

# Machine Learning

## DSECL ZG565

Dr. Chetana Gavankar, Ph.D,  
IIT Bombay-Monash University Australia  
[Chetana.gavankar@pilani.bits-pilani.ac.in](mailto:Chetana.gavankar@pilani.bits-pilani.ac.in)



**BITS** Pilani

Pilani Campus



## **Lecture No. – 10 | Convolutional Neural Network**

**Date – 18/01/2020**

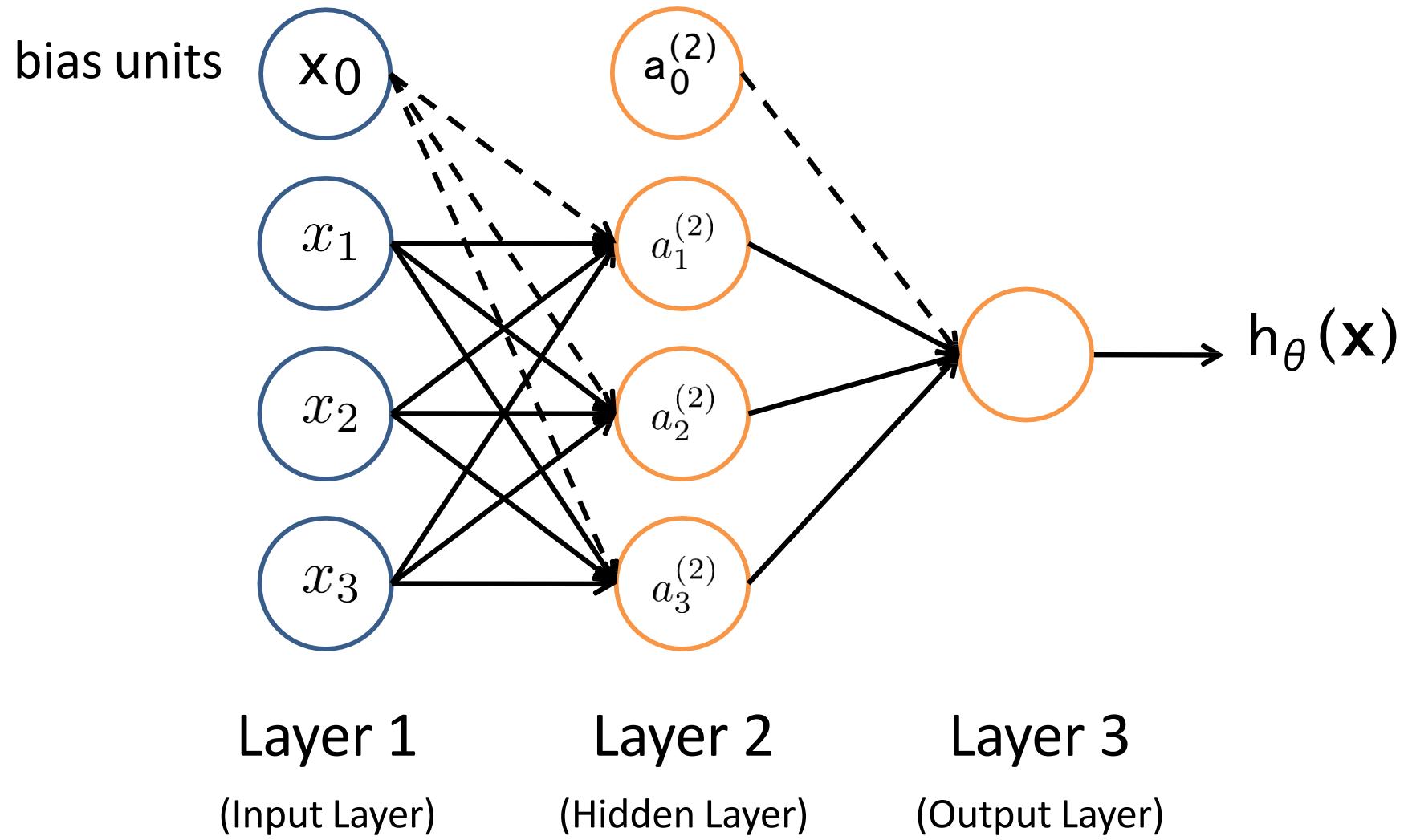
**Time – 9:00 AM – 11:00 AM**

These slides are prepared by the instructor, with grateful acknowledgement of Tom Mitchell, Mitesh Khapra and many others who made their course materials freely available online.

# Session Content

- Back propagation Algorithm (Andrew Ng Notes and Tom Mitchell chapter 4)
  - Training Neural Network
  - Hidden Layer Representations
  - Overfitting in ANN
  - Early Stopping
  - Drop out
  - Parameter sharing
- Convolutional Neural Network

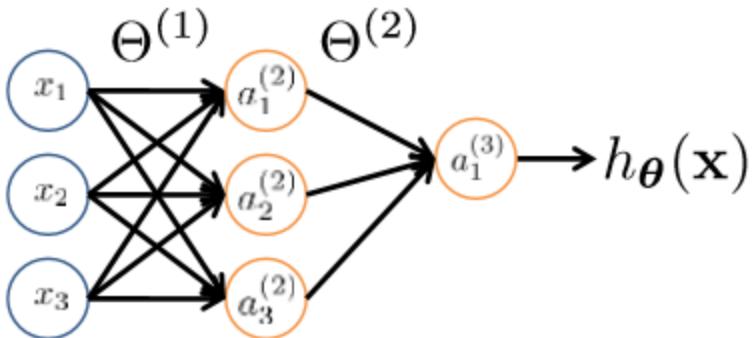
# Neural Network



# Feed-Forward Process

- Input layer units are set by some exterior function (think of these as **sensors**), which causes their output links to be **activated** at the specified level
- Working forward through the network, the **input function** of each unit is applied to compute the input value
  - Usually this is just the weighted sum of the activation on the links feeding into this node
- The **activation function** transforms this input function into a final value
  - Typically this is a **nonlinear** function, often a **sigmoid** function corresponding to the “threshold” of that node

# Neural Network



$a_i^{(j)}$  = “activation” of unit  $i$  in layer  $j$   
 $\Theta^{(j)}$  = weight matrix controlling function mapping from layer  $j$  to layer  $j + 1$

$$a_1^{(2)} = g(\Theta_{10}^{(1)}x_0 + \Theta_{11}^{(1)}x_1 + \Theta_{12}^{(1)}x_2 + \Theta_{13}^{(1)}x_3)$$

$$a_2^{(2)} = g(\Theta_{20}^{(1)}x_0 + \Theta_{21}^{(1)}x_1 + \Theta_{22}^{(1)}x_2 + \Theta_{23}^{(1)}x_3)$$

$$a_3^{(2)} = g(\Theta_{30}^{(1)}x_0 + \Theta_{31}^{(1)}x_1 + \Theta_{32}^{(1)}x_2 + \Theta_{33}^{(1)}x_3)$$

$$h_\Theta(x) = a_1^{(3)} = g(\Theta_{10}^{(2)}a_0^{(2)} + \Theta_{11}^{(2)}a_1^{(2)} + \Theta_{12}^{(2)}a_2^{(2)} + \Theta_{13}^{(2)}a_3^{(2)})$$

If network has  $s_j$  units in layer  $j$  and  $s_{j+1}$  units in layer  $j+1$ ,  
then  $\Theta^{(j)}$  has dimension  $s_{j+1} \times (s_j + 1)$ .

$$\Theta^{(1)} \in \mathbb{R}^{3 \times 4} \quad \Theta^{(2)} \in \mathbb{R}^{1 \times 4}$$

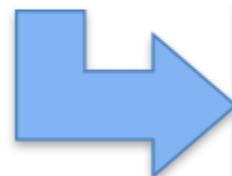
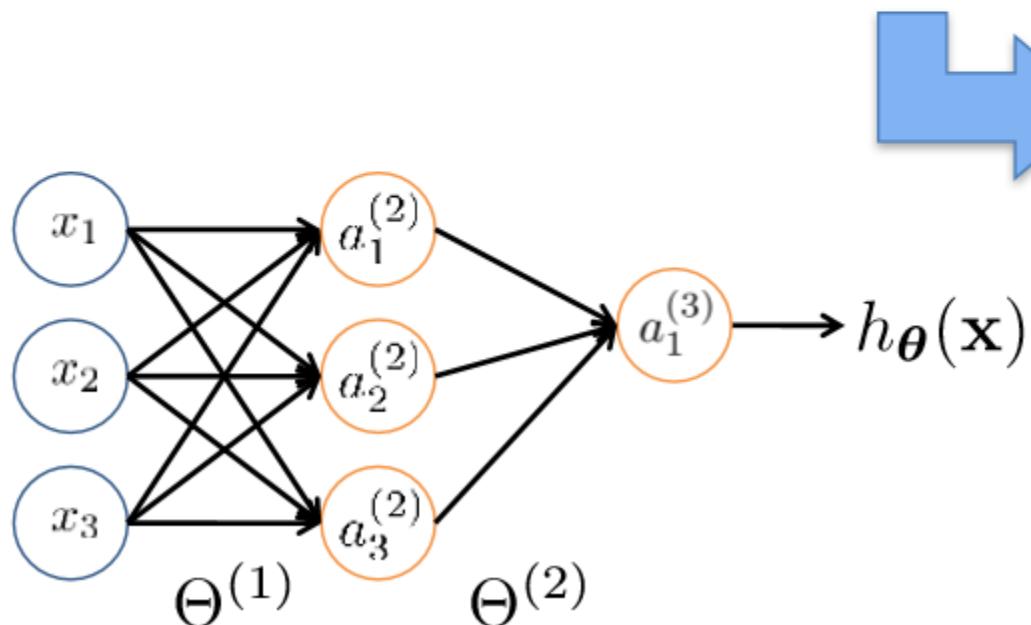
# Vectorization

$$a_1^{(2)} = g \left( \Theta_{10}^{(1)} x_0 + \Theta_{11}^{(1)} x_1 + \Theta_{12}^{(1)} x_2 + \Theta_{13}^{(1)} x_3 \right) = g \left( z_1^{(2)} \right)$$

$$a_2^{(2)} = g \left( \Theta_{20}^{(1)} x_0 + \Theta_{21}^{(1)} x_1 + \Theta_{22}^{(1)} x_2 + \Theta_{23}^{(1)} x_3 \right) = g \left( z_2^{(2)} \right)$$

$$a_3^{(2)} = g \left( \Theta_{30}^{(1)} x_0 + \Theta_{31}^{(1)} x_1 + \Theta_{32}^{(1)} x_2 + \Theta_{33}^{(1)} x_3 \right) = g \left( z_3^{(2)} \right)$$

$$h_{\Theta}(\mathbf{x}) = g \left( \Theta_{10}^{(2)} a_0^{(2)} + \Theta_{11}^{(2)} a_1^{(2)} + \Theta_{12}^{(2)} a_2^{(2)} + \Theta_{13}^{(2)} a_3^{(2)} \right) = g \left( z_1^{(3)} \right)$$



Feed-Forward Steps:

$$\mathbf{z}^{(2)} = \Theta^{(1)} \mathbf{x}$$

$$\mathbf{a}^{(2)} = g(\mathbf{z}^{(2)})$$

$$\text{Add } a_0^{(2)} = 1$$

$$\mathbf{z}^{(3)} = \Theta^{(2)} \mathbf{a}^{(2)}$$

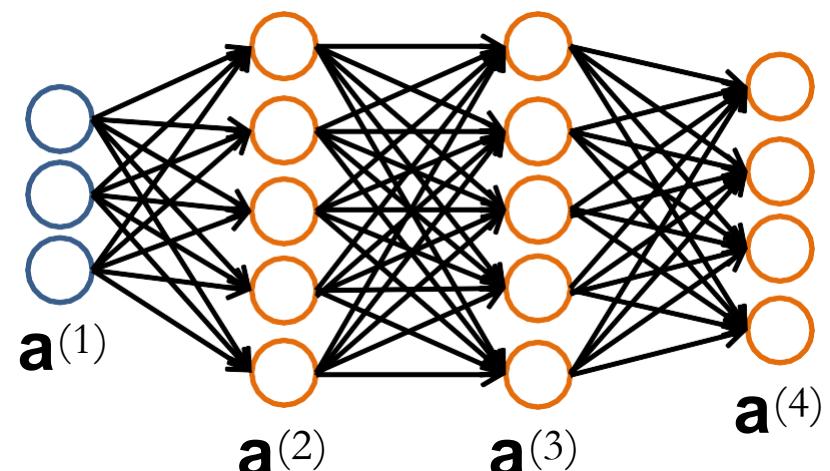
$$h_{\Theta}(\mathbf{x}) = \mathbf{a}^{(3)} = g(\mathbf{z}^{(3)})$$

# Forward Propagation

- Given one labeled training instance  $(\mathbf{x}, y)$ :

## Forward Propagation

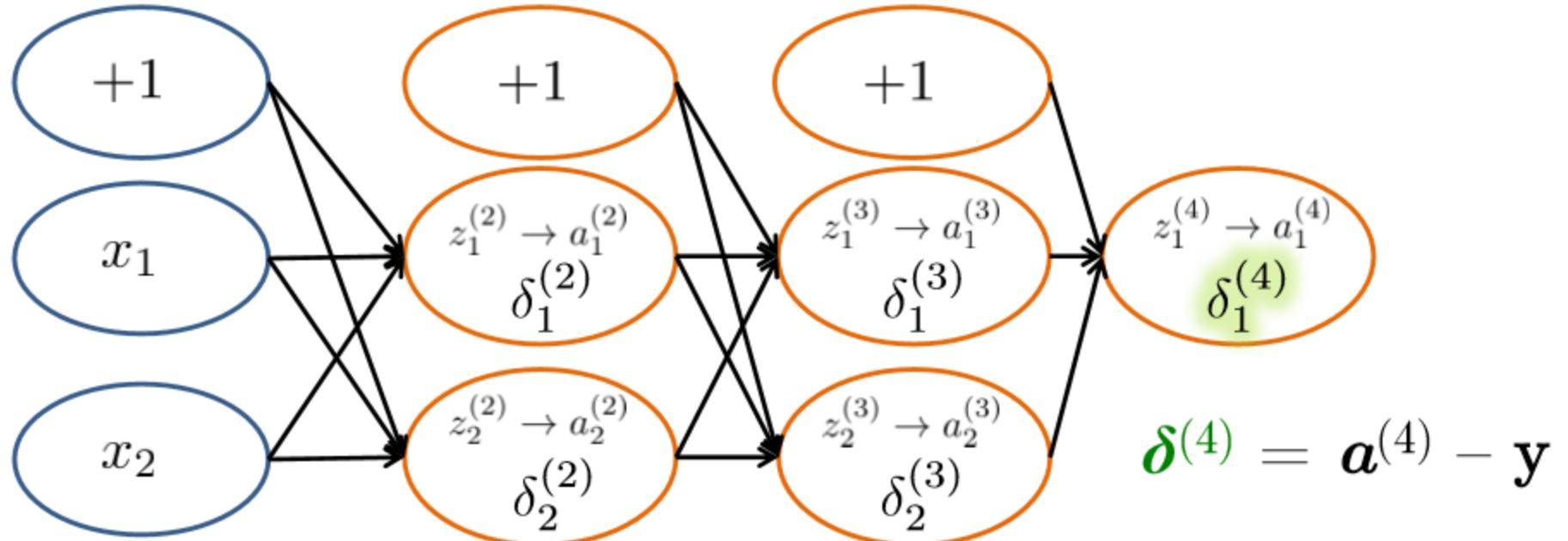
- $\mathbf{a}^{(1)} = \mathbf{x}$
- $\mathbf{z}^{(2)} = \Theta^{(1)}\mathbf{a}^{(1)}$
- $\mathbf{a}^{(2)} = g(\mathbf{z}^{(2)})$  [add  $a_0^{(2)}$ ]
- $\mathbf{z}^{(3)} = \Theta^{(2)}\mathbf{a}^{(2)}$
- $\mathbf{a}^{(3)} = g(\mathbf{z}^{(3)})$  [add  $a_0^{(3)}$ ]
- $\mathbf{z}^{(4)} = \Theta^{(3)}\mathbf{a}^{(3)}$
- $\mathbf{a}^{(4)} = h_{\Theta}(\mathbf{x}) = g(\mathbf{z}^{(4)})$



# Backpropagation Intuition

- Each hidden node  $j$  is “responsible” for some fraction of the error  $\delta_j^{(l)}$  in each of the output nodes to which it connects
- $\delta_j^{(l)}$  is divided according to the strength of the connection between hidden node and the output node
- Then, the “blame” is propagated back to provide the error values for the hidden layer

# Backpropagation Intuition

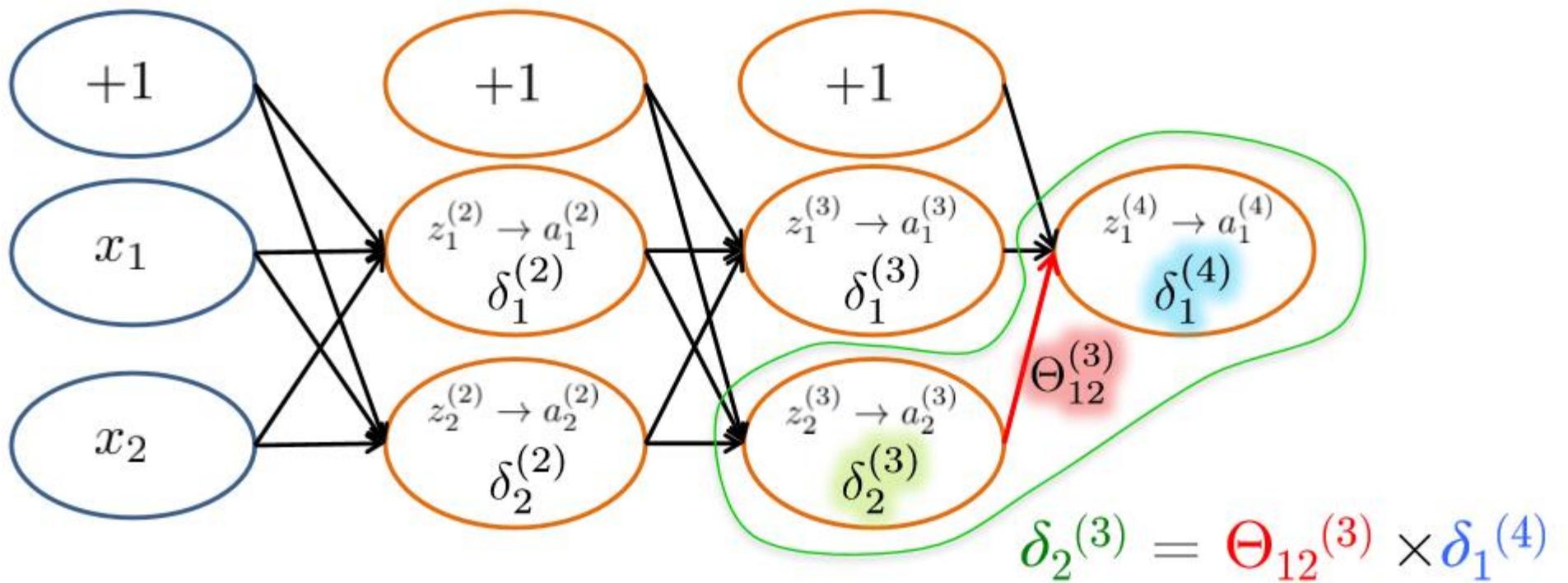


$\delta_j^{(l)}$  = “error” of node  $j$  in layer  $l$

Formally,  $\delta_j^{(l)} = \frac{\partial}{\partial z_j^{(l)}} \text{cost}(\mathbf{x}_i)$

where  $\text{cost}(\mathbf{x}_i) = y_i \log h_\Theta(\mathbf{x}_i) + (1 - y_i) \log(1 - h_\Theta(\mathbf{x}_i))$

# Backpropagation Intuition

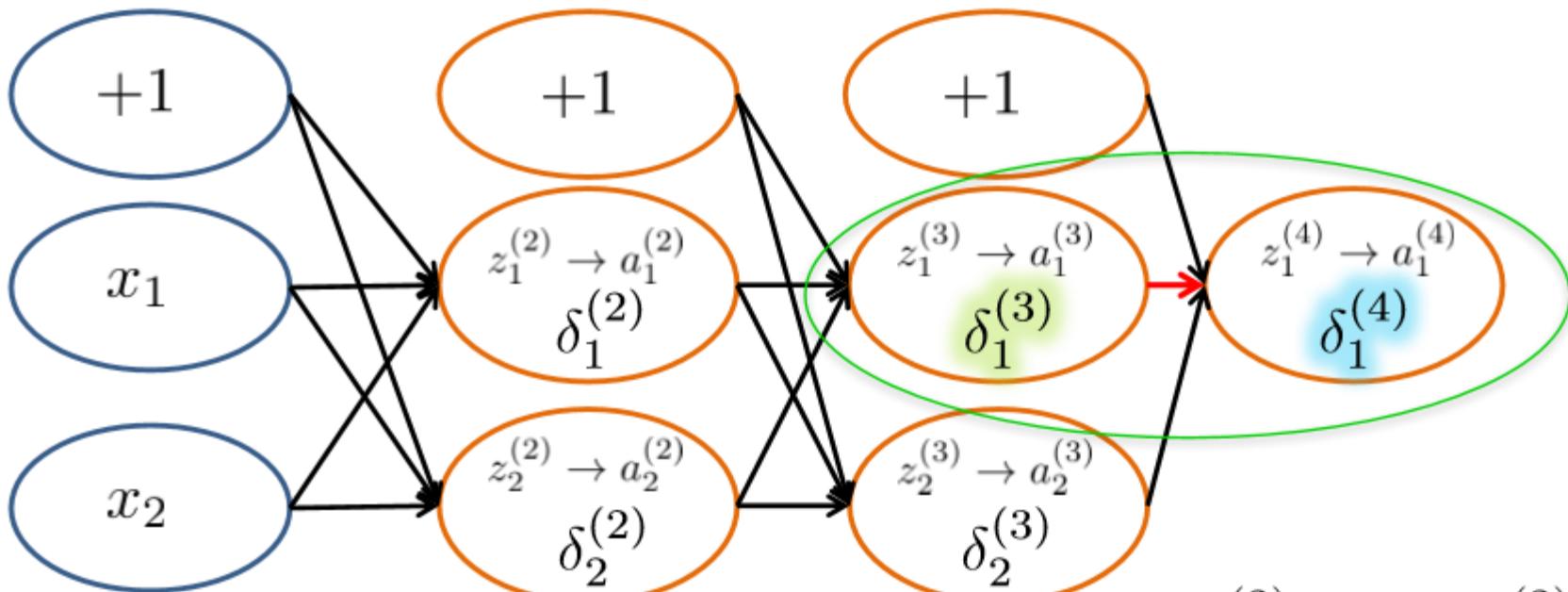


$\delta_j^{(l)}$  = “error” of node  $j$  in layer  $l$

Formally,  $\delta_j^{(l)} = \frac{\partial}{\partial z_j^{(l)}} \text{cost}(\mathbf{x}_i)$

where  $\text{cost}(\mathbf{x}_i) = y_i \log h_\Theta(\mathbf{x}_i) + (1 - y_i) \log(1 - h_\Theta(\mathbf{x}_i))$

# Backpropagation Intuition



$$\delta_2^{(3)} = \Theta_{12}^{(3)} \times \delta_1^{(4)}$$

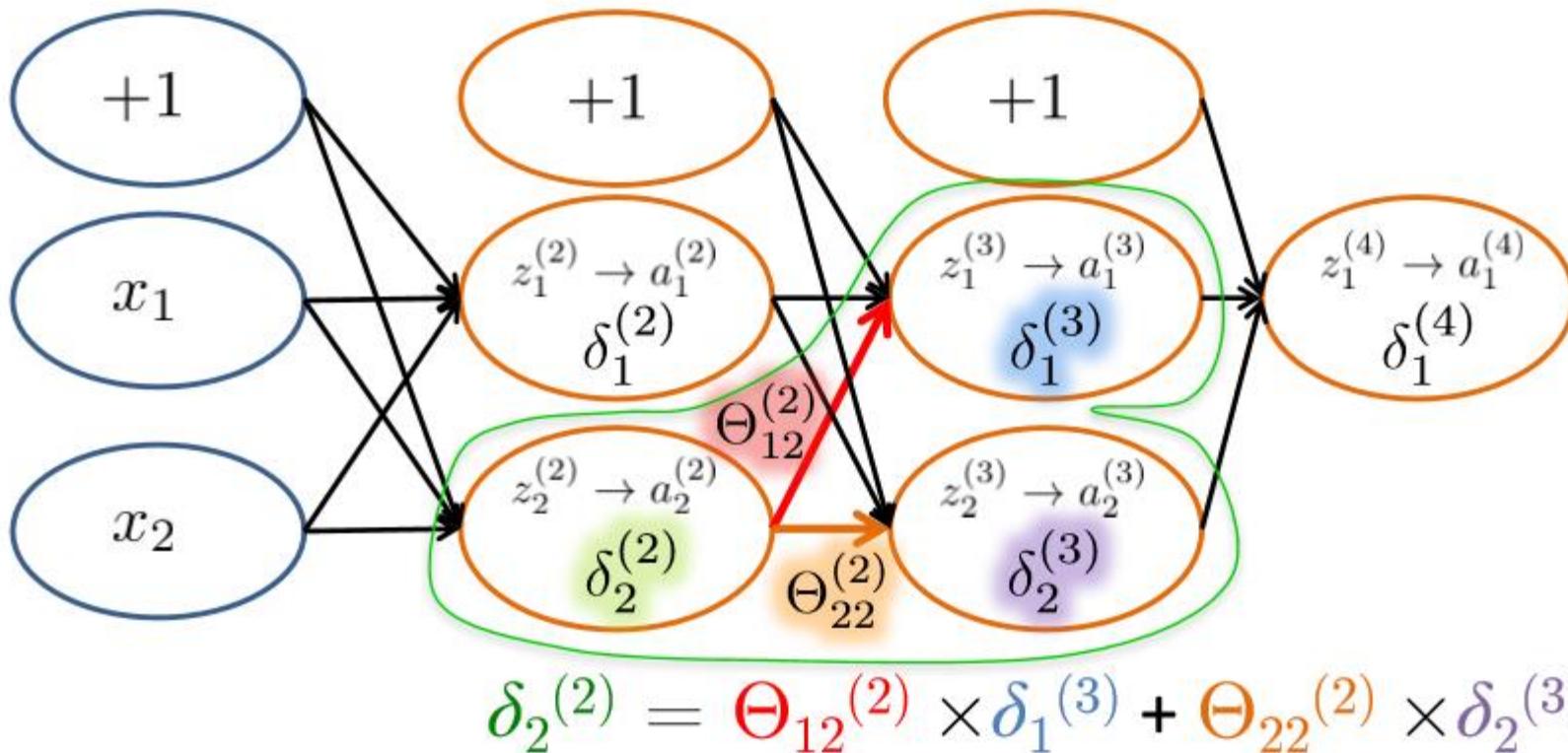
$$\delta_1^{(3)} = \Theta_{11}^{(3)} \times \delta_1^{(4)}$$

$\delta_j^{(l)}$  = “error” of node  $j$  in layer  $l$

Formally,  $\delta_j^{(l)} = \frac{\partial}{\partial z_j^{(l)}} \text{cost}(\mathbf{x}_i)$

where  $\text{cost}(\mathbf{x}_i) = y_i \log h_\Theta(\mathbf{x}_i) + (1 - y_i) \log(1 - h_\Theta(\mathbf{x}_i))$

# Backpropagation Intuition



$\delta_j^{(l)}$  = “error” of node  $j$  in layer  $l$

Formally,  $\delta_j^{(l)} = \frac{\partial}{\partial z_j^{(l)}} \text{cost}(\mathbf{x}_i)$

where  $\text{cost}(\mathbf{x}_i) = y_i \log h_\Theta(\mathbf{x}_i) + (1 - y_i) \log(1 - h_\Theta(\mathbf{x}_i))$

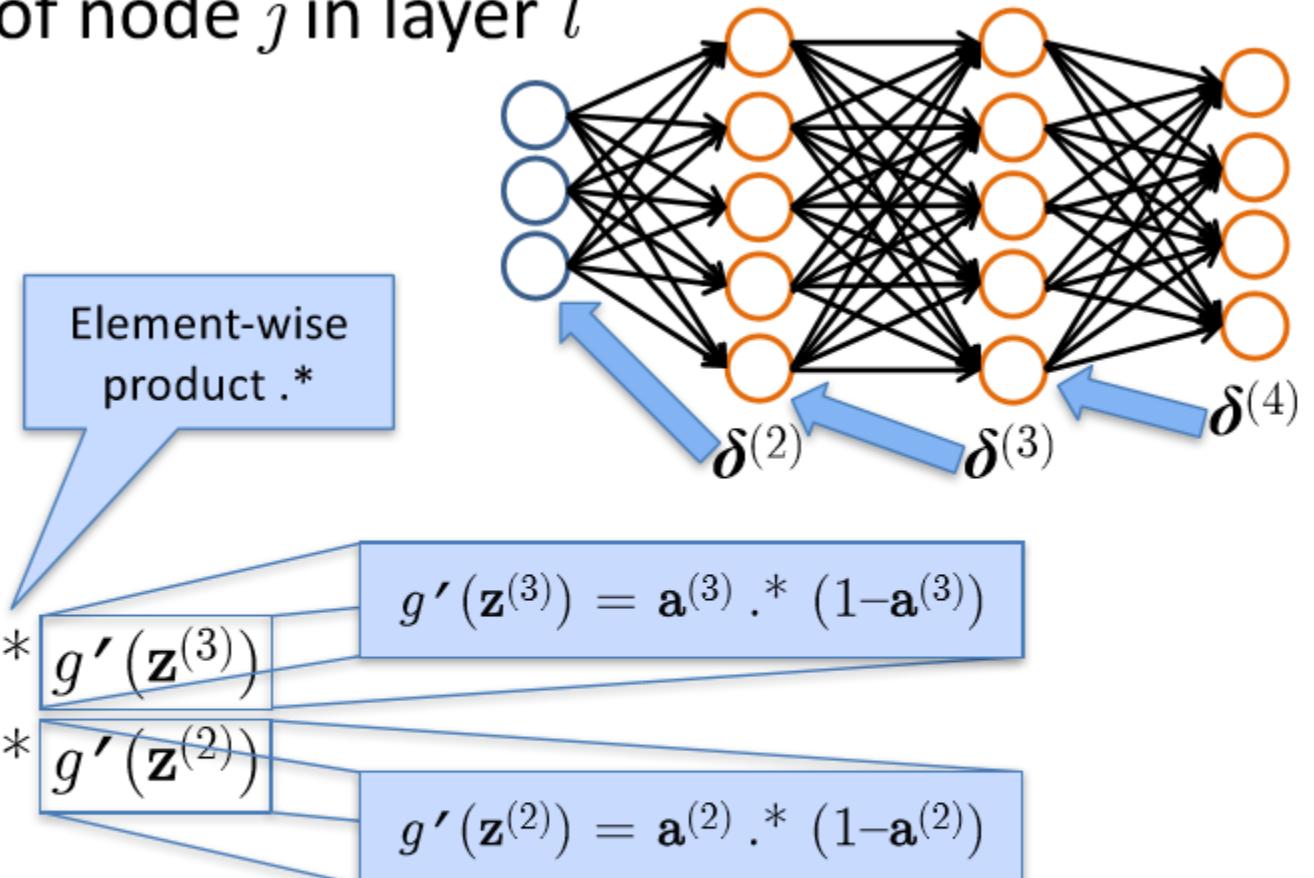
# Backpropagation: Gradient Computation

Let  $\delta_j^{(l)}$  = “error” of node  $j$  in layer  $l$

(#layers  $L = 4$ )

## Backpropagation

- $\delta^{(4)} = a^{(4)} - y$
- $\delta^{(3)} = (\Theta^{(3)})^\top \delta^{(4)} .*$
- $\delta^{(2)} = (\Theta^{(2)})^\top \delta^{(3)} .*$
- (No  $\delta^{(1)}$ )



$$\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = a_j^{(l)} \delta_i^{(l+1)} \quad (\text{ignoring } \lambda; \text{ if } \lambda = 0)$$

# Backpropagation

Given: training set  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$

Initialize all  $\Theta^{(l)}$  randomly (NOT to 0!)

Loop // each iteration is called an epoch

Set  $\Delta_{ij}^{(l)} = 0 \quad \forall l, i, j$  (Used to accumulate gradient)

For each training instance  $(\mathbf{x}_i, y_i)$ :

Set  $\mathbf{a}^{(1)} = \mathbf{x}_i$

Compute  $\{\mathbf{a}^{(2)}, \dots, \mathbf{a}^{(L)}\}$  via forward propagation

Compute  $\boldsymbol{\delta}^{(L)} = \mathbf{a}^{(L)} - y_i$

Compute errors  $\{\boldsymbol{\delta}^{(L-1)}, \dots, \boldsymbol{\delta}^{(2)}\}$

Compute gradients  $\Delta_{ij}^{(l)} = \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$

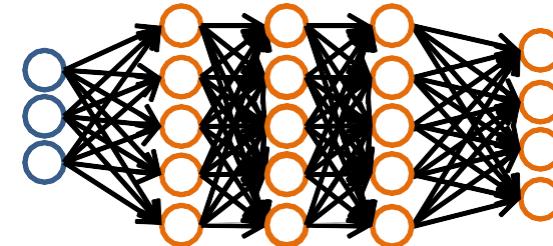
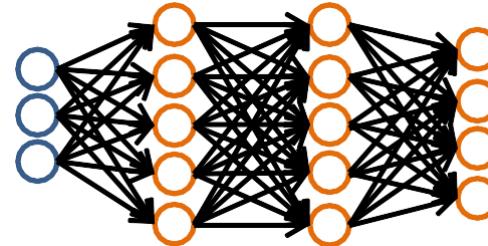
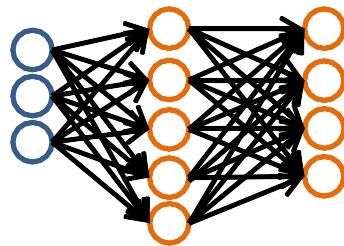
Compute avg regularized gradient  $D_{ij}^{(l)} = \begin{cases} \frac{1}{n} \Delta_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)} & \text{if } j \neq 0 \\ \frac{1}{n} \Delta_{ij}^{(l)} & \text{otherwise} \end{cases}$

Update weights via gradient step  $\Theta_{ij}^{(l)} = \Theta_{ij}^{(l)} - \alpha D_{ij}^{(l)}$

Until weights converge or max #epochs is reached

# Training a Neural Network

Pick a network architecture (connectivity pattern between nodes)



- # input units = # of features in dataset
- # output units = # classes

**Reasonable default:** 1 hidden layer

- or if >1 hidden layer, have same # hidden units in every layer (usually the more the better)

# Training a Neural Network

1. Randomly initialize weights
2. Implement forward propagation to get  $h_{\Theta}(\mathbf{x}_i)$  for any instance  $\mathbf{x}_i$
3. Implement code to compute cost function  $J(\Theta)$
4. Implement backprop to compute partial derivatives
$$\frac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$$
5. Use gradient checking to compare  $\frac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$  computed using backpropagation vs. the numerical gradient estimate.
  - Then, disable gradient checking code
6. Use gradient descent with backprop to fit the network

## More on Backpropagation

---

- Gradient descent over entire *network* weight vector
- Easily generalized to arbitrary directed graphs
- Will find a local, not necessarily global error minimum
  - In practice, often works well (can run multiple times)
- Often include weight *momentum*  $\alpha$ 
$$\Delta w_{i,j}(n) = \eta \delta_j x_{i,j} + \alpha \Delta w_{i,j}(n - 1)$$
- Minimizes error over *training* examples
  - Will it generalize well to subsequent examples?
- Training can take thousands of iterations → slow!
- Using network after training is very fast

# Expressive Capabilities of ANNs

---

Boolean functions:

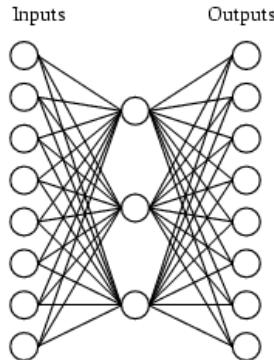
- Every boolean function can be represented by network with single hidden layer
- but might require exponential (in number of inputs) hidden units

Continuous functions:

- Every bounded continuous function can be approximated with arbitrarily small error, by network with one hidden layer [Cybenko 1989; Hornik et al. 1989]
- Any function can be approximated to arbitrary accuracy by a network with two hidden layers [Cybenko 1988].

# Learning Hidden Layer Representations

---



A target function:

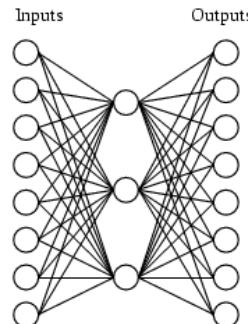
Input	Output
10000000	$\rightarrow$ 10000000
01000000	$\rightarrow$ 01000000
00100000	$\rightarrow$ 00100000
00010000	$\rightarrow$ 00010000
00001000	$\rightarrow$ 00001000
00000100	$\rightarrow$ 00000100
00000010	$\rightarrow$ 00000010
00000001	$\rightarrow$ 00000001

Can this be learned??

# Learning Hidden Layer Representations

---

A network:

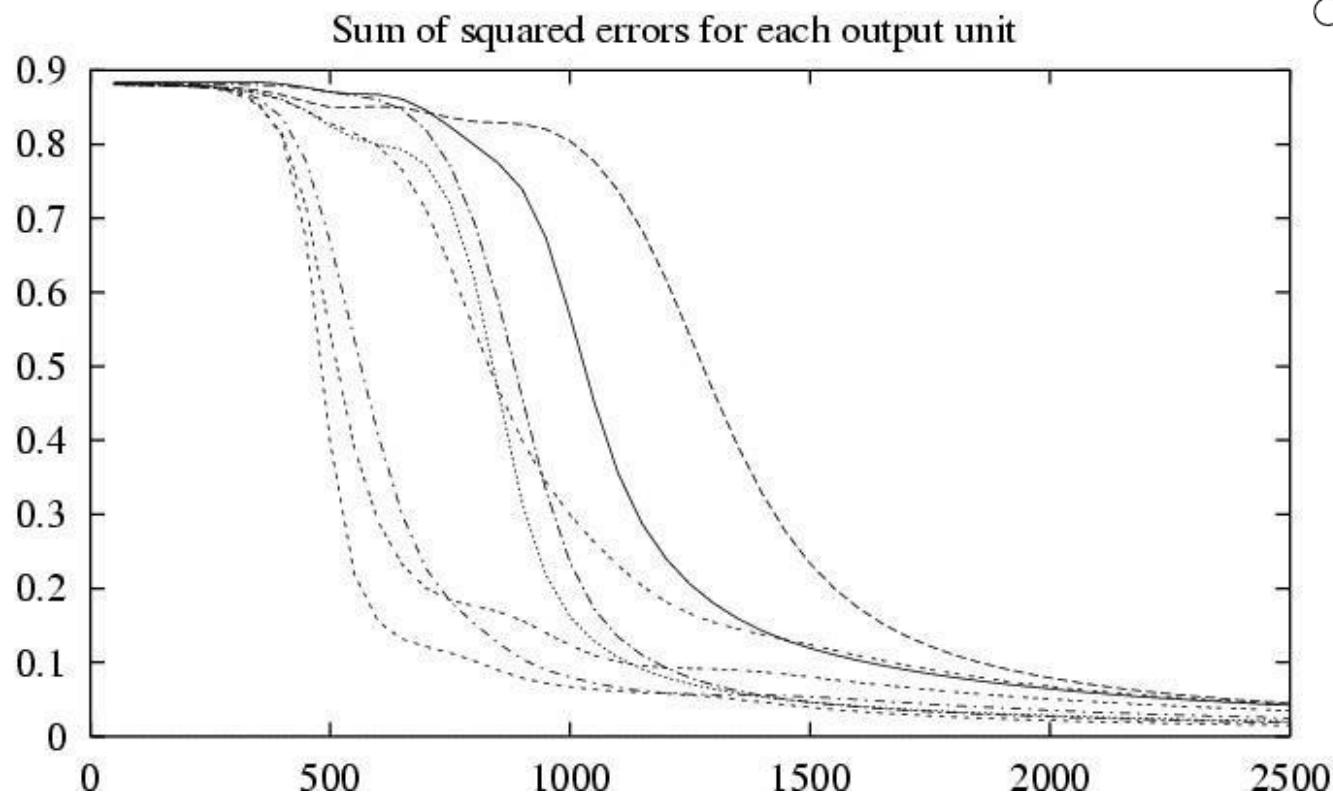
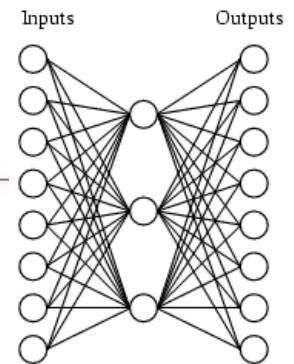


Learned hidden layer representation:

Input	Hidden Values	Output
10000000	→ .89 .04 .08	→ 10000000
01000000	→ .01 .11 .88	→ 01000000
00100000	→ .01 .97 .27	→ 00100000
00010000	→ .99 .97 .71	→ 00010000
00001000	→ .03 .05 .02	→ 00001000
00000100	→ .22 .99 .99	→ 00000100
00000010	→ .80 .01 .98	→ 00000010
00000001	→ .60 .94 .01	→ 00000001

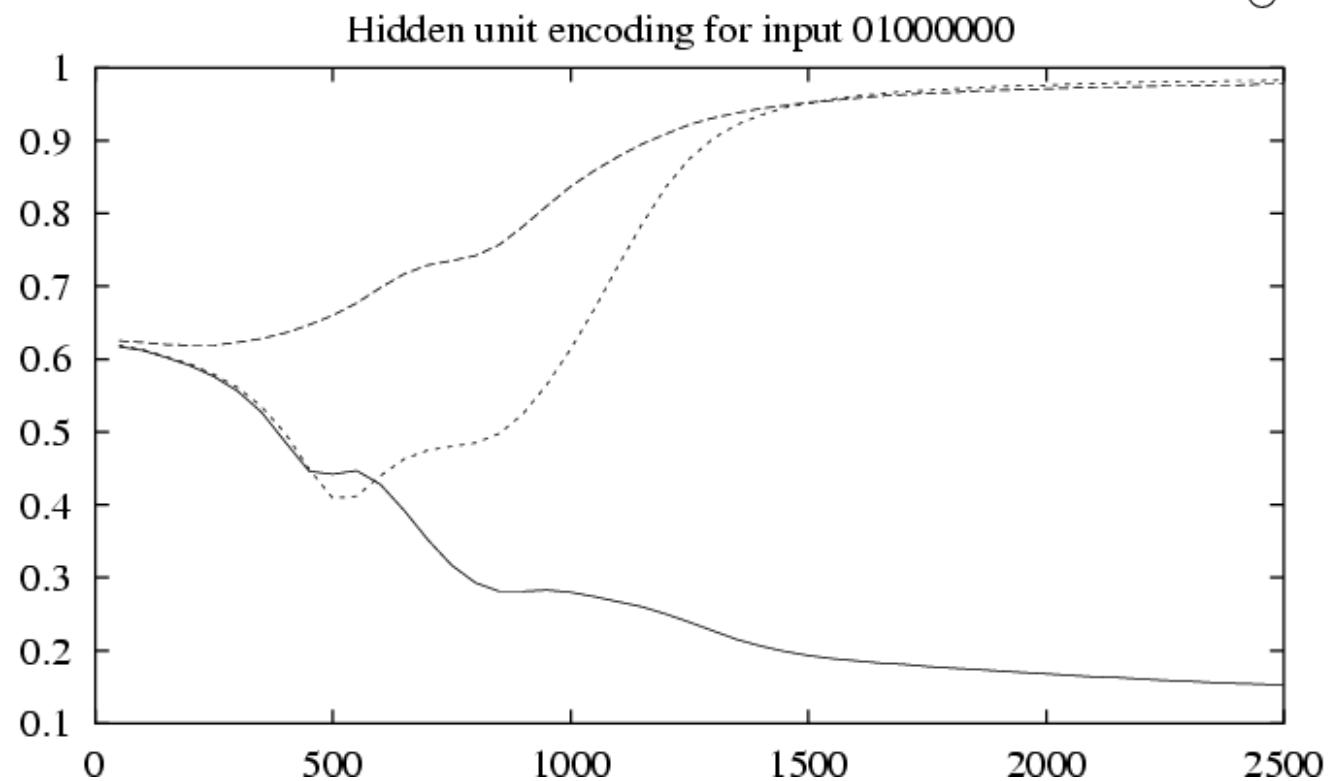
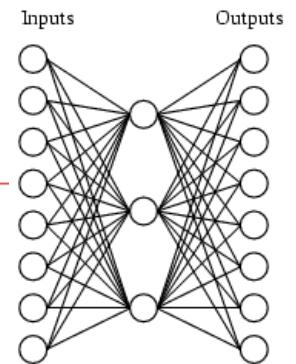
# Training

---



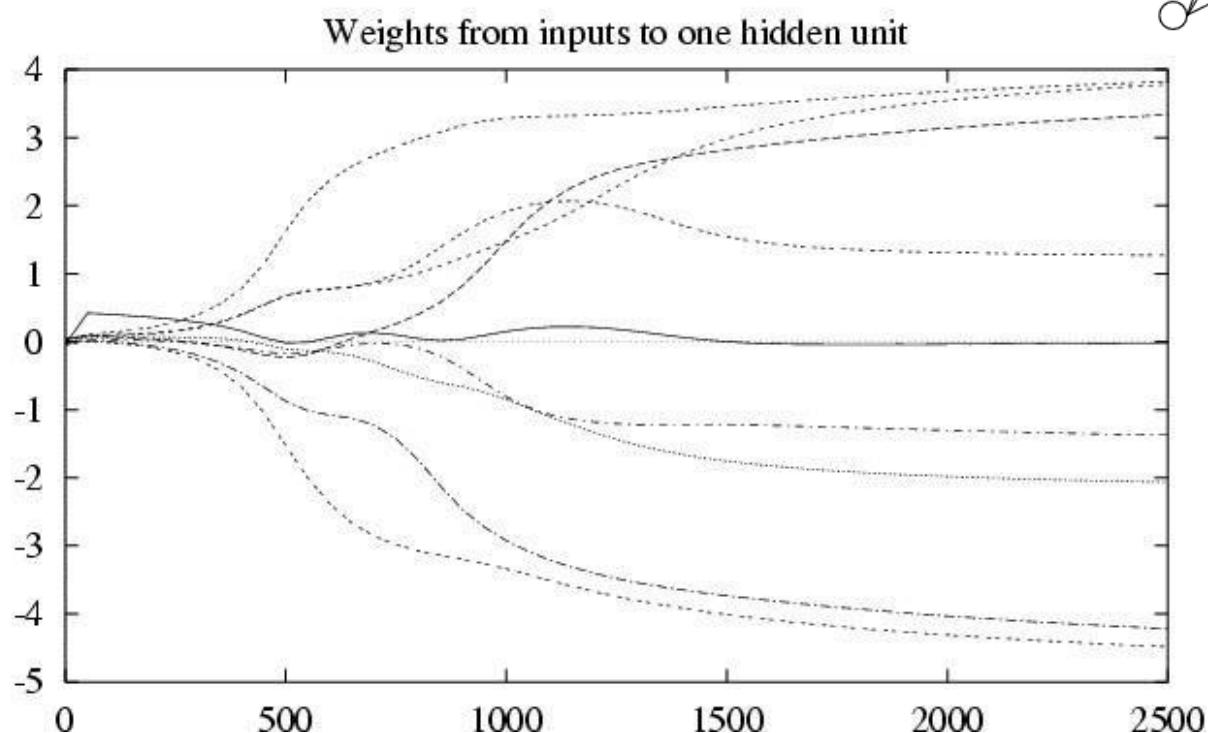
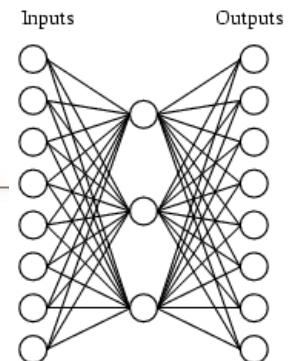
# Training

---



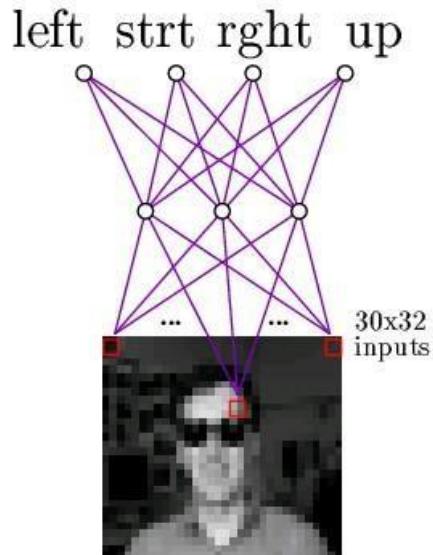
# Training

---



# Neural Nets for Face Recognition

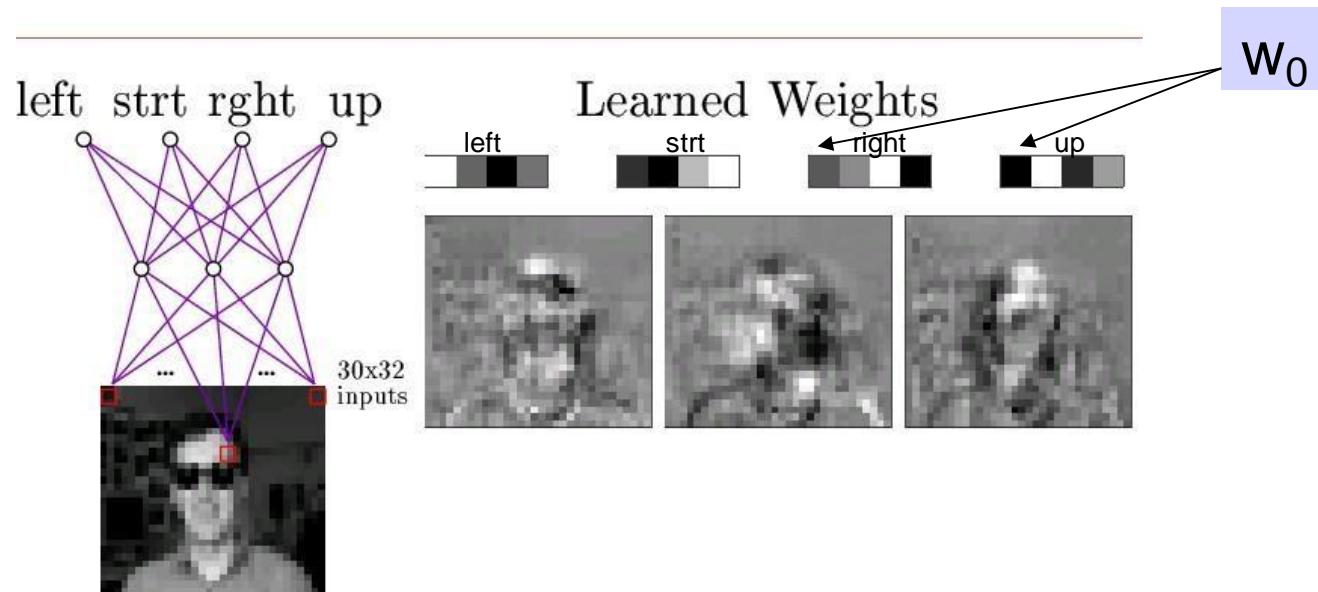
---



Typical input images

90% accurate learning head pose, and recognizing 1-of-20 faces

# Learned Hidden Unit Weights

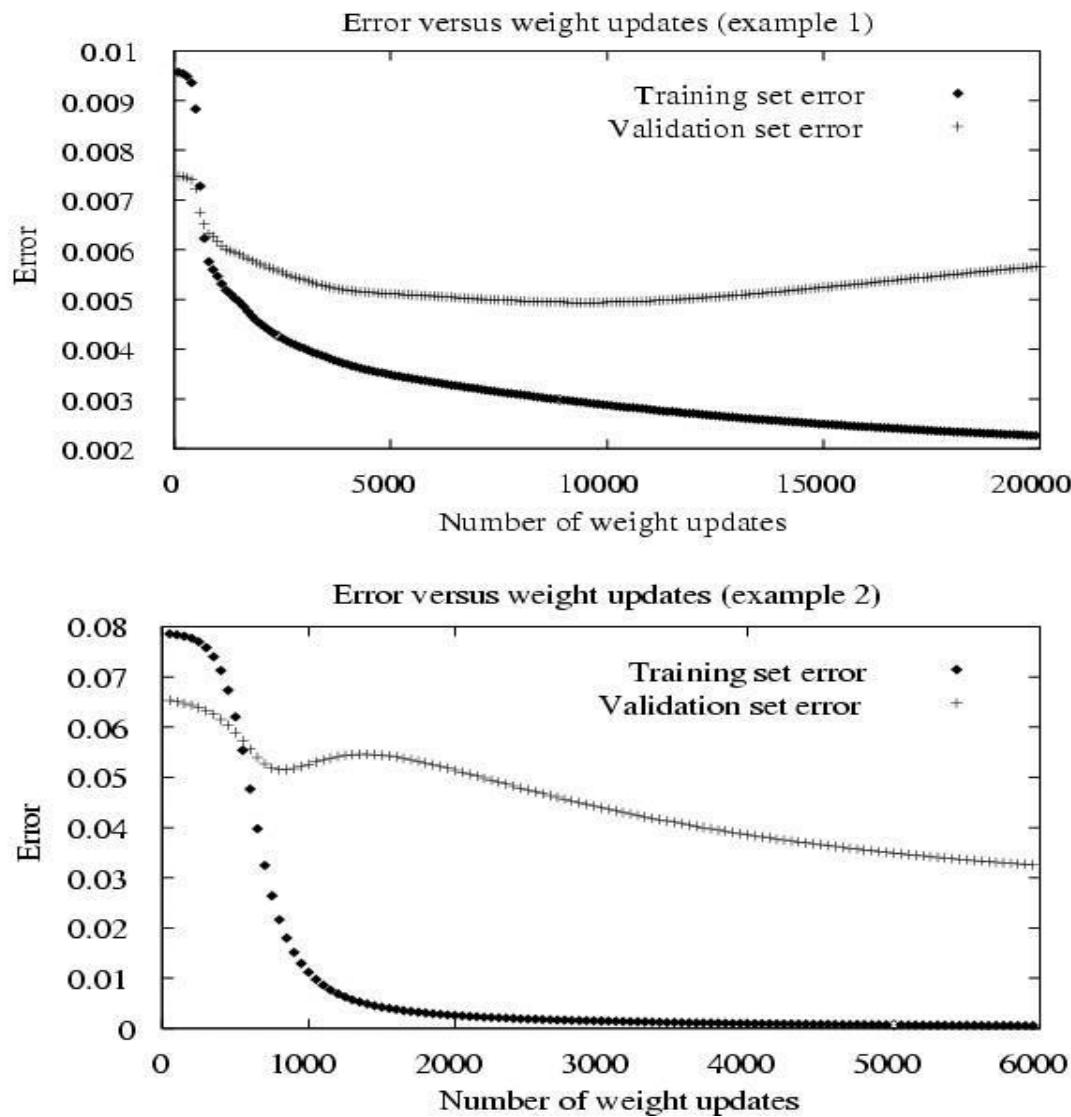


Typical input images

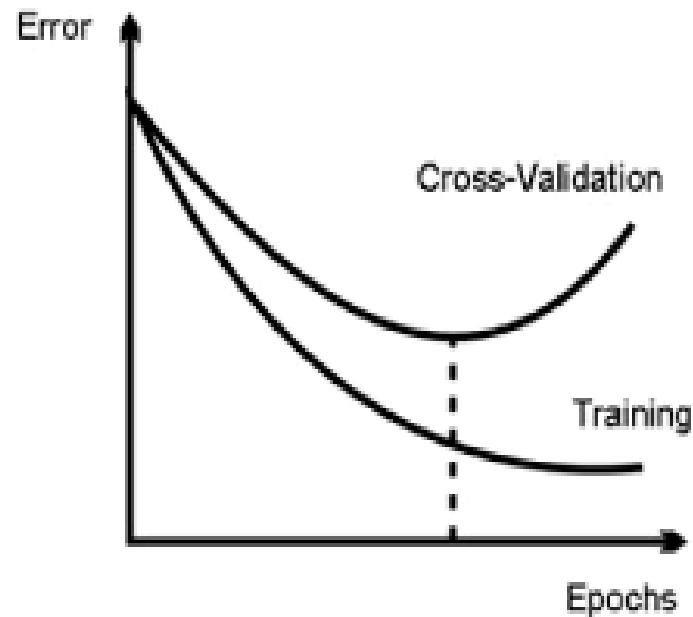
<http://www.cs.cmu.edu/~tom/faces.html>

# Overfitting in ANNs

---



# Early Stopping



**Figure 2:** Profiles for training and cross-validation errors.

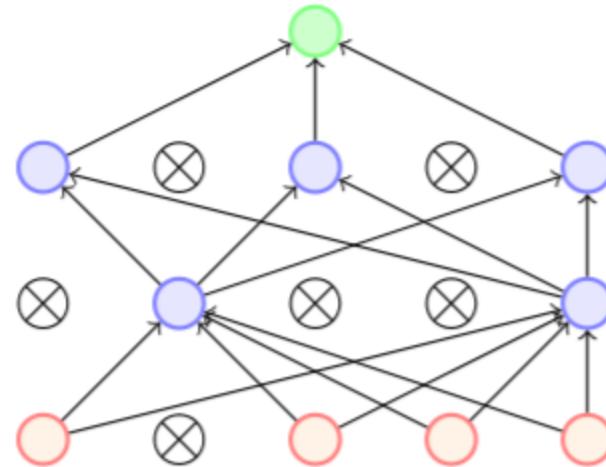
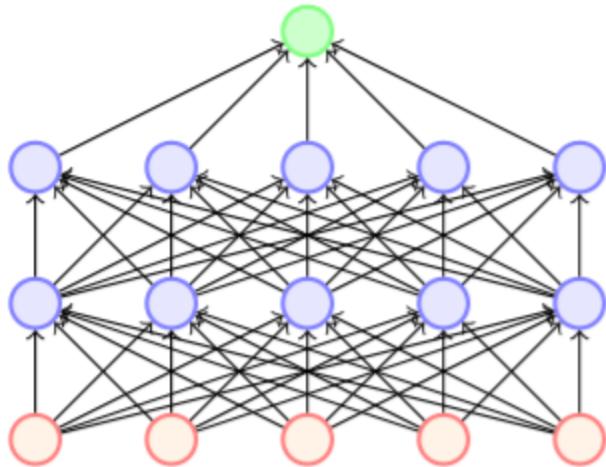
If we let a complex model train long enough on a given data set it can eventually learn the data exactly.

Given data that isn't represented in the training set, the model will perform poorly when analyzing the data (overfitting).

**How is the sweet spot for training located?**

When the error on the training set begins to deviate from the error on the validation set, a threshold can be set to determine the early stopping condition and the ideal number of epochs to train.

# Dropout



- Dropout refers to dropping out units
- Temporarily remove a node and all its incoming/outgoing connections resulting in a thinned network
- Each node is retained with a fixed probability (typically  $p = 0.5$ ) for hidden nodes and  $p = 0.8$  for visible nodes

# Dropouts

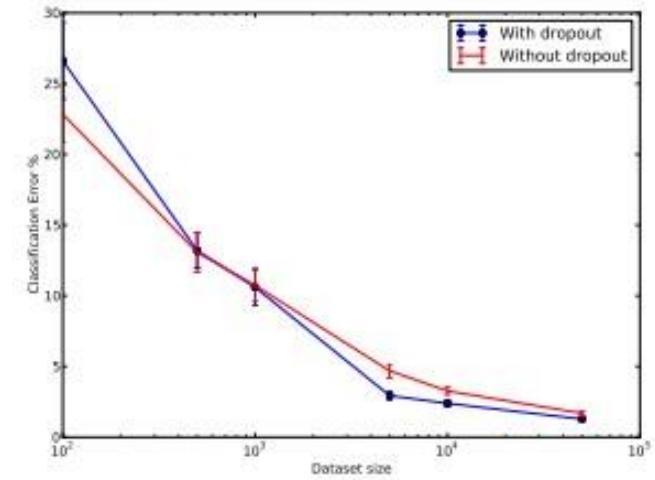
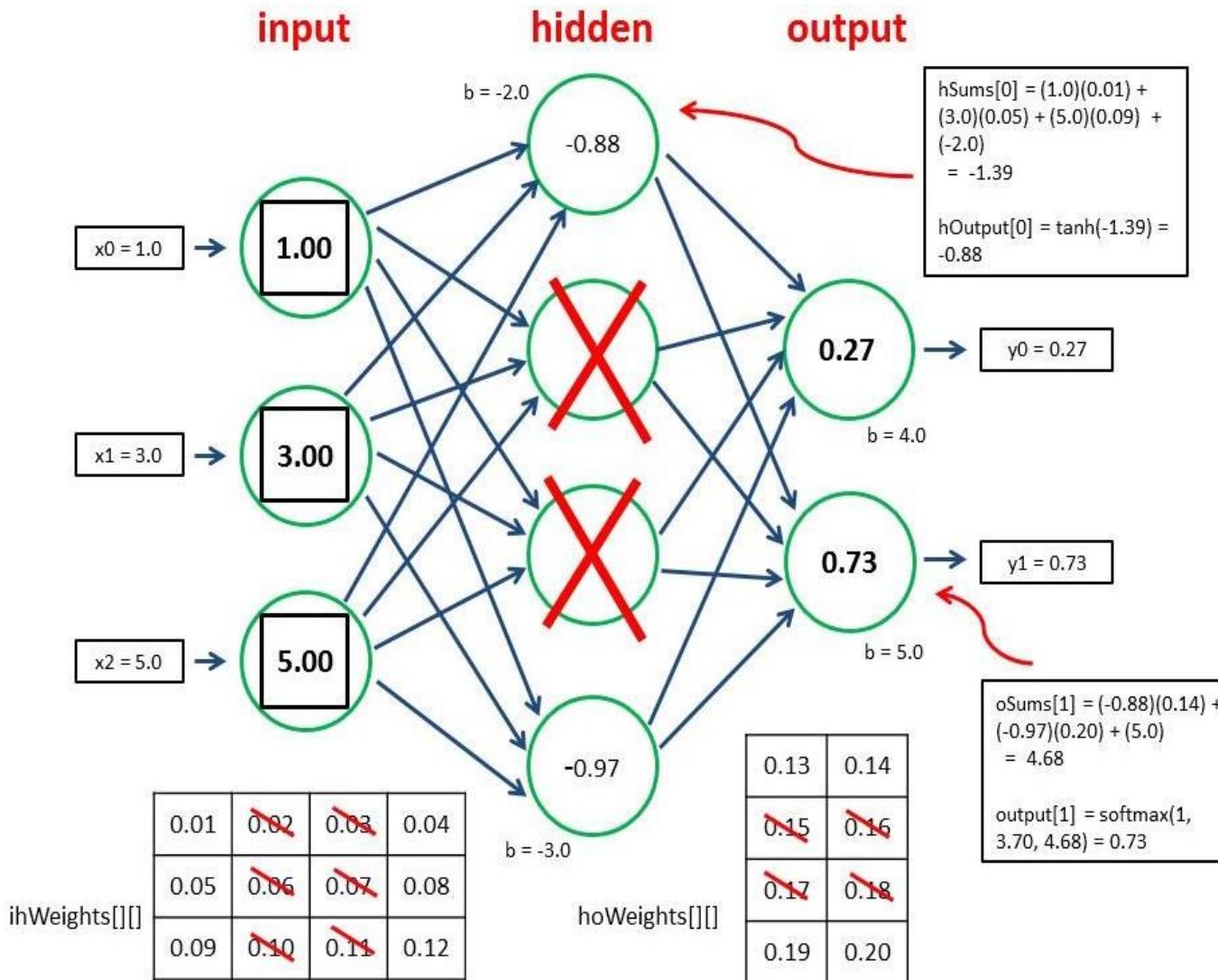


Figure 10: Effect of varying data set size.

# Parameter Sharing

---

Parameter sharing is where we:

- Force sets of parameters to be equal

- Because we interpret various models or model components as sharing a unique set of parameters
- Only a subset of the parameters needs to be stored in memory
  - In a CNN significant reduction in the memory footprint of the model

Sharing parameters gives the network the ability to look for a given feature everywhere in the data(image), rather than in just a certain area. This is extremely useful when the object of interest could be anywhere in the data(image).

It is uncommon to have training data where useful features will usually always be in the same area, so this is not seen often.

# Convolutional Neural Network

# Convolution



$$s_t = \sum_{a=0}^{\infty} x_{t-a} w_{-a} = (x * w)_t$$

input                      filter  
                        convolution

A mathematical equation illustrating convolution. It shows the output  $s_t$  as a weighted sum of past input measurements  $x_{t-a}$  and a filter  $w_{-a}$ . The result is labeled as a convolution of the input  $x$  and the filter  $w$  at time  $t$ .

- Suppose we are tracking the position of an aeroplane using a laser sensor at discrete time intervals
- Now suppose our sensor is noisy
- To obtain a less noisy estimate we would like to average several measurements
- More recent measurements are more important so we would like to take a weighted average

# Convolution



- We can think of images as 2D inputs
- We would now like to use a 2D filter ( $m \times n$ )
- First let us see what the 2D formula looks like
- This formula looks at all the preceding neighbours  $(i - a, j - b)$

$$S_{ij} = (I * K)_{ij} = \sum_{a=0}^{m-1} \sum_{b=0}^{n-1} I_{i-a, j-b} K_{a,b}$$

# Convolution



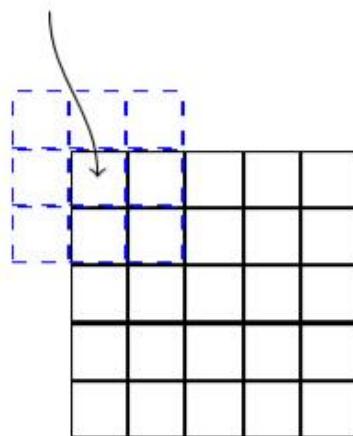
$$S_{ij} = (I * K)_{ij} = \sum_{a=0}^{m-1} \sum_{b=0}^{n-1} I_{i+a, j+b} K_{a,b}$$

- We can think of images as 2D inputs
- We would now like to use a 2D filter ( $m \times n$ )
- First let us see what the 2D formula looks like
- This formula looks at all the preceding neighbours ( $i - a, j - b$ )
- In practice, we use the following formula which looks at the succeeding neighbours

# Convolution

$$S_{ij} = (I * K)_{ij} = \sum_{a=\left\lfloor -\frac{m}{2} \right\rfloor}^{\left\lfloor \frac{m}{2} \right\rfloor} \sum_{b=\left\lfloor -\frac{n}{2} \right\rfloor}^{\left\lfloor \frac{n}{2} \right\rfloor} I_{i-a, j-b} K_{\frac{m}{2}+a, \frac{n}{2}+b}$$

pixel of interest



- For the rest of the discussion we will use the following formula for convolution
- In other words we will assume that the kernel is centered on the pixel of interest
- So we will be looking at both preceding and succeeding neighbors

# Let us apply this idea to a toy example and see the results

Input

a	b	c	d
e	f	g	h
i	j	k	A

Kernel

w	x
y	z

Output

aw+bx+ey+fz		

# Let us apply this idea to a toy example and see the results

Input

a	b	c	d
e	f	g	h
i	j	k	A

Kernel

w	x
y	z

Output

aw+bx+ey+fz	bw+cx+fy+gz	

Let us apply this idea to a toy example and see the results

Input

a	b	c	d
e	f	g	h
i	j	k	A

Kernel

w	x
y	z

Output

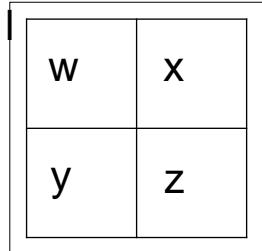
aw+bx+ey+fz	bw+cx+fy+gz	cw+dx+gy+hz

# Let us apply this idea to a toy example and see the results

Input

a	b	c	d
e	f	g	h
i	j	k	A

Kerne



Output

$aw+bx+ey+fz$	$bw+cx+fy+gz$	$cw+dx+gy+hz$
$ew+fx+iy+jz$		

# Let us apply this idea to a toy example and see the results

Input

a	b	c	d
e	f	g	h
i	j	k	A

Kernel

w	x
y	z

Output

$aw+bx+ey+fz$	$bw+cx+fy+gz$	$cw+dx+gy+hz$
$ew+fx+iy+jz$	$fw+gx+jy+kz$	

# Let us apply this idea to a toy example and see the results

Input

a	b	c	d
e	f	g	h
i	j	k	A

Kernel

w	x
y	z

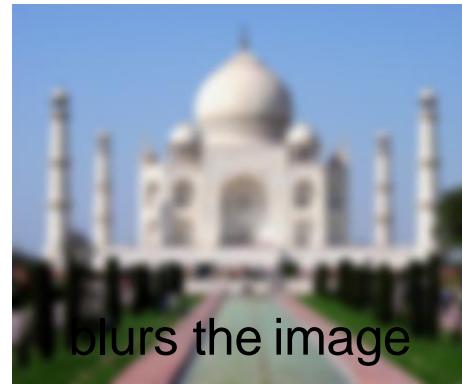
Output

$aw+bx+ey+fz$	$bw+cx+fy+gz$	$cw+dx+gy+hz$
$ew+fx+iy+jz$	$fw+gx+jy+kz$	$gw+hx+ky+Az$

# Example of kernel: Blur



$$* \begin{matrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{matrix} =$$



# Example of kernel: Edge detection

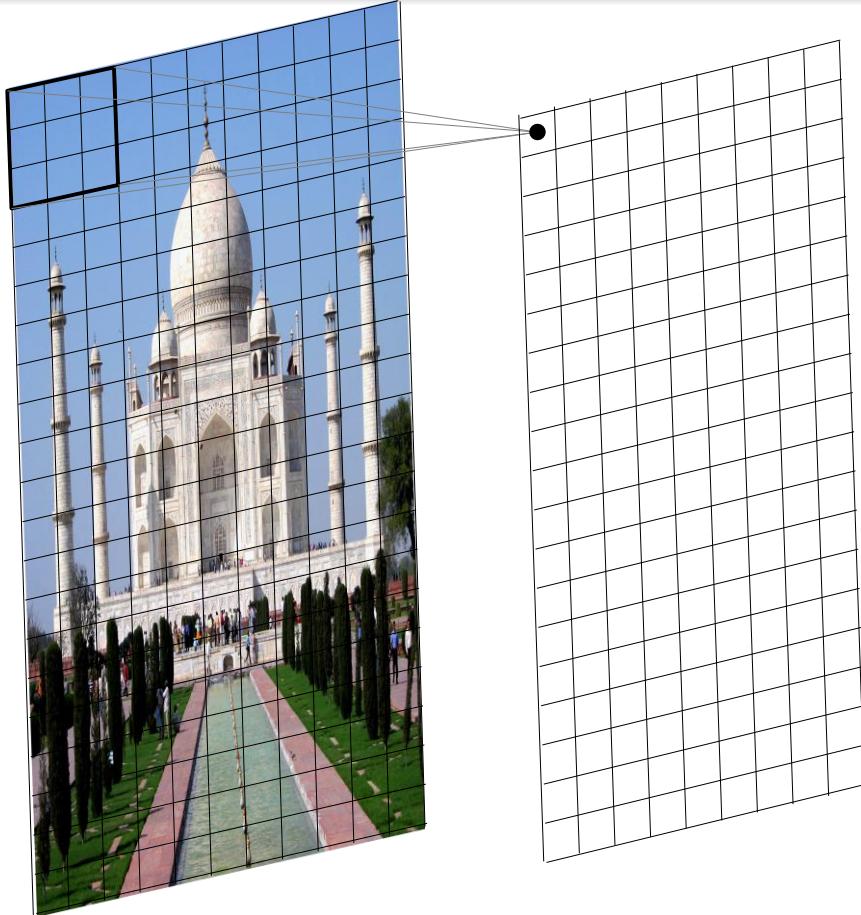


$$\begin{array}{ccc} * & \begin{matrix} 1 & 1 & 1 \\ 1 & -8 & 1 \\ 1 & 1 & 1 \end{matrix} & = \end{array}$$



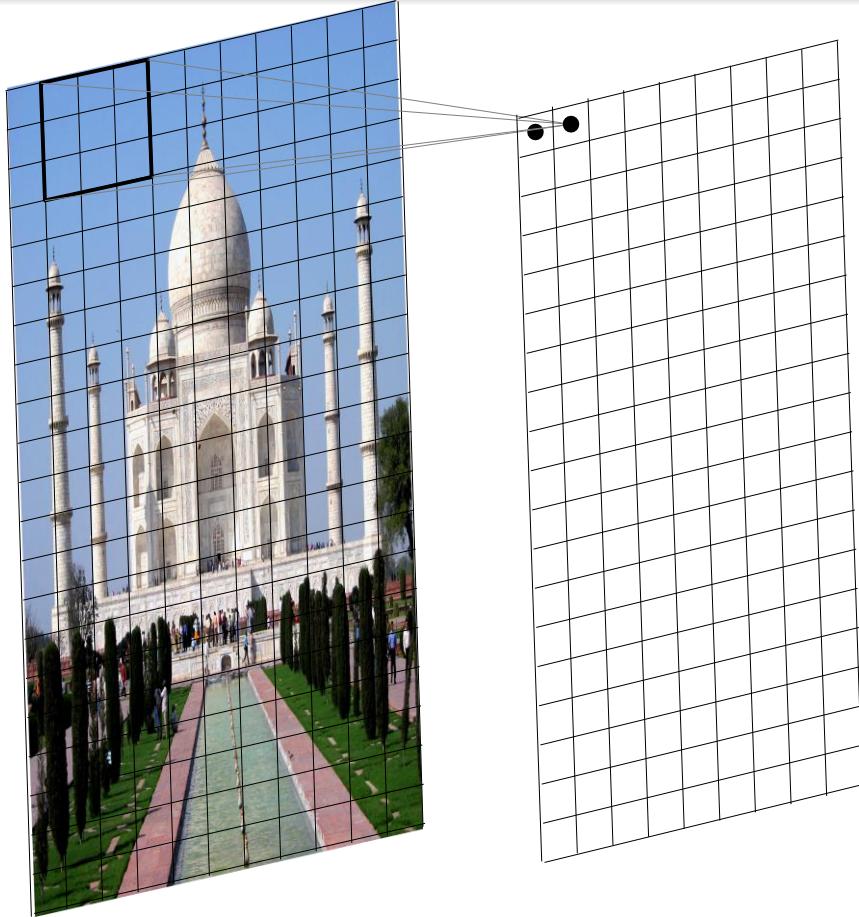
edges

# Convolution



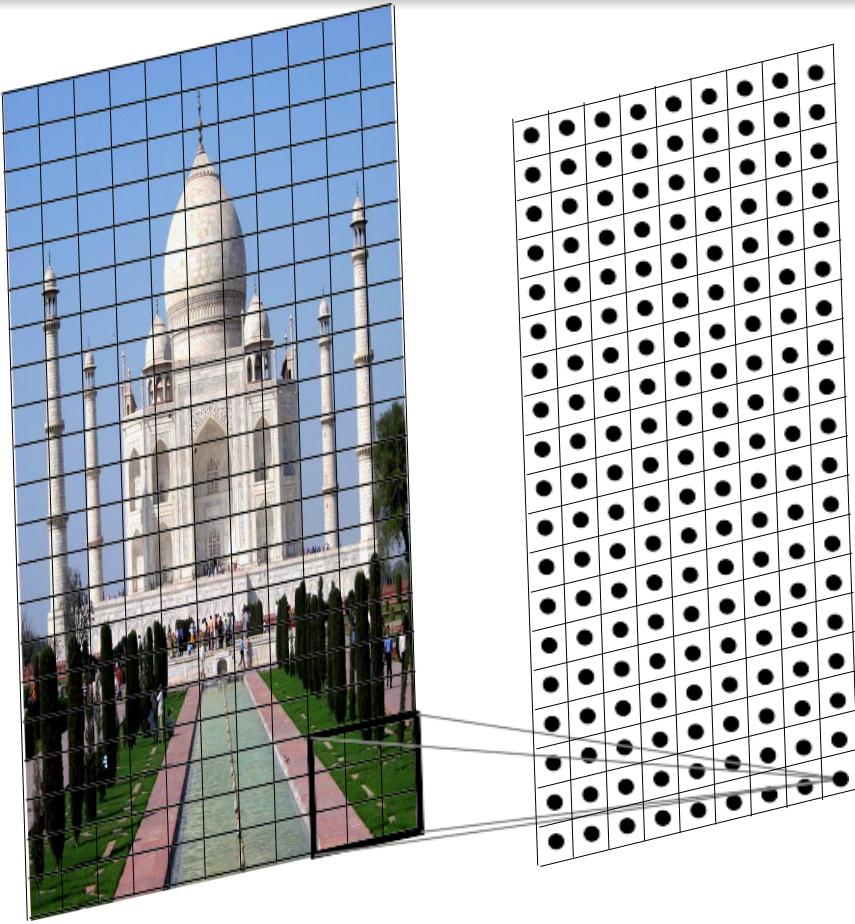
- We just slide the kernel over the input image
- Each time we slide the kernel we get one value in the output

# Convolution

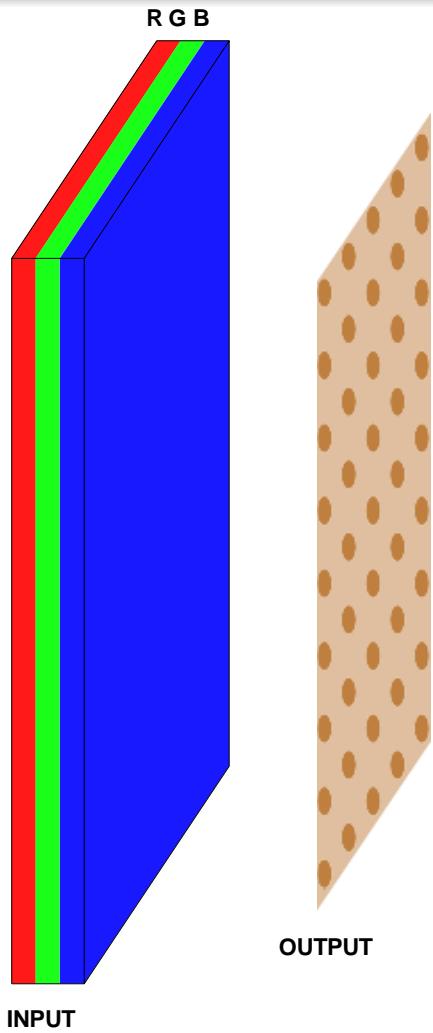


- We just slide the kernel over the input image
- Each time we slide the kernel we get one value in the output

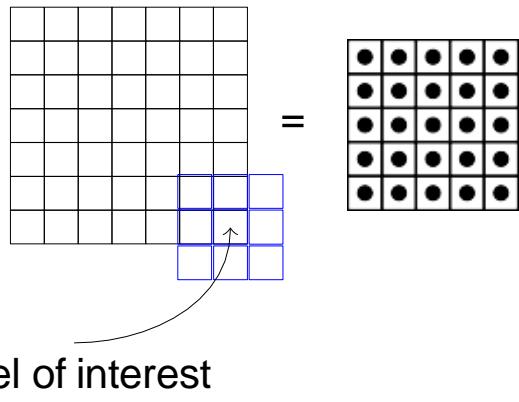
# 2D convolutions applied to images



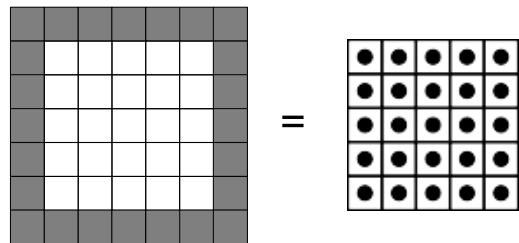
- We just slide the kernel over the input image
- Each time we slide the kernel we get one value in the output
- The resulting output is called a feature map.
- We can use multiple filters to get multiple feature maps.



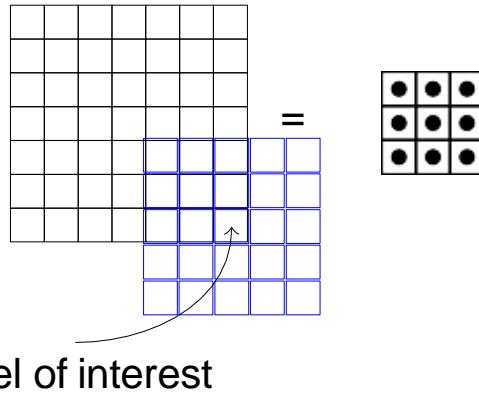
- What would a 3D filter look like?
- It will be 3D and we will refer to it as a volume
- Once again we will slide the volume over the 3D input and compute the convolution operation
- We will assume that the filter always extends to the depth of the image
  
- In effect, we are doing a 2D convolution operation on a 3D input (because the filter moves along the height and the width but not along the depth)
- As a result the output will be 2D (only width and height, no depth)
- Once again we can apply multiple filters to get multiple feature maps



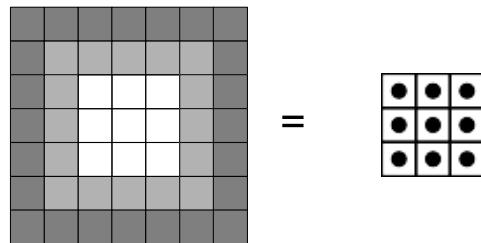
- Let us compute the dimension ( $W_2, H_2$ ) of the output
- Notice that we can't place the kernel at the corners as it will cross the input boundary



- Let us compute the dimension ( $W_2, H_2$ ) of the output
- Notice that we can't place the kernel at the corners as it will cross the input boundary
- This is true for all the shaded points (the kernel crosses the input boundary)
- This results in an output which is of smaller dimensions than the input
- As the size of the kernel increases, this becomes true for even more pixels



- Let us compute the dimension ( $W_2, H_2$ ) of the output
- Notice that we can't place the kernel at the corners as it will cross the input boundary
- This is true for all the shaded points (the kernel crosses the input boundary)
- This results in an output which is of smaller dimensions than the input
- As the size of the kernel increases, this becomes true for even more pixels
- For example, let's consider a  $5 \times 5$  kernel
- We have an even smaller output now



*In general,*  $W_2 = W_1 - F + 1$

$$H_2 = H_1 - F + 1$$

*We will refine this formula further*

- Let us compute the dimension ( $W_2, H_2$ ) of the output
- Notice that we can't place the kernel at the corners as it will cross the input boundary
- This is true for all the shaded points (the kernel crosses the input boundary)
- This results in an output which is of smaller dimensions than the input
- As the size of the kernel increases, this becomes true for even more pixels
- For example, let's consider a  $5 \times 5$  kernel
- We have an even smaller output now

$$\begin{array}{ccccccccc}
 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & & & & & & & & 0 \\
 0 & & & & & & & & 0 \\
 0 & & & & & & & & 0 \\
 0 & & & & & & & & 0 \\
 0 & & & & & & & & 0 \\
 0 & & & & & & & & 0 \\
0 & & & & & & & & 0 \\
0 & & & & & & & & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0
\end{array} = 
 \begin{array}{ccccccccc}
\bullet & \bullet \\
\bullet & \bullet \\
\bullet & \bullet \\
\bullet & \bullet \\
\bullet & \bullet \\
\bullet & \bullet \\
\bullet & \bullet \\
\bullet & \bullet
\end{array}$$

- What if we want the output to be of same size as the input?
- We can use something known as padding
- Pad the inputs with appropriate number of 0 inputs so that you can now apply the kernel at the corners
- Let us use pad  $P = 1$  with a  $3 \times 3$  kernel
- This means we will add one row and one column of 0 inputs at the top, bottom, left and right

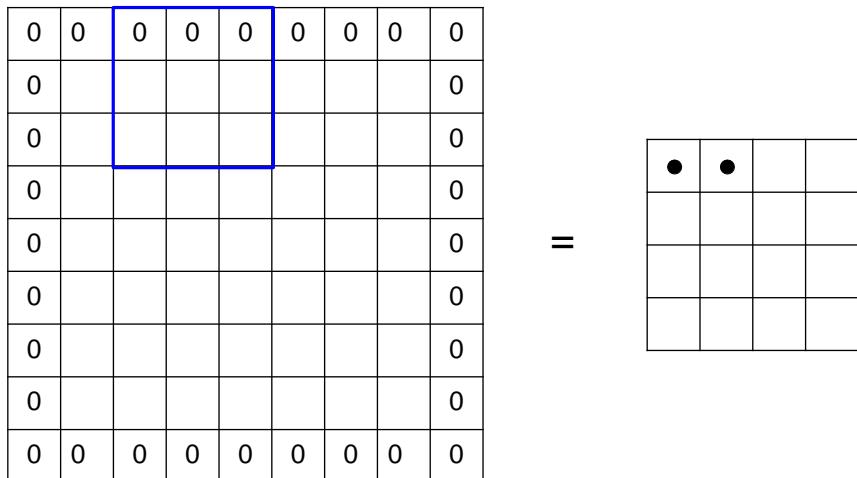
We now have,  $W_2 = W_1 - F + 2P + 1$   $H_2 = H_1 - F + 2P + 1$

0	0	0	0	0	0	0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0

=

•			

- What does the stride S do?
- It defines the intervals at which the filter is applied (here  $S = 2$ )
- Here, we are essentially skipping every 2nd pixel which will again result in an output which is of smaller dimensions

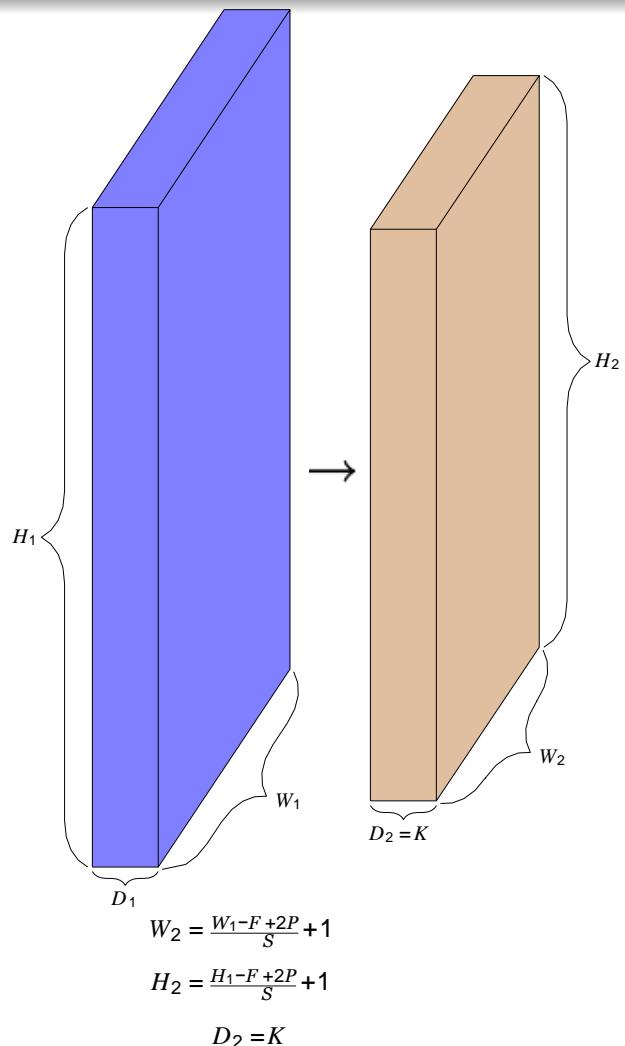


- What does the stride S do?
- It defines the intervals at which the filter is applied (here  $S = 2$ )
- Here, we are essentially skipping every 2nd pixel which will again result in an output which is of smaller dimensions

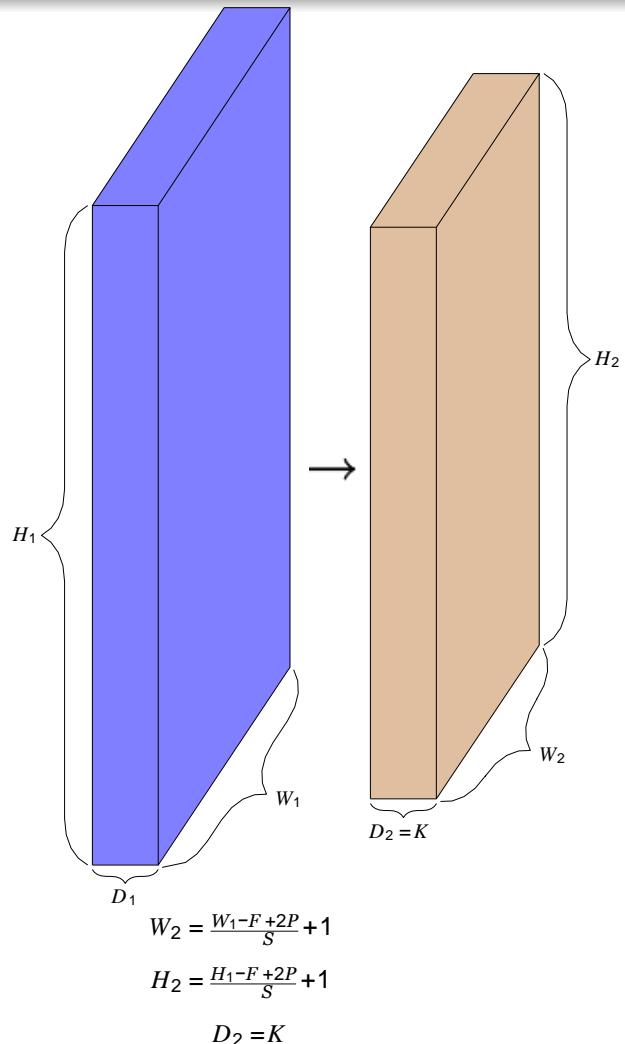
*So what should our final formula look like,*

$$W = \frac{W_1 - F + 2P}{S} + 1$$

$$H = \frac{H_1 - F + 2P}{S} + 1$$

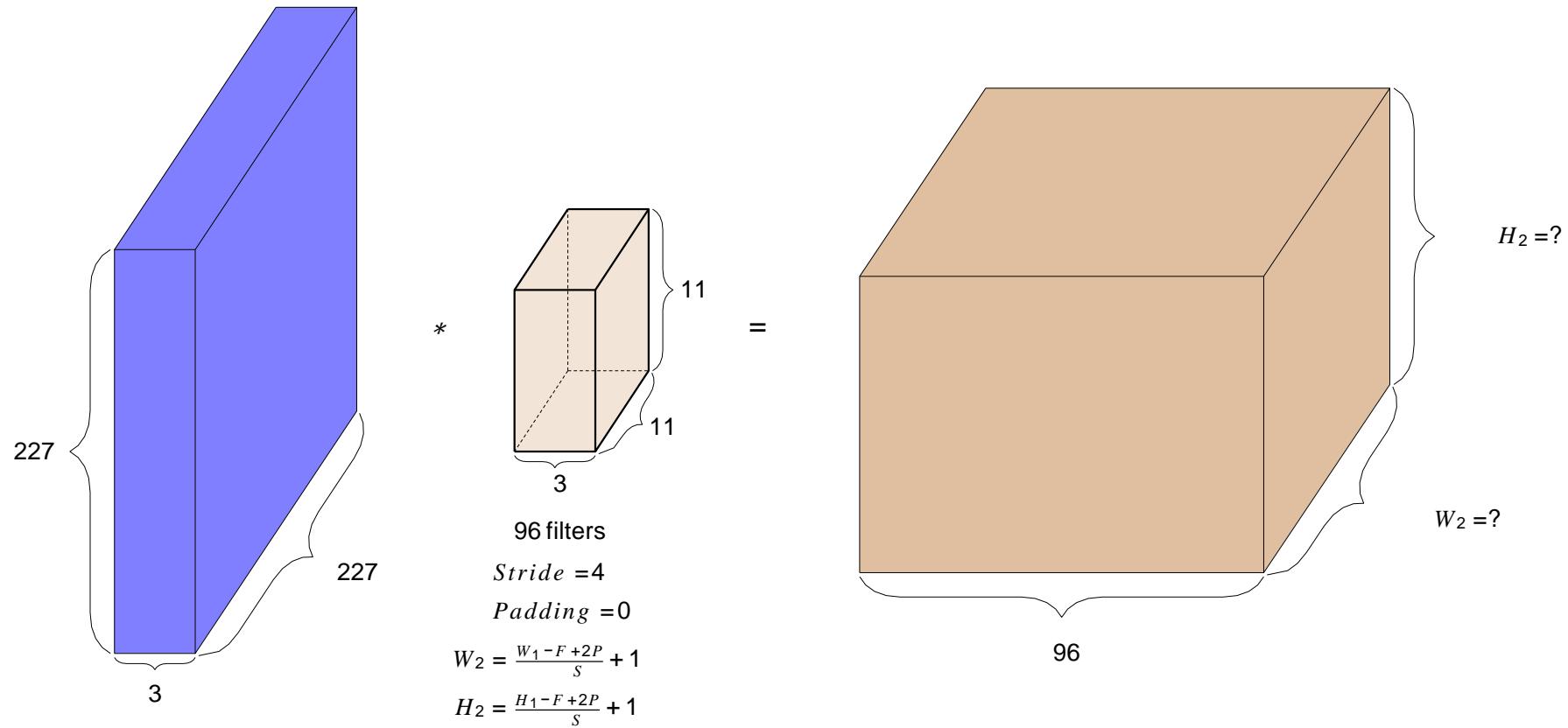


- Finally, coming to the depth of the output.
- Each filter gives us one 2D output.
- $K$  filters will give us  $K$  such 2D outputs

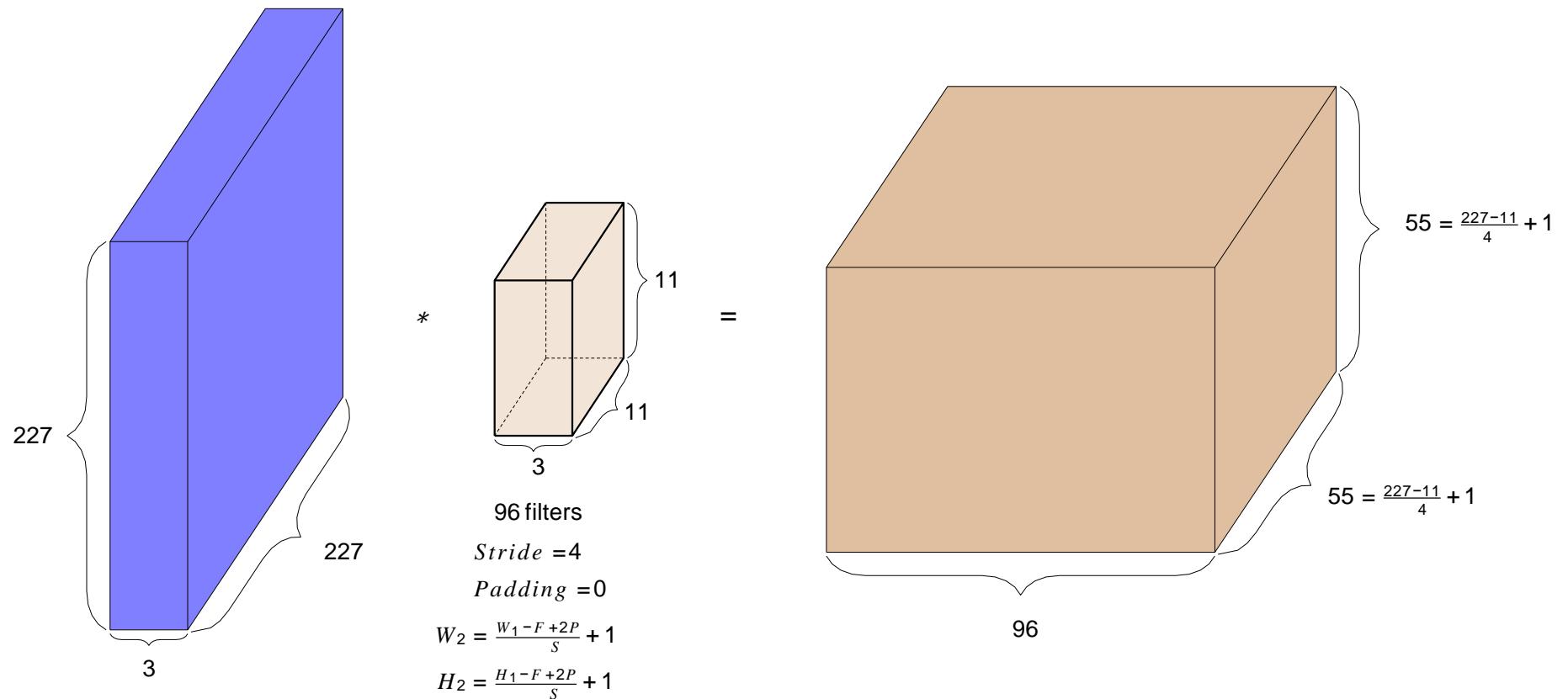


- Finally, coming to the depth of the output.
- Each filter gives us one 2D output.
- $K$  filters will give us  $K$  such 2D outputs
- We can think of the resulting output as  $K \times W_2 \times H_2$  volume
- Thus  $D_2 = K$

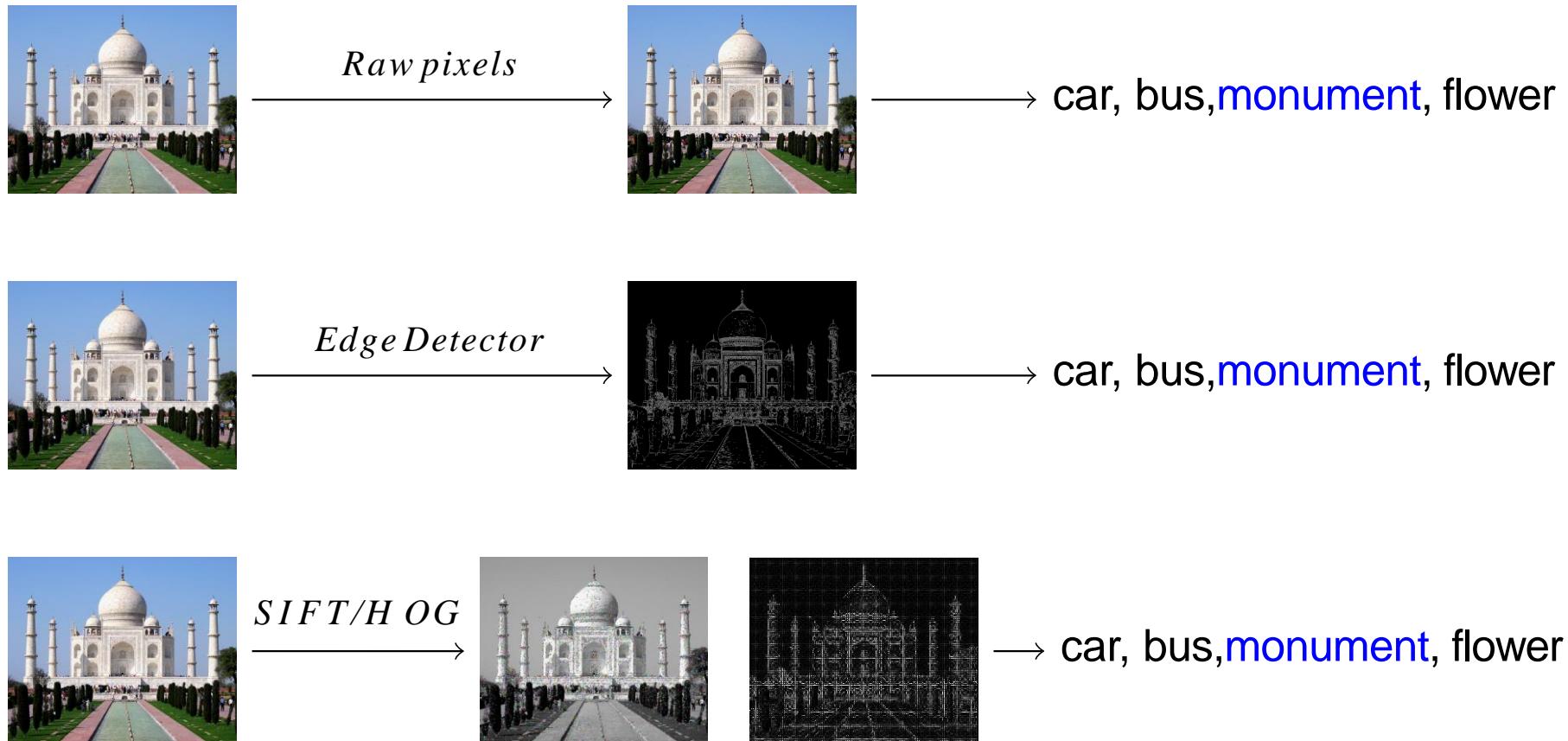
# Example

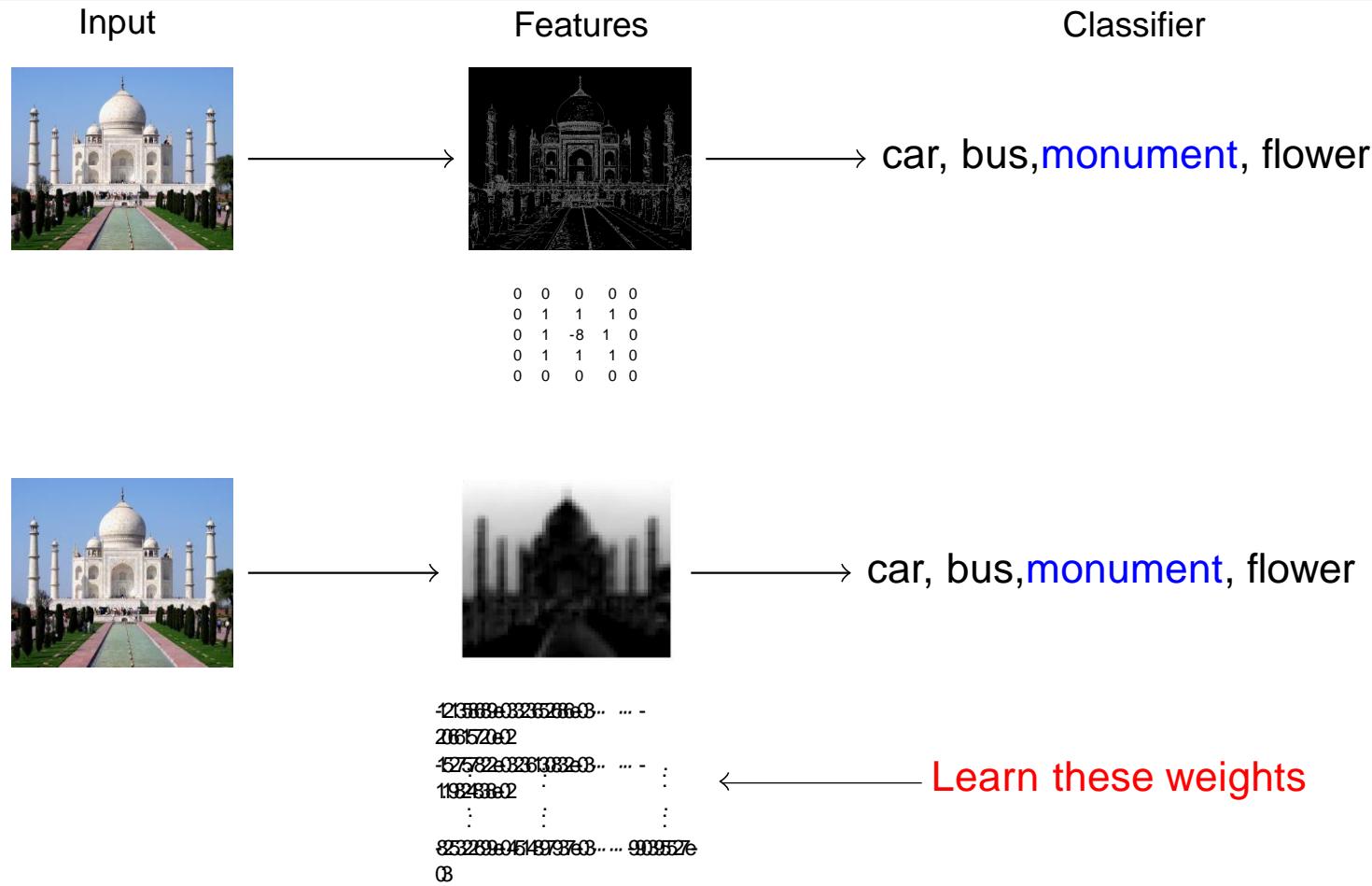


# Example

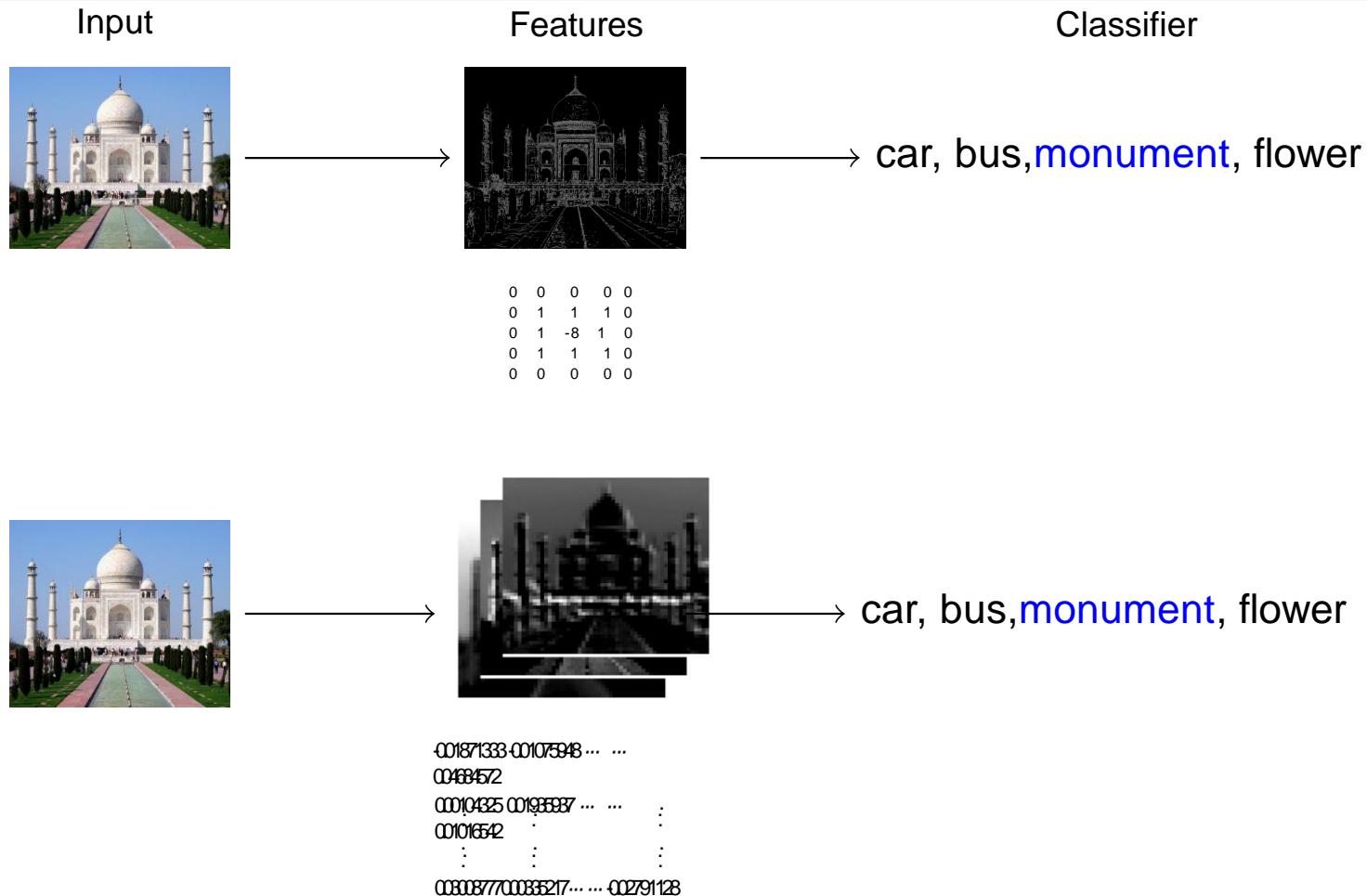


# Features



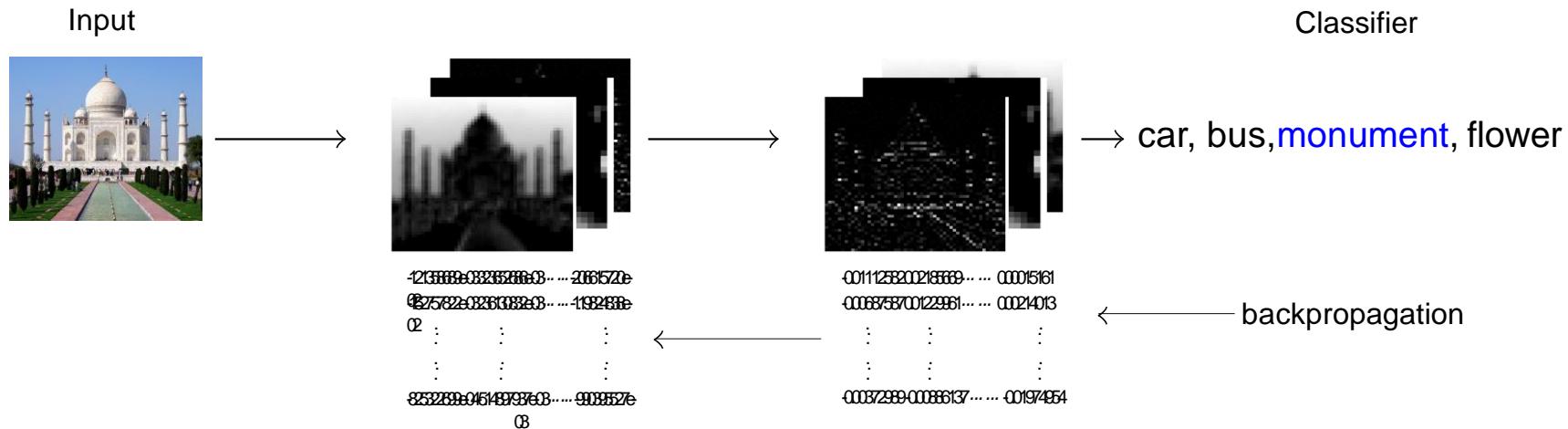


■ Instead of using handcrafted kernels such as edge detectors **can we learn meaningful kernels/filters in addition to learning the weights of the classifier?**

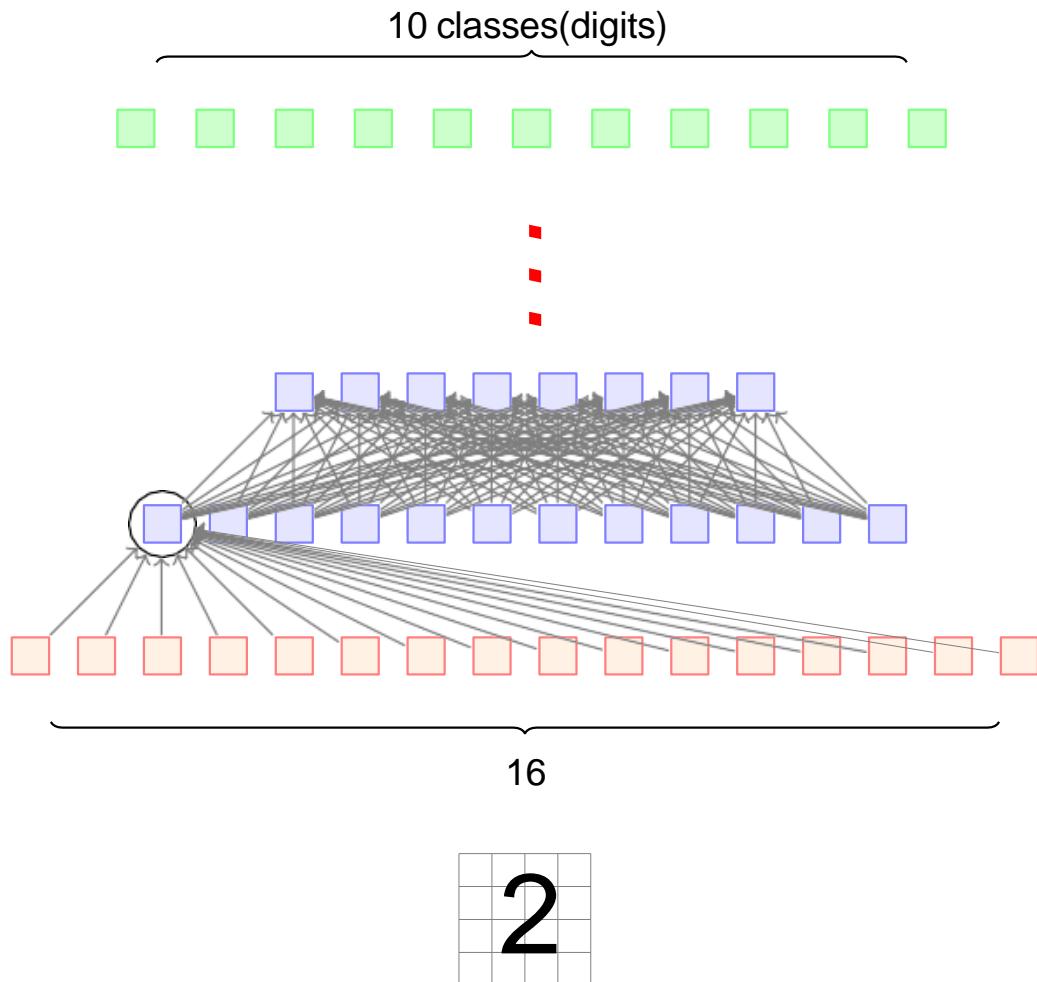


- Even better: Instead of using handcrafted kernels (such as edge detectors) can we learn **multiple meaningful kernels/filters** in addition to learning the weights of the classifier?

# Convolutional Neural Network

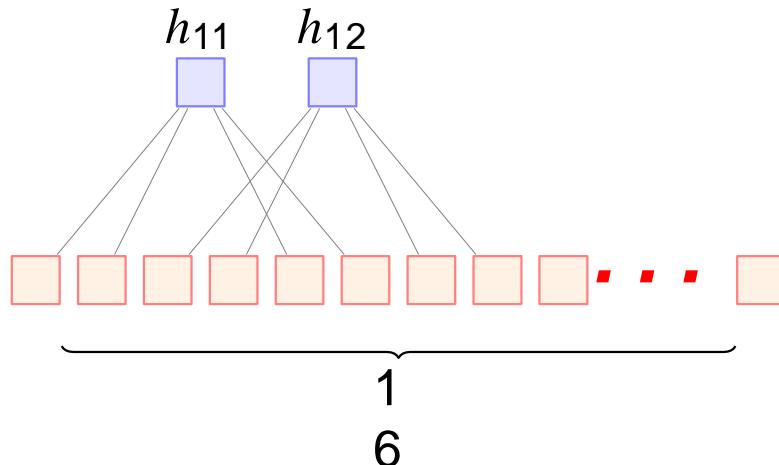


- Can we learn multiple **layers** of meaningful **kernels/filters** in addition to learning the weights of the classifier?
- Yes, we can !
- Simply by treating these kernels as parameters and learning them in addition to the weights of the classifier (using back propagation)
- Such a network is called a Convolutional Neural Network.



- This is what a regular feed-forward neural network will look like
- There are many dense connections here
- For example all the 16 input neurons are contributing to the computation of  $h_{11}$
- Contrast this to what happens in the case of convolution

# Sparse Connectivity

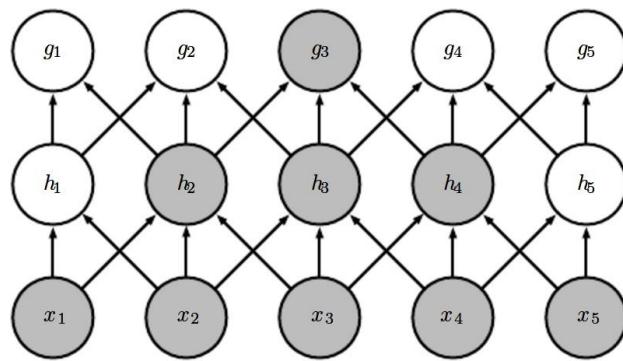


The diagram shows a 4x4 input matrix with a handwritten digit '2' and a 3x3 weight matrix. An asterisk (\*) indicates multiplication, and an equals sign (=) indicates the result. The result is a single blue square, representing the output of the neuron  $h_{11}$ .

- Only a few local neurons participate in the computation of  $h_{11}$
- For example, only pixels 1, 2, 5, 6 contribute to  $h_{11}$
- The connections are much sparser
- We are taking advantage of the structure of the image (interactions)
- This **sparse connectivity** reduces the number of parameters in the model

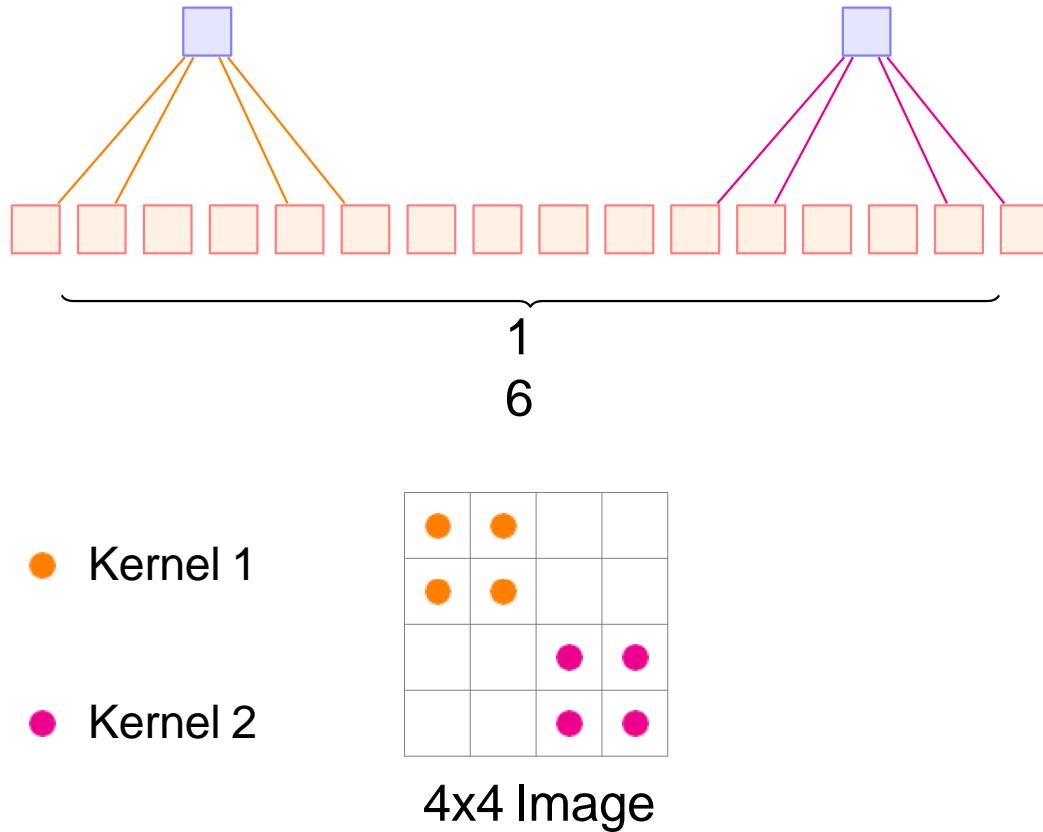
# Sparse Connectivity

- But is sparse connectivity really good thing ?
- Aren't we losing information (by losing interactions between some input pixels)
- Well, not really
- The two highlighted neurons ( $x_1$  &  $x_5$ ) do not interact in *layer 1*
- But they indirectly contribute to the computation of  $g_3$  and hence interact indirectly



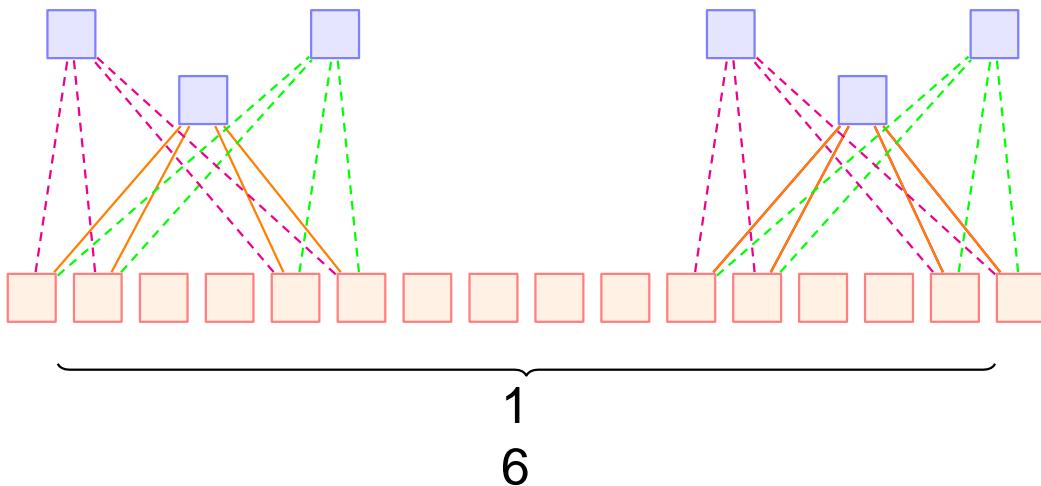
\*Goodfellow-et-al-2016

# Weight-sharing

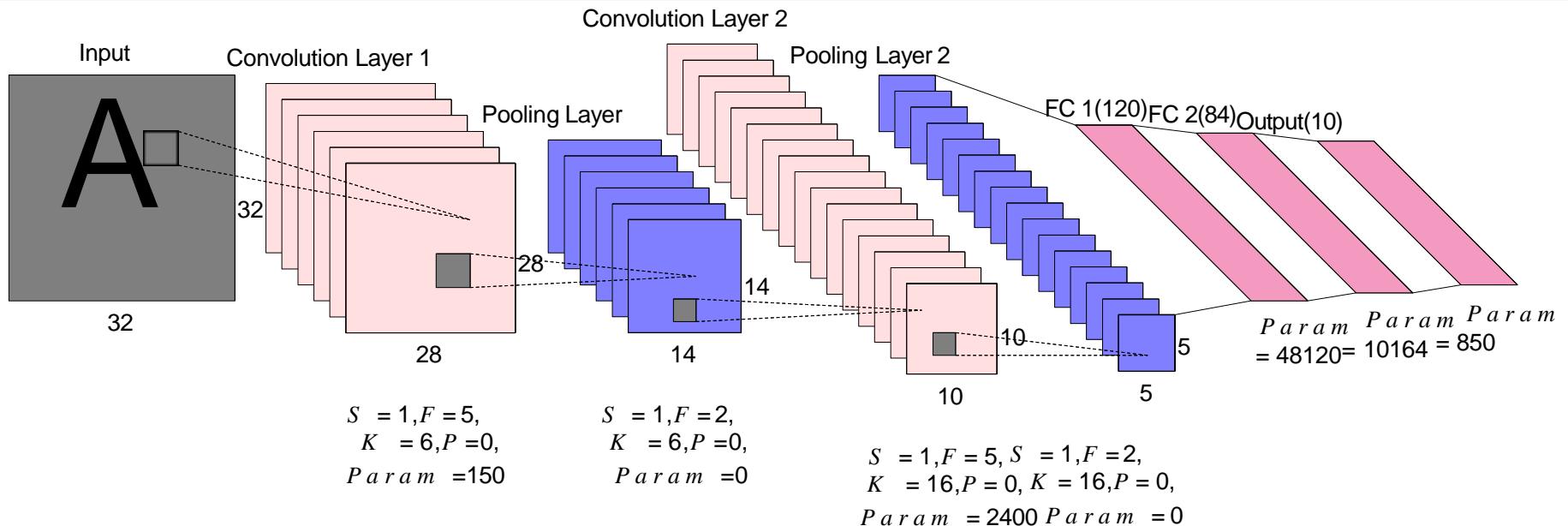


- Another characteristic of CNNs is **weight sharing**
- Consider the following network
- Do we want the kernel weights to be different for different portions of the image?
- Imagine that we are trying to learn a kernel that detects edges
- Shouldn't we be applying the same kernel at all the portions of the image?

# Weight-sharing

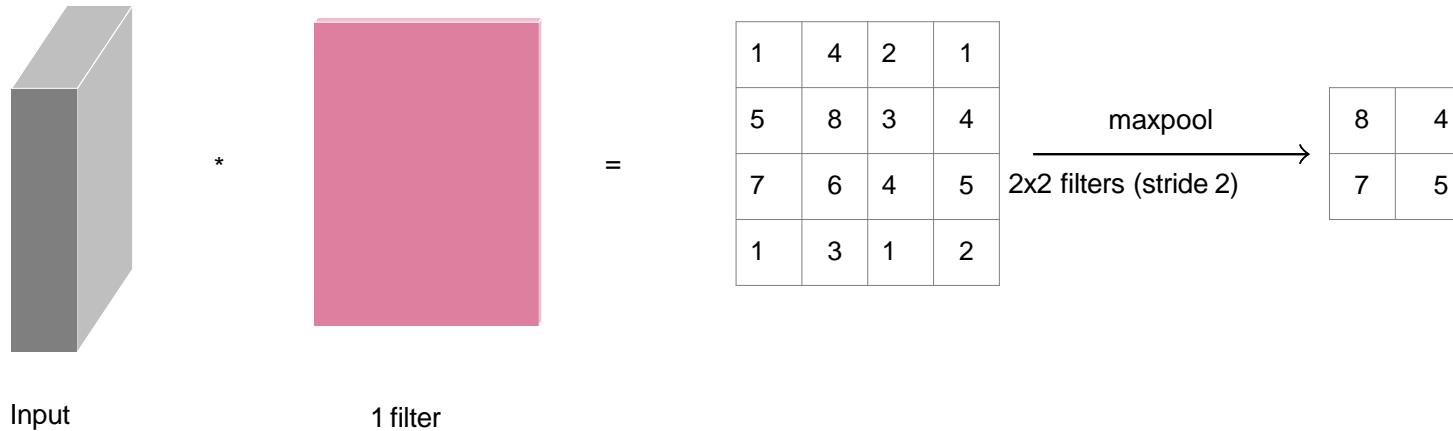


- In other words shouldn't the *orange* and *pink* kernels be the same
- Yes, indeed
- This would make the job of learning easier(instead of trying to learn the same weights/kernels at different locations again and again)
- But does that mean we can have only one kernel?
- No, we can have many such kernels but the kernels will be shared by all locations in the image
- This is called “weight sharing”

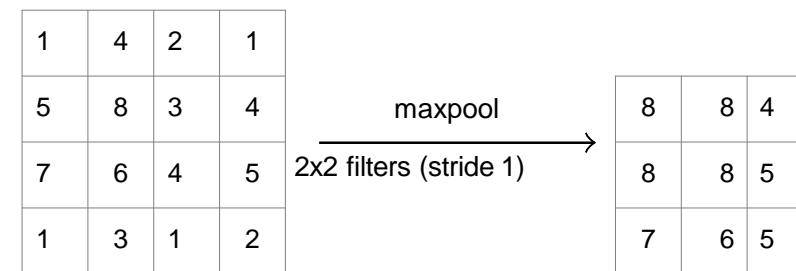
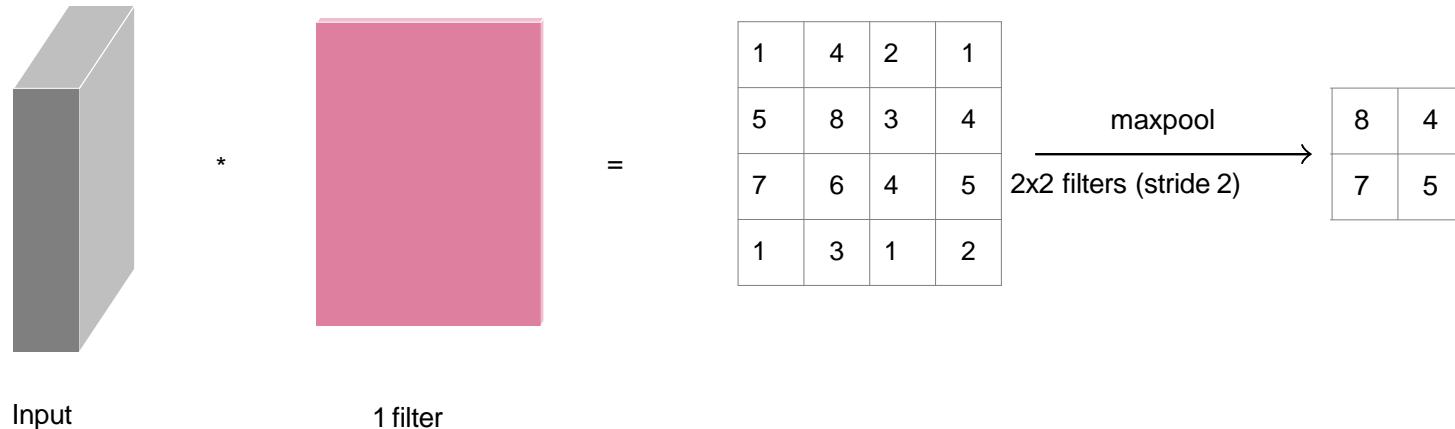


- It has alternate convolution and pooling layers
- What does a pooling layer do?
- Let us see

# Pooling



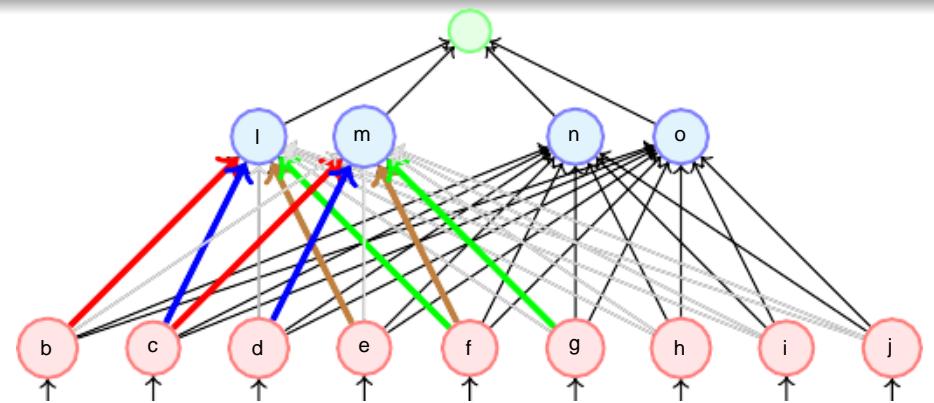
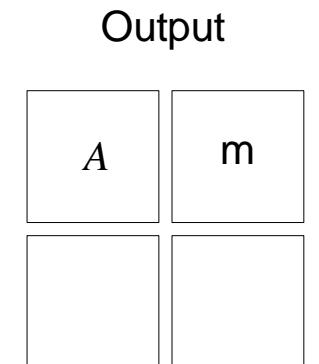
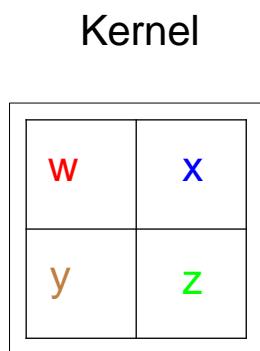
# Pooling



- Instead of max pooling we can also do average pooling

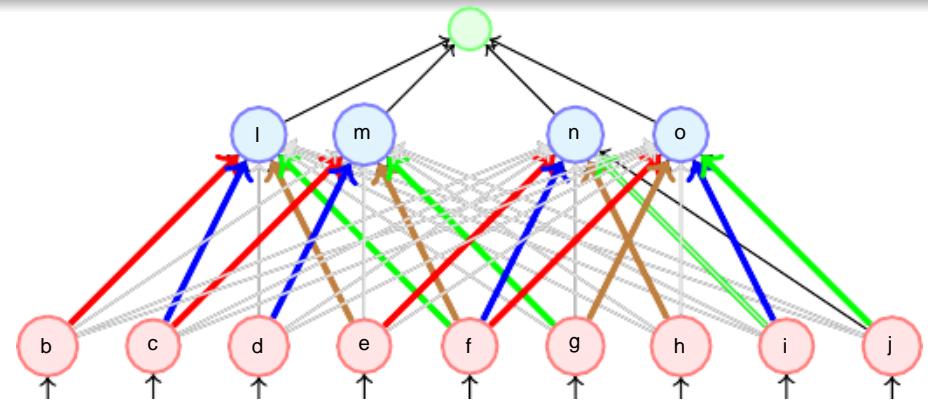
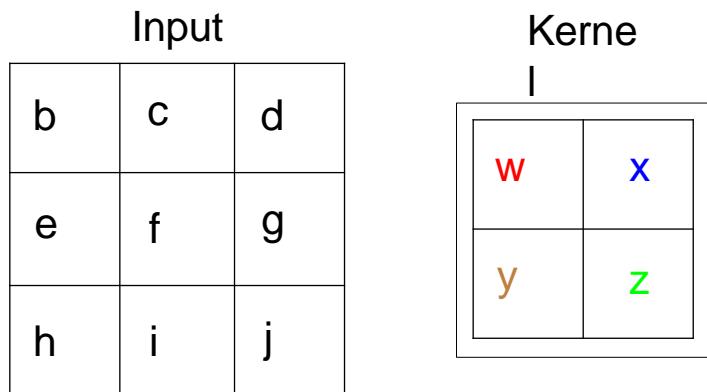
# Training CNN

Input		
b	c	d
e	f	g
h	i	j



- A CNN can be implemented as a feedforward neural network
- wherein only a few weights(in color) are active
- the rest of the weights (in gray) are zero

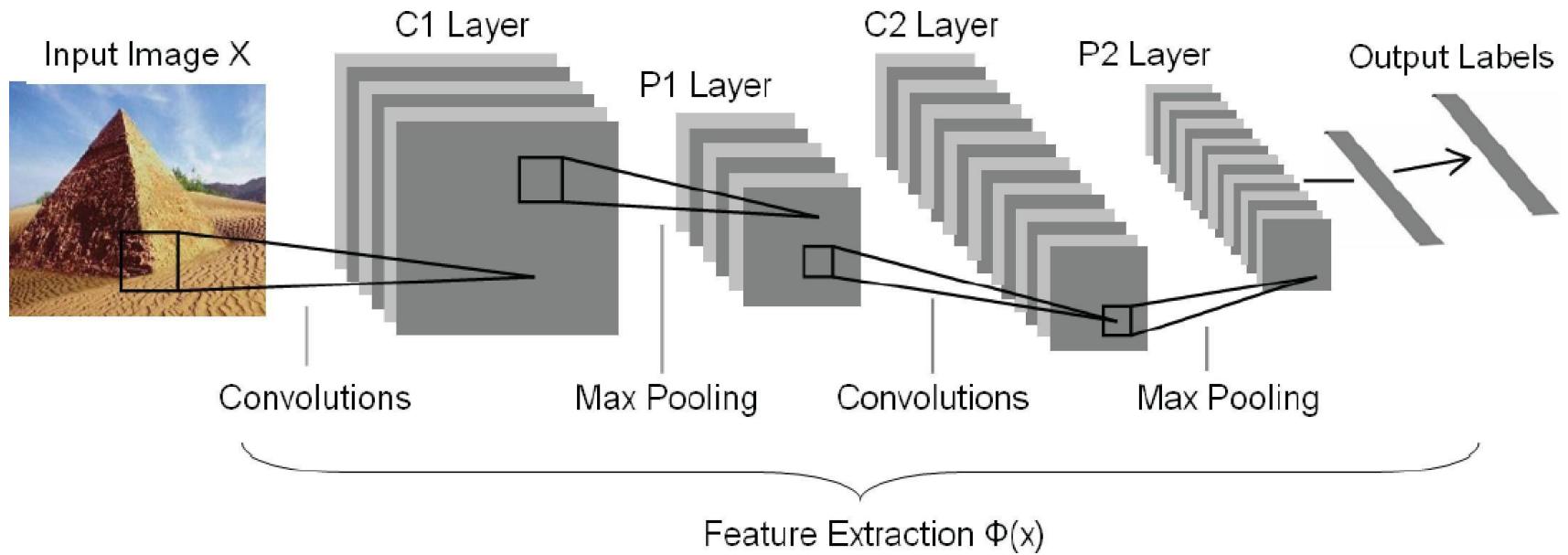
# Training CNN



- We can thus train a convolution neural network using backpropagation by thinking of it as a feedforward neural network with sparse connections
- A CNN can be implemented as a feedforward neural network
- wherein only a few weights(in color) are active
- the rest of the weights (in gray) are zero

# Convolutional Neural Nets for Image Recognition

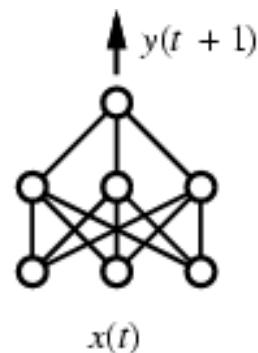
[Le Cun, 1992]



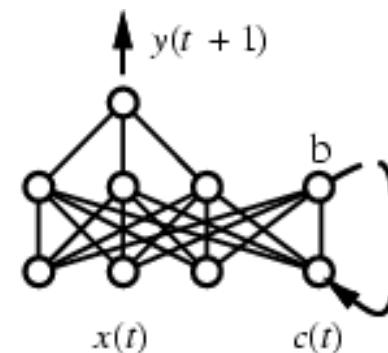
- specialized architecture: mix different types of units, not completely connected, motivated by primate visual cortex
- many shared parameters, stochastic gradient training
- very successful! now many specialized architectures for vision, speech, translation, ...

# Recurrent Networks: Time Series

- Suppose we want to predict next state of world
  - and it depends on history of unknown length
  - e.g., robot with forward-facing sensors trying to predict next sensor reading as it moves and turns
- Idea: use hidden layer in network to capture state history



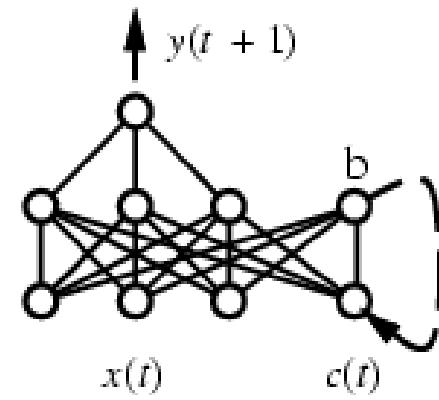
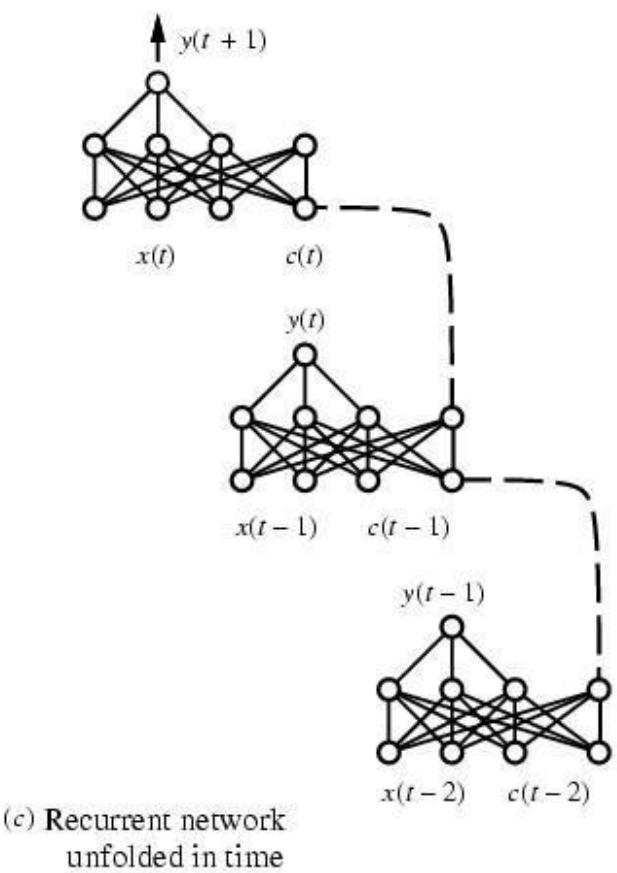
(a) Feedforward network



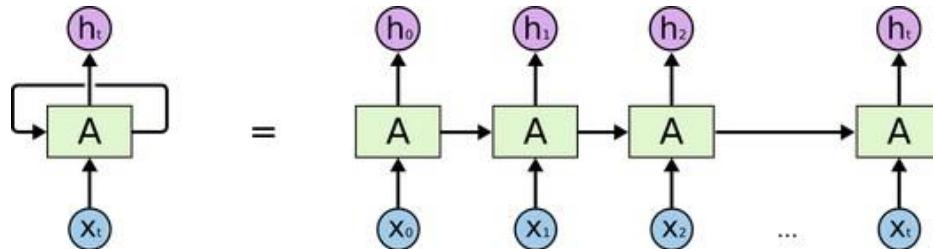
(b) Recurrent network

# Recurrent Networks on Time Series

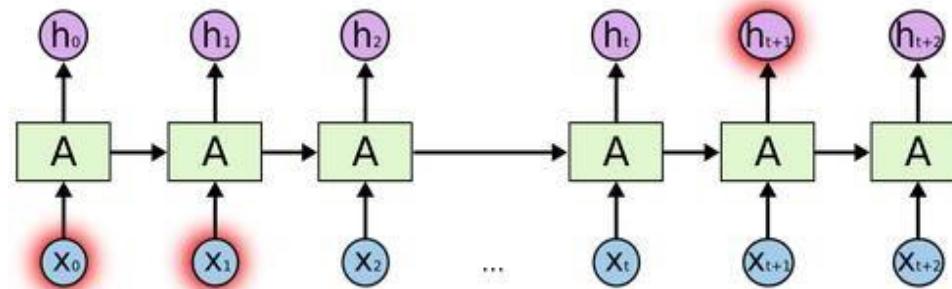
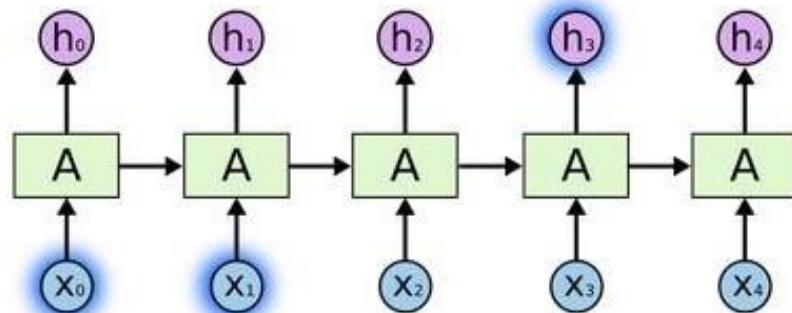
How can we train recurrent net??



# Recurrent Neural Networks



An unrolled recurrent neural network.



## Why RNNs:

Traditional neural networks can't have persistence.

For example, imagine you want to classify what kind of event is happening at every point in a movie. It's unclear how a traditional neural network could use its reasoning about previous events in the film to inform later ones.

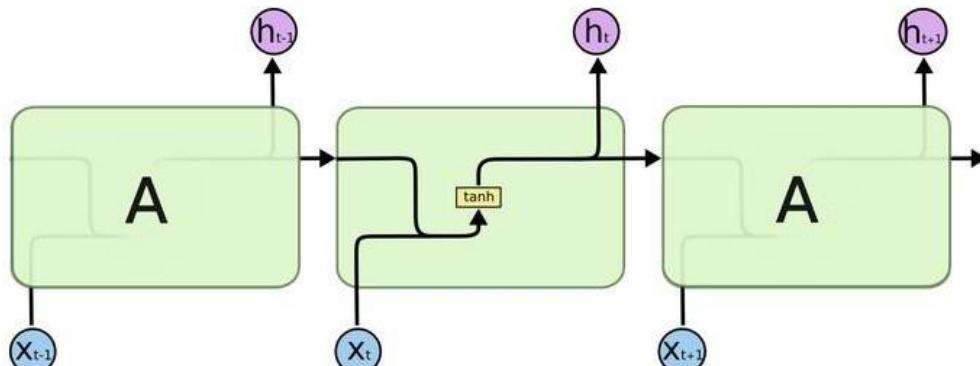
RNNs address this issue. They are networks with loops in them, allowing information to persist.

## Problem of Long Term Dependency:

Consider trying to predict the last word in the text "I grew up in France... I speak fluent *French*." Recent information suggests that the next word is probably the name of a language, but if we want to narrow down which language, we need the context of France, from further back. It's entirely possible for the gap between the relevant information and the point where it is needed to become very large.

Unfortunately, as that gap grows, RNNs become unable to learn to connect the information.

# Long Short Term Memory Networks

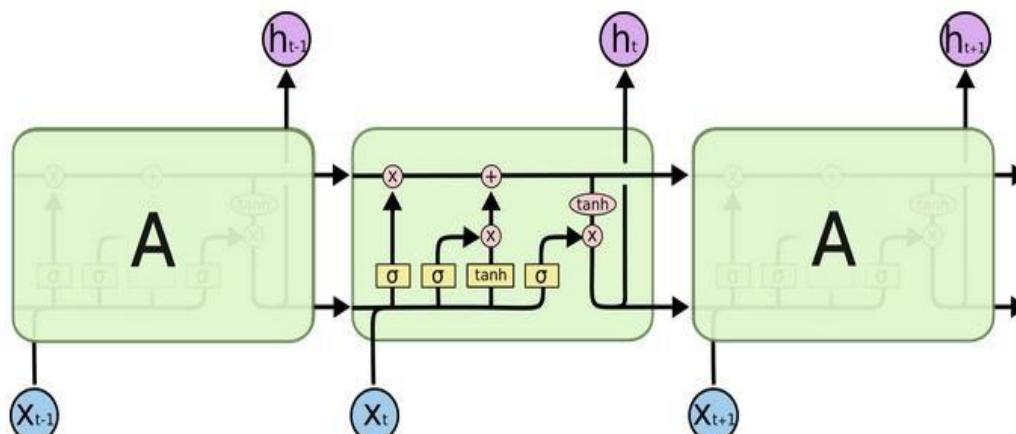


The repeating module in a standard RNN contains a single layer.

All recurrent neural networks have the form of a chain of repeating modules of neural network. In standard RNNs, this repeating module will have a very simple structure, such as a single tanh layer.

LSTMs also have this chain like structure, but the repeating module has a different structure. Instead of having a single neural network layer, there are four, interacting in a very special way.

- Forget gate layer
- Input gate layer
- Control gate layer
- Output gate layer



The repeating module in an LSTM contains four interacting layers.

The LSTM have the ability to remove or add information to the cell state, carefully regulated by structures called **gates**.

The sigmoid layer outputs numbers between 0 and 1, describing how much of each component should be let through. A value of 0 means “let nothing through,” while a value of 1 means “let everything through!”

# Artificial Neural Networks: Summary

- Highly non-linear regression/classification
- Hidden layers learn intermediate representations
- Potentially millions of parameters to estimate
- Stochastic gradient descent, local minima problems
- Deep networks have produced real progress in many fields
  - computer vision
  - speech recognition
  - mapping images to text
  - recommender systems
  - ...
- They learn very useful non-linear representations

# Good References for understanding Neural Network

Mitesh Khapra

<https://www.youtube.com/watch?v=yw8xwS15Pf4>

Visualization of CNN

<https://www.youtube.com/watch?v=cNBBNAxC8I4>