the key points on the face.

Sliding window implementation - high computation cost

Convolutional implementation of sliding windows

YOLO (You Only Look Once) algorithm : You apply the localization algorithm to each of the grids.

Anchor boxes for multiple objects, where each object in the training image is assigned to grid cell that contains the object's mid-point & anchor box for the grid cell with highest IoU.