

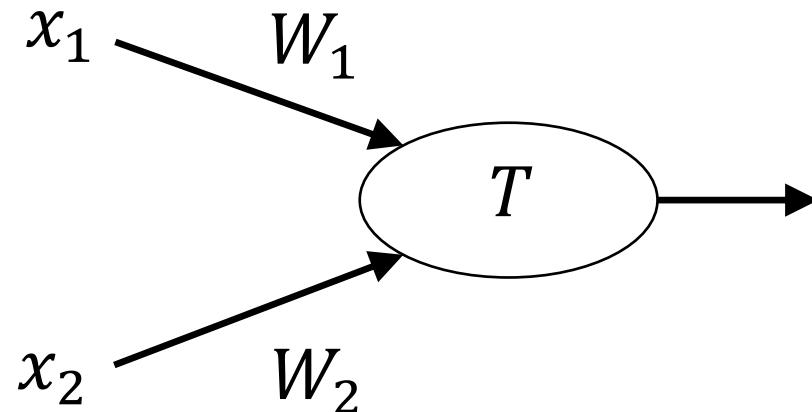
# **Neural Networks**

Chetan Arora

Disclaimer: The contents of these slides are taken from various publicly available resources such as research papers, talks and lectures. To be used for the purpose of classroom teaching, and academic dissemination only.



# Simple network (AND)



$$\text{output} = \begin{cases} 1 & \text{if } \sum_{i=1}^n W_i x_i > T \\ -1 & \text{otherwise} \end{cases}$$



# Simple network (AND)

U

$\omega_0$

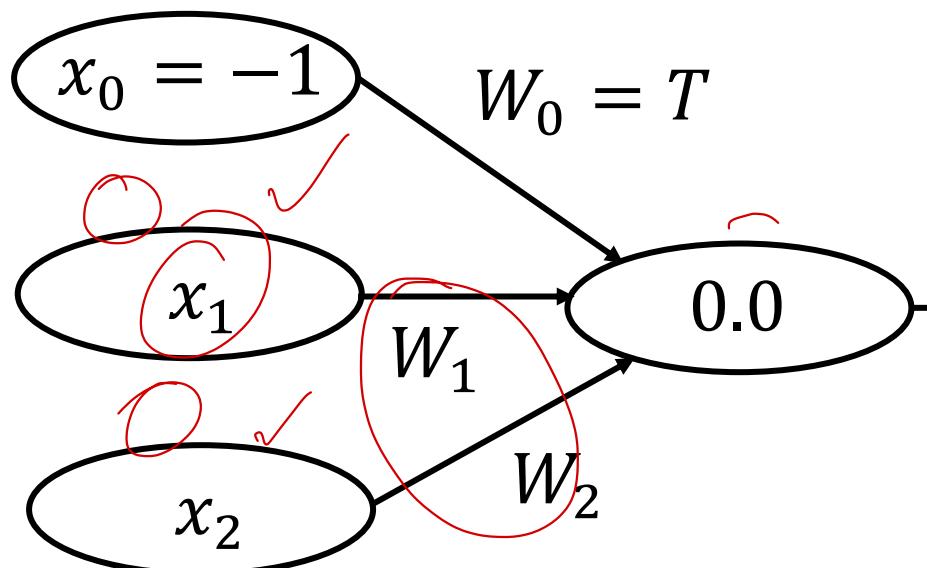
$\omega_1$

$\omega_2$

1

$v_r$

$n$

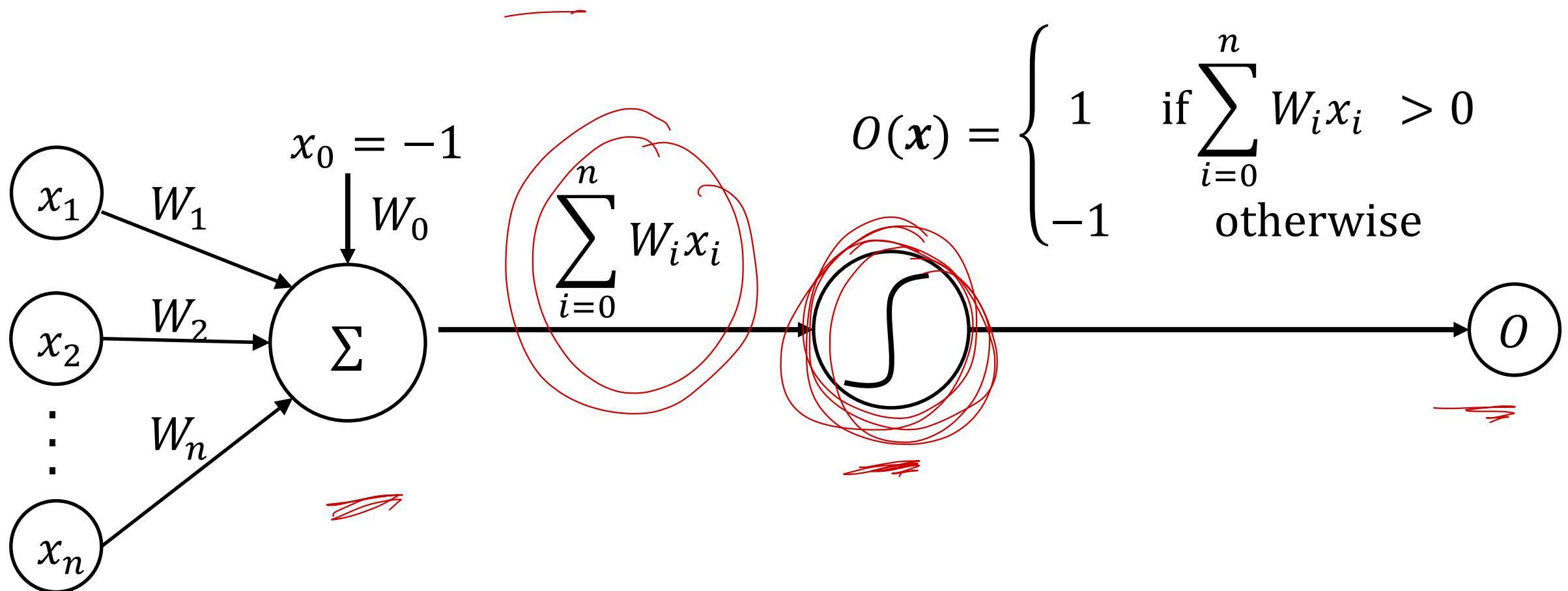


output = 
$$\begin{cases} 1 & \text{if } \sum_{i=0}^n W_i x_i > 0 \\ -1 & \text{otherwise} \end{cases}$$



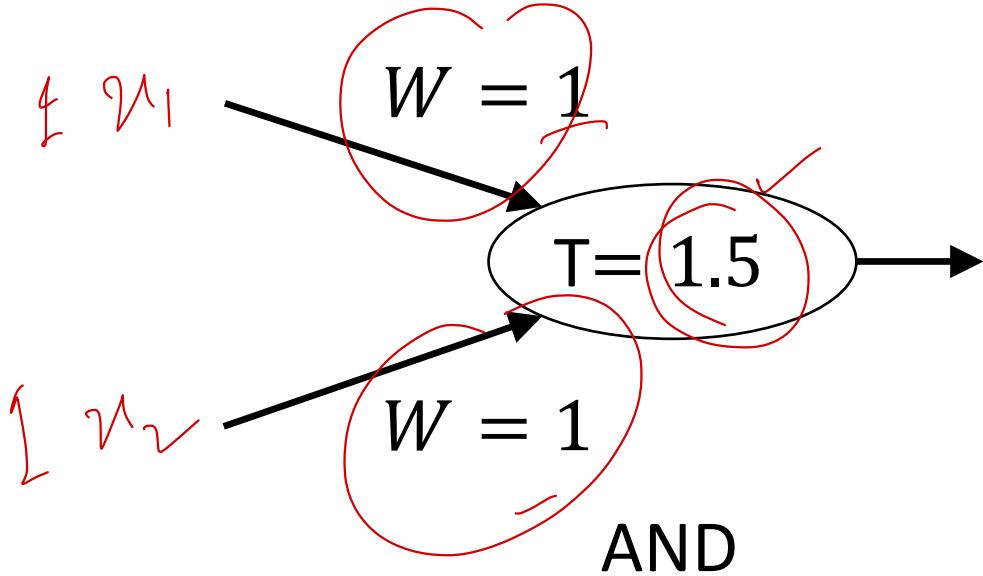
# Perceptron

- Linear threshold unit (LTU)



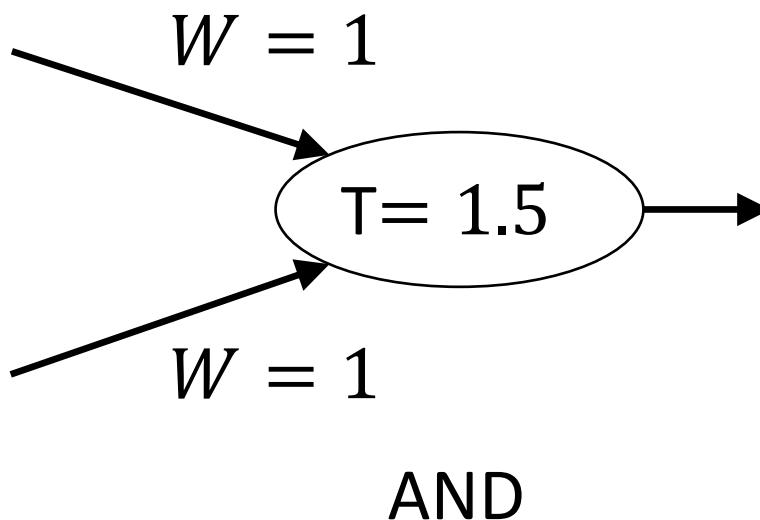


# Perceptron

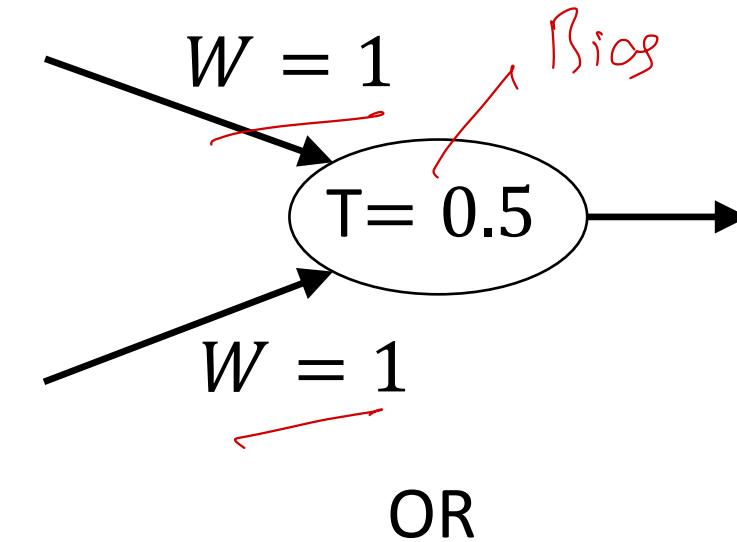




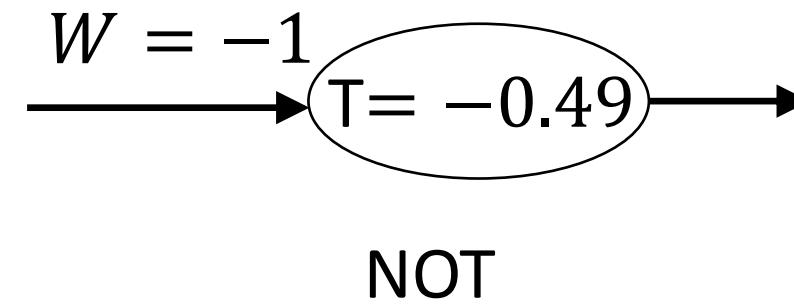
# Perceptron



AND



OR

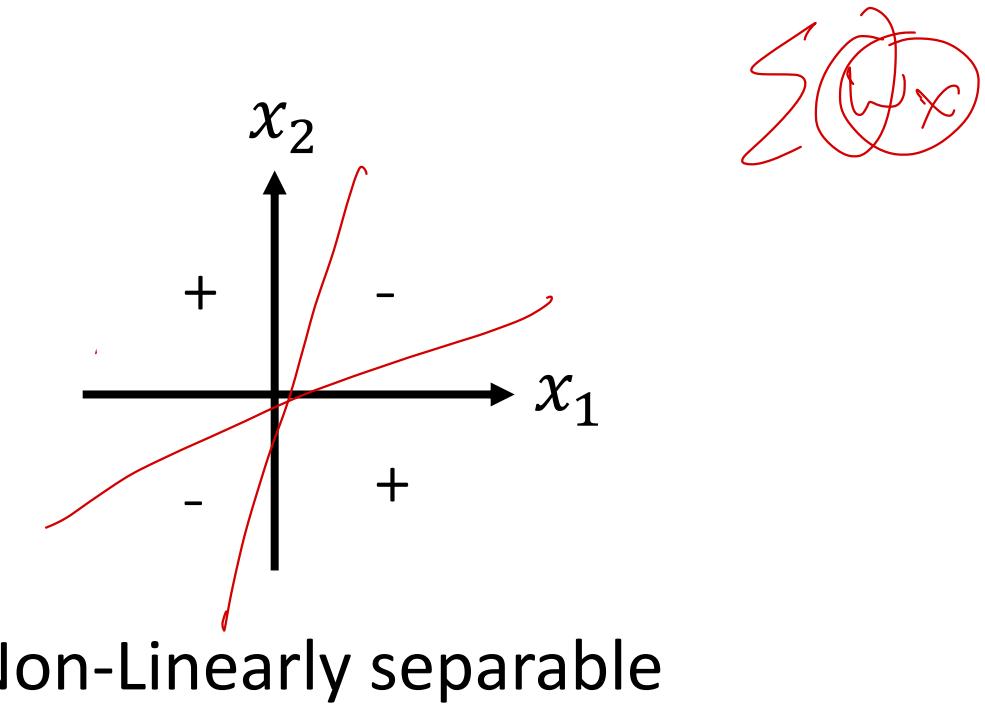
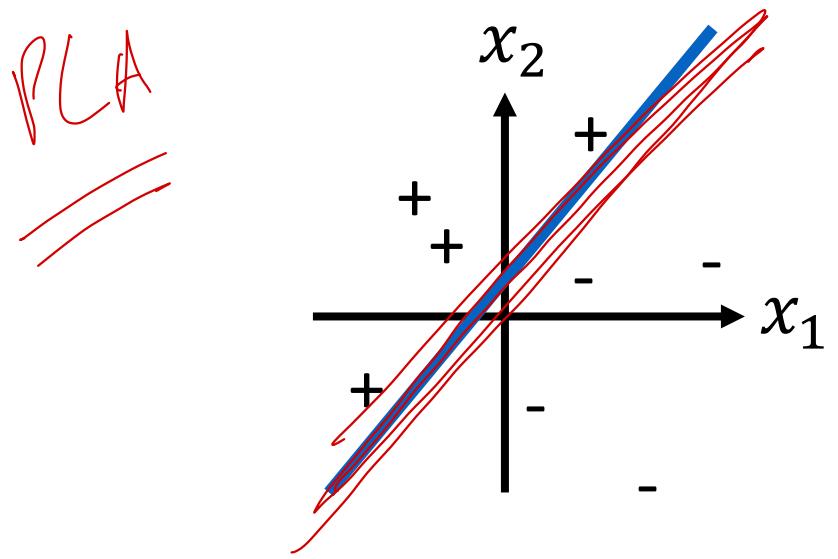


NOT



# Decision Surface of a Perceptron

- Perceptron can represent many useful functions: AND, OR, NOT etc.
- But functions that are not linearly separable (e.g.  $XOR$ ) are not representable.





# Solution in 1980s: Multilayer Perceptrons

- Removes many limitations of single-layer networks

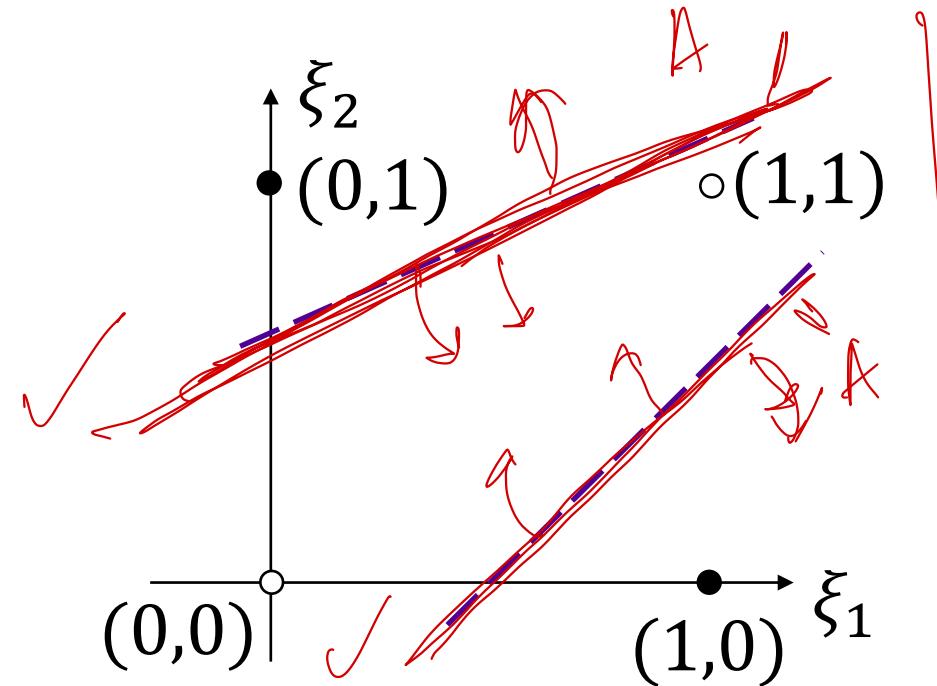
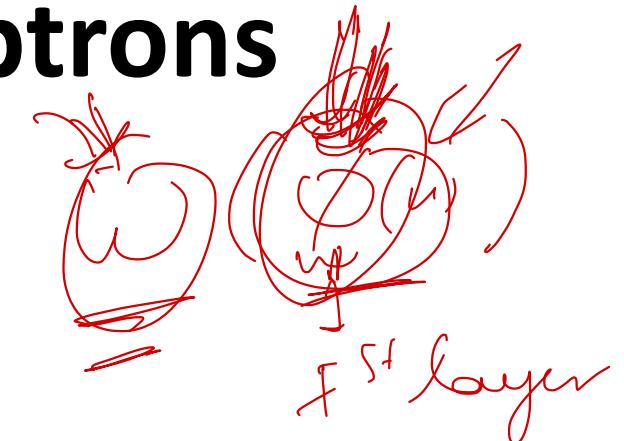
- Can solve  $XOR$

Interpretability

- **Exercise:** Draw a two-layer perceptron that computes the XOR function

- 2 binary inputs  $\xi_1$  and  $\xi_2$
- 1 binary output
- One “hidden” layer
- Find the appropriate weights and threshold

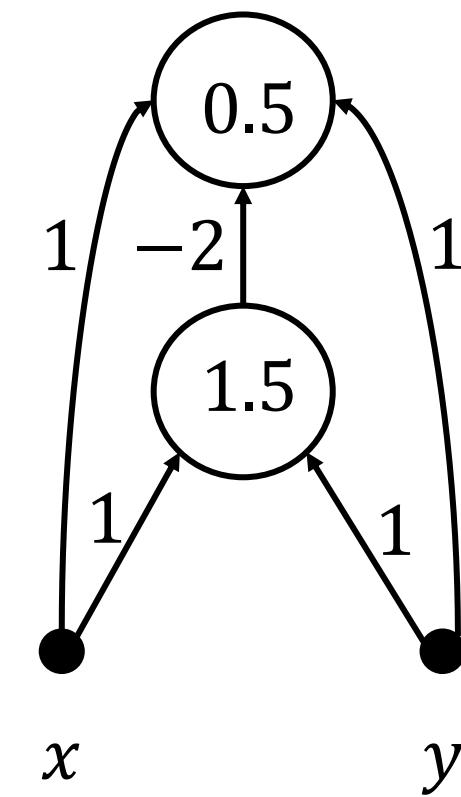
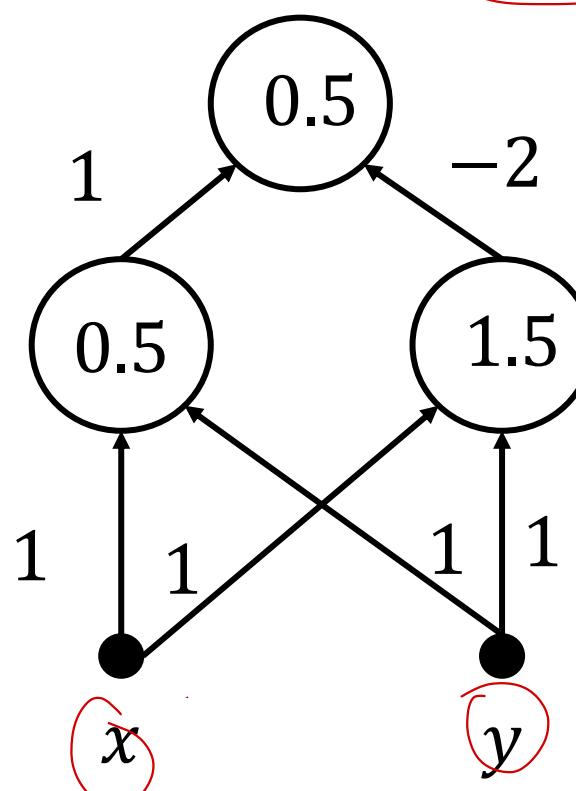
$$\omega \phi(r) \quad f(r)$$





# Solution in 1980s: Multilayer Perceptrons

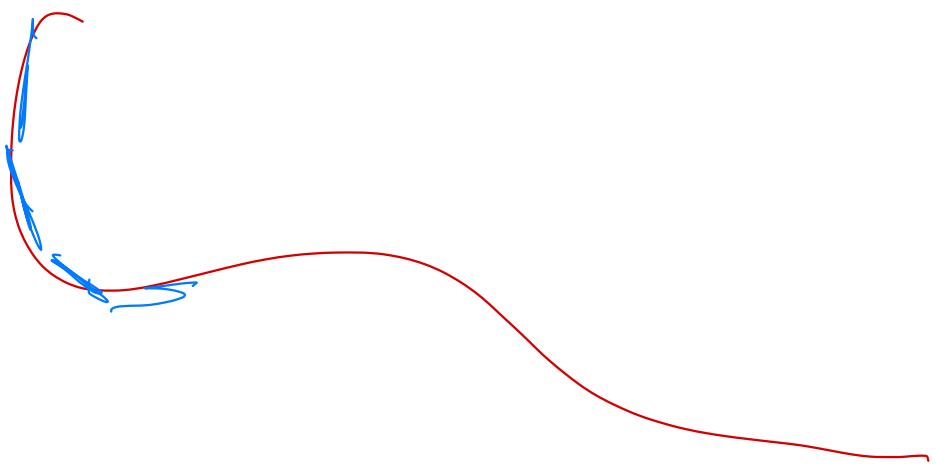
- Examples of two-layer perceptrons that compute  $XOR$
- E.g. Right side network
  - Output is 1 if and only if  $x + y - 2(x + y - 1.5) > 0 \rightarrow 0.5 > 0$





# Different Non-Linearly Separable Problems

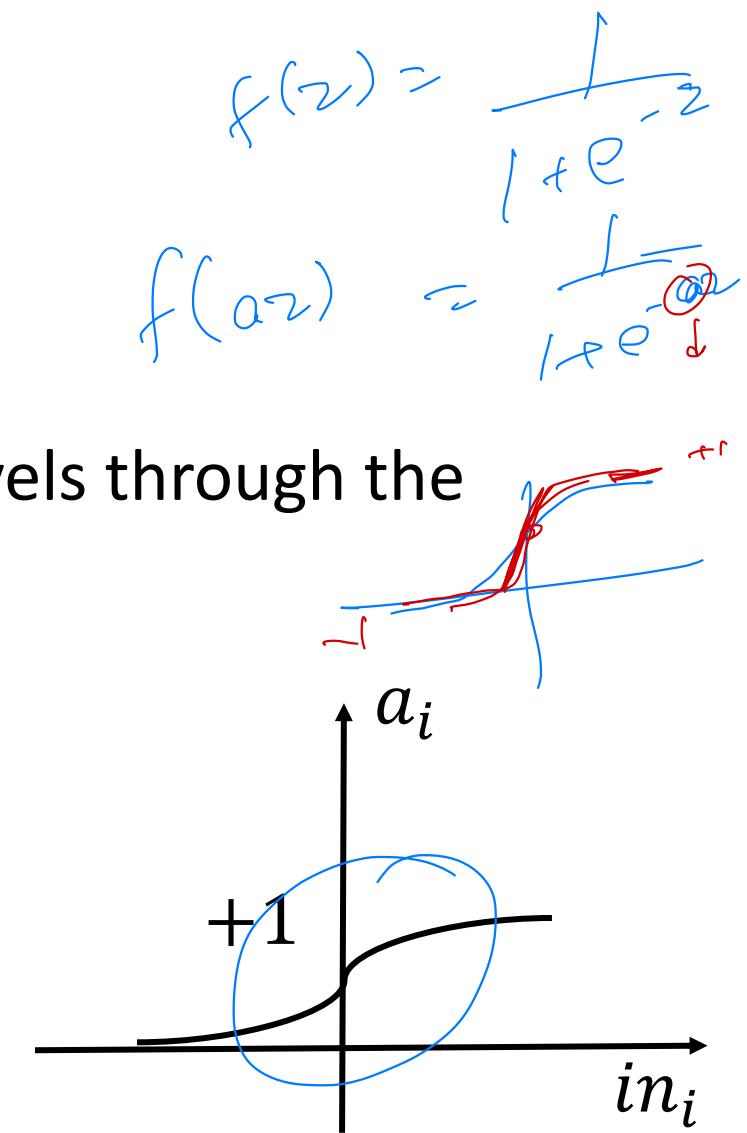
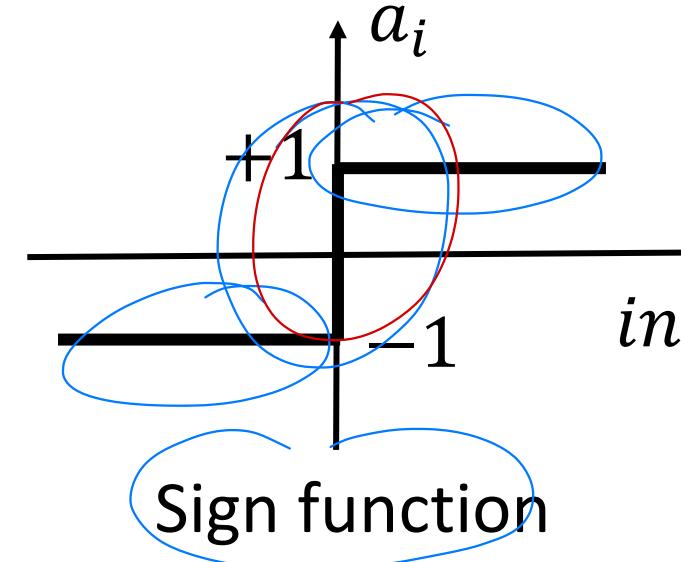
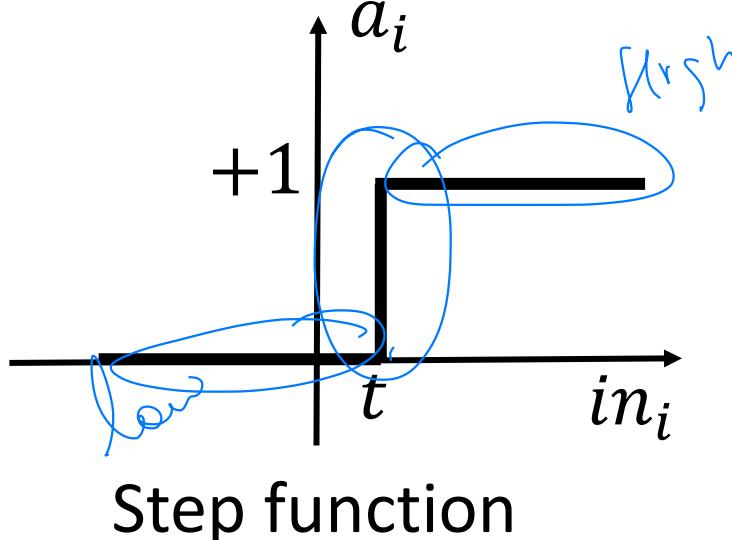
Structure	Types of Decision Regions	Exclusive-OR Problem	Classes with Meshed regions	Most General Region Shapes
Single-Layer	Half Plane Bounded By Hyperplane			
Two Layer	Convex Open Or Closed Regions			
Three Layer	Arbitrary (Complexity Limited by No. of Nodes)			





# Activation functions

- Transforms neuron's input into output.
- Features of activation functions:
  - A squashing effect is required
    - Prevents accelerating growth of activation levels through the network.
  - Simple and easy to calculate



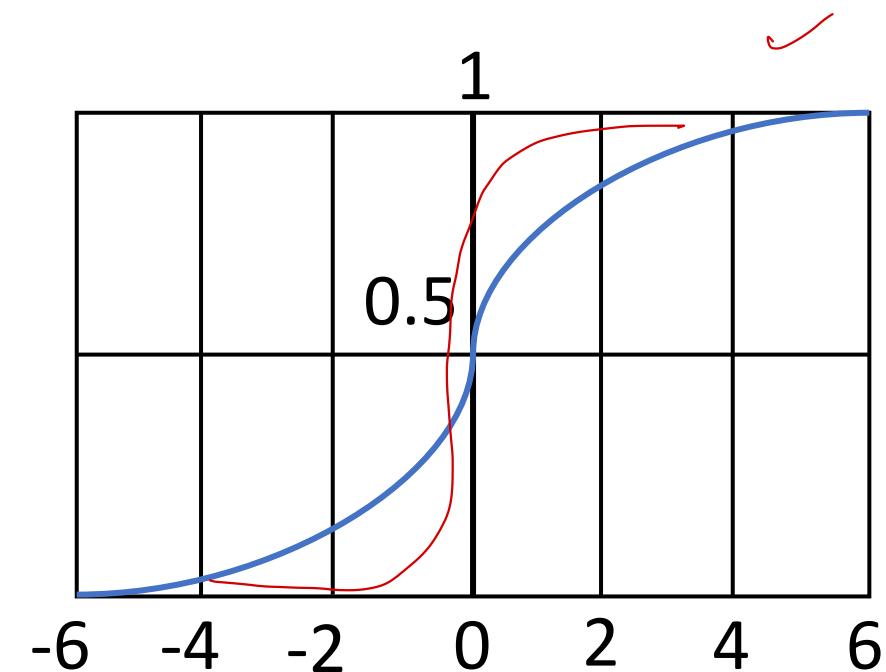
Sigmoid function



# Standard activation functions

- The hard-limiting threshold function
  - Corresponds to the biological paradigm
    - either fires or not
- Sigmoid functions ('S'-shaped curves)
  - The logistic function
  - The hyperbolic tangent (symmetrical)
  - Both functions have a simple differential
  - Only the shape is important

$$\phi(x) = \frac{1}{1 + e^{-x}}$$



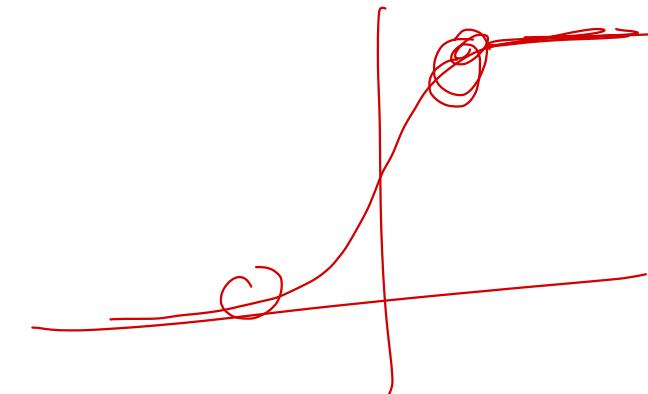


# Standard activation functions

✓

- tanh

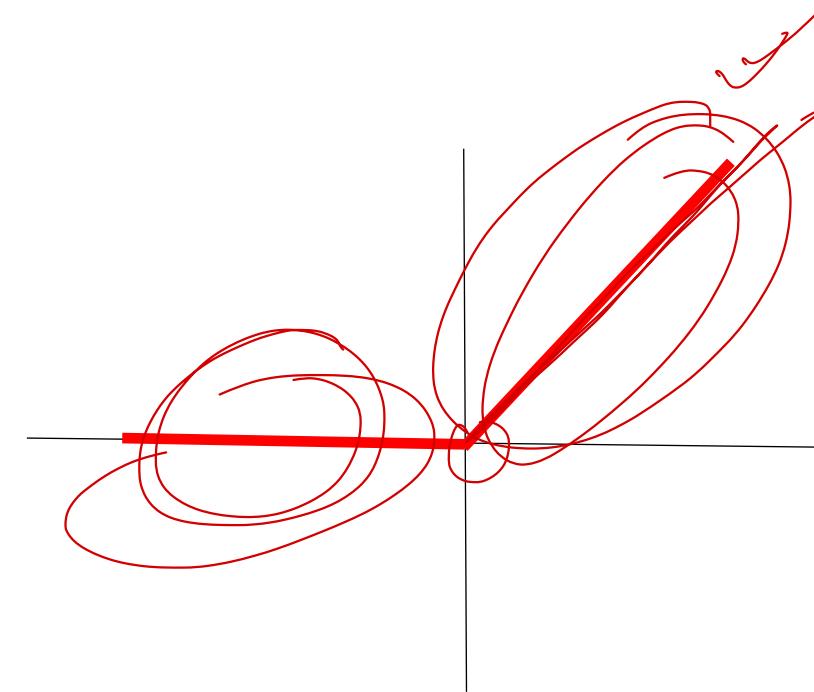
$$\frac{\exp(z) - \exp(-z)}{\exp(z) + \exp(-z)}$$



✓

- ReLU

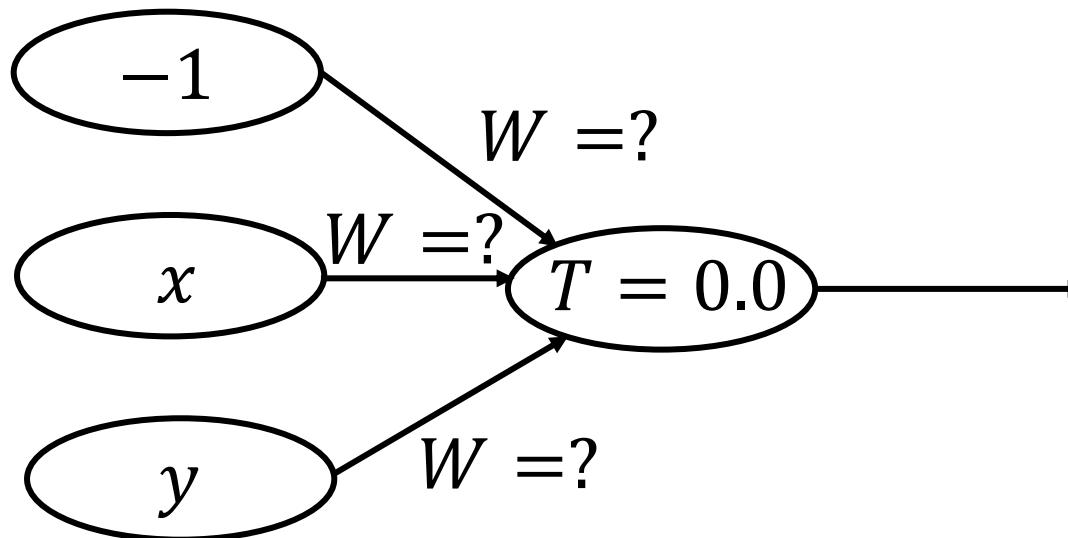
$$relu(z) = \max(0, z)$$





# Training Perceptrons

- What are the weight values?



For AND

A	B	Output
0	0	0
0	1	0
1	0	0
1	1	1



# Training Algorithms

- Perceptron Learning Algorithm: for linearly separable samples.
- Backpropagation: Coming soon.

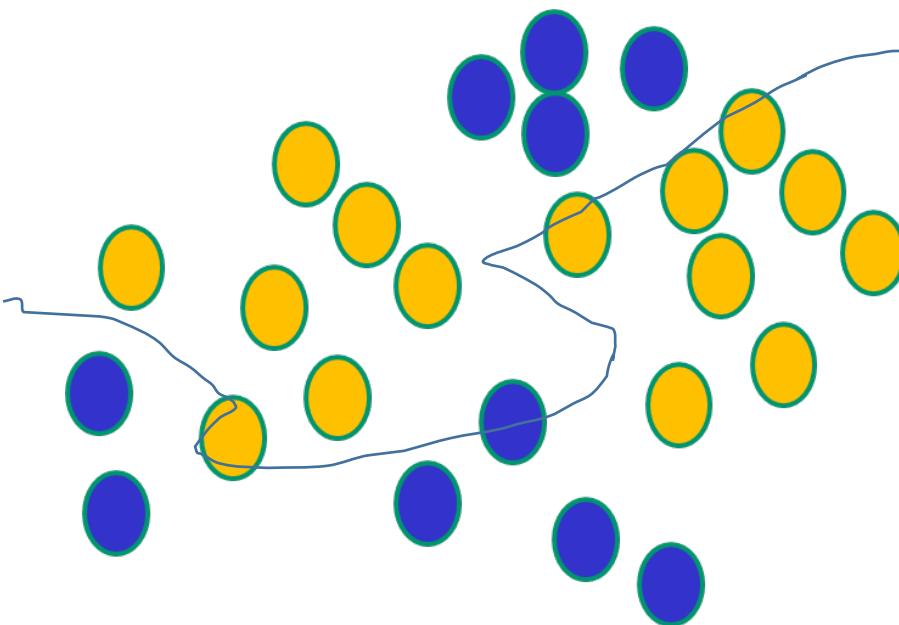
# Neural Networks as Feature detectors

$$\omega(\phi(v))$$



# Neural Networks: Two Perspectives

Use nonlinear  $f(x)$  to draw complex boundaries, but keep the data unchanged

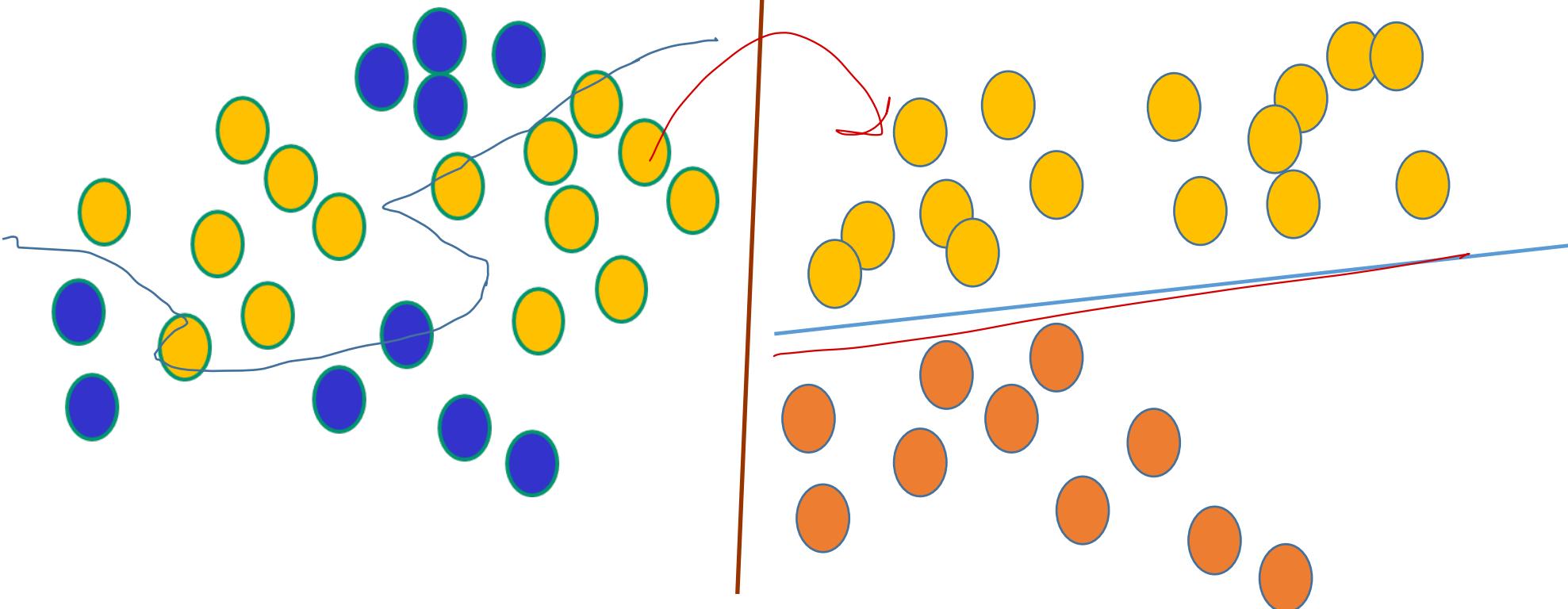




# Neural Networks: Two Perspectives

Use nonlinear  $f(x)$  to draw complex boundaries, but keep the data unchanged

Find linear boundaries only, but transform the data first in a way that makes linear boundaries OK





# Neural Networks as Feature detectors

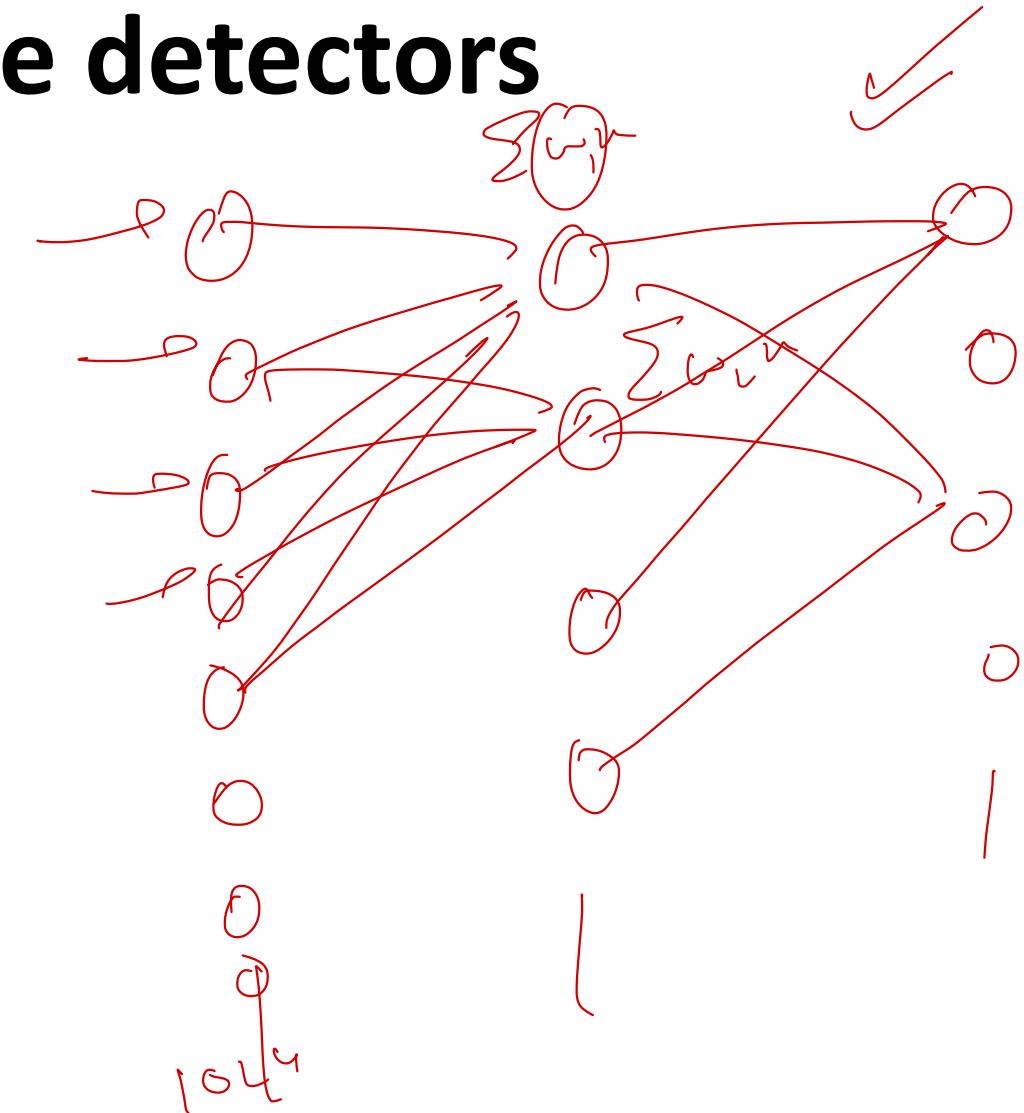
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9

MLP

Handwritten images

MNIST

72x72





# Neural Networks as Feature detectors

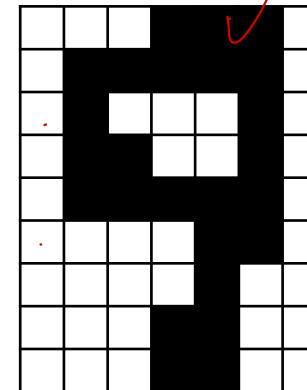
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9

Handwritten images

■ +1

□ -1

9



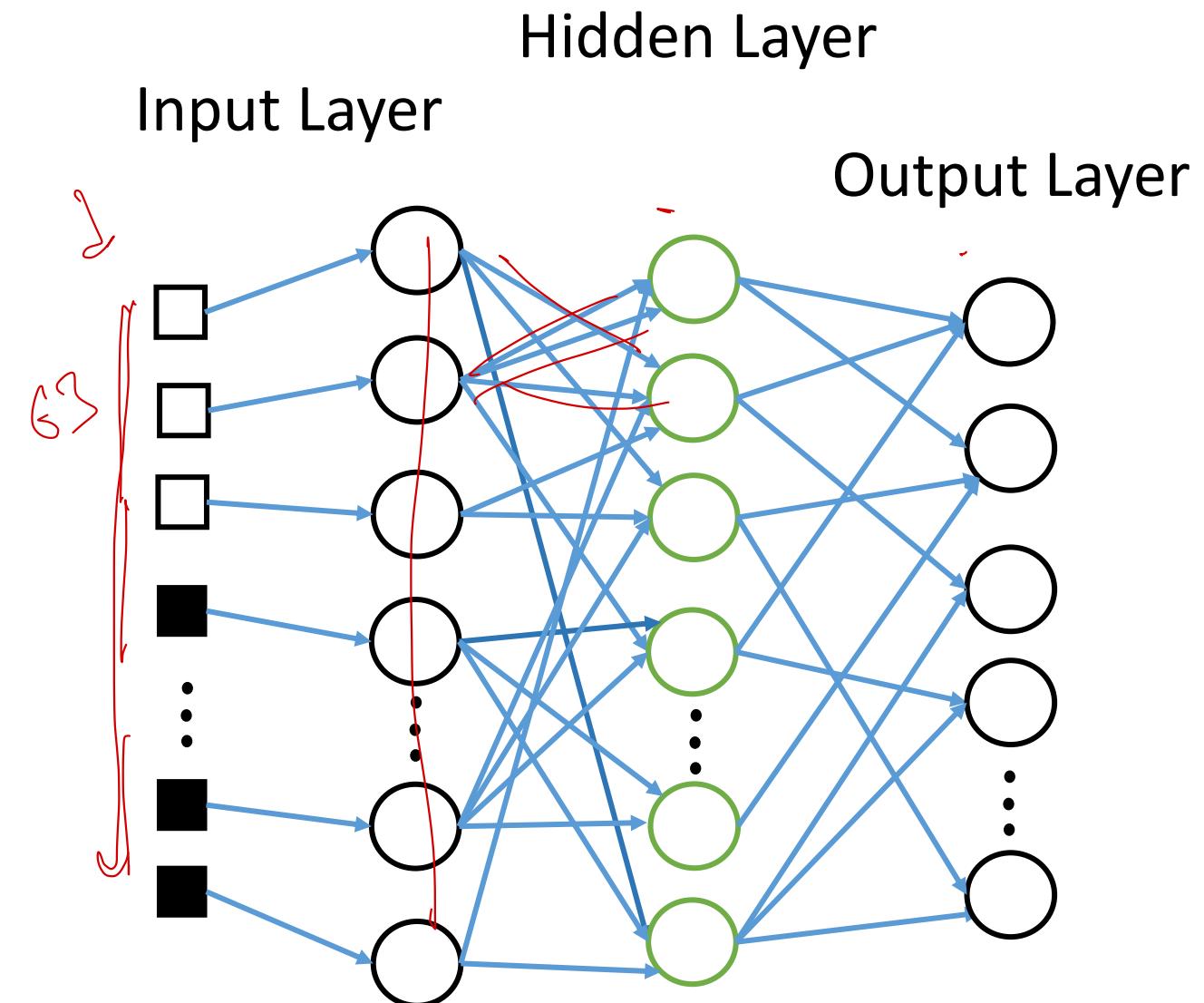
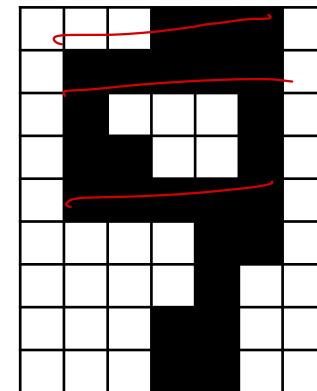


# Neural Networks as Feature detectors

0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9

Handwritten images

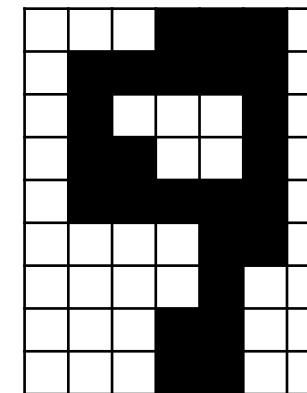
■ +1  
□ -1



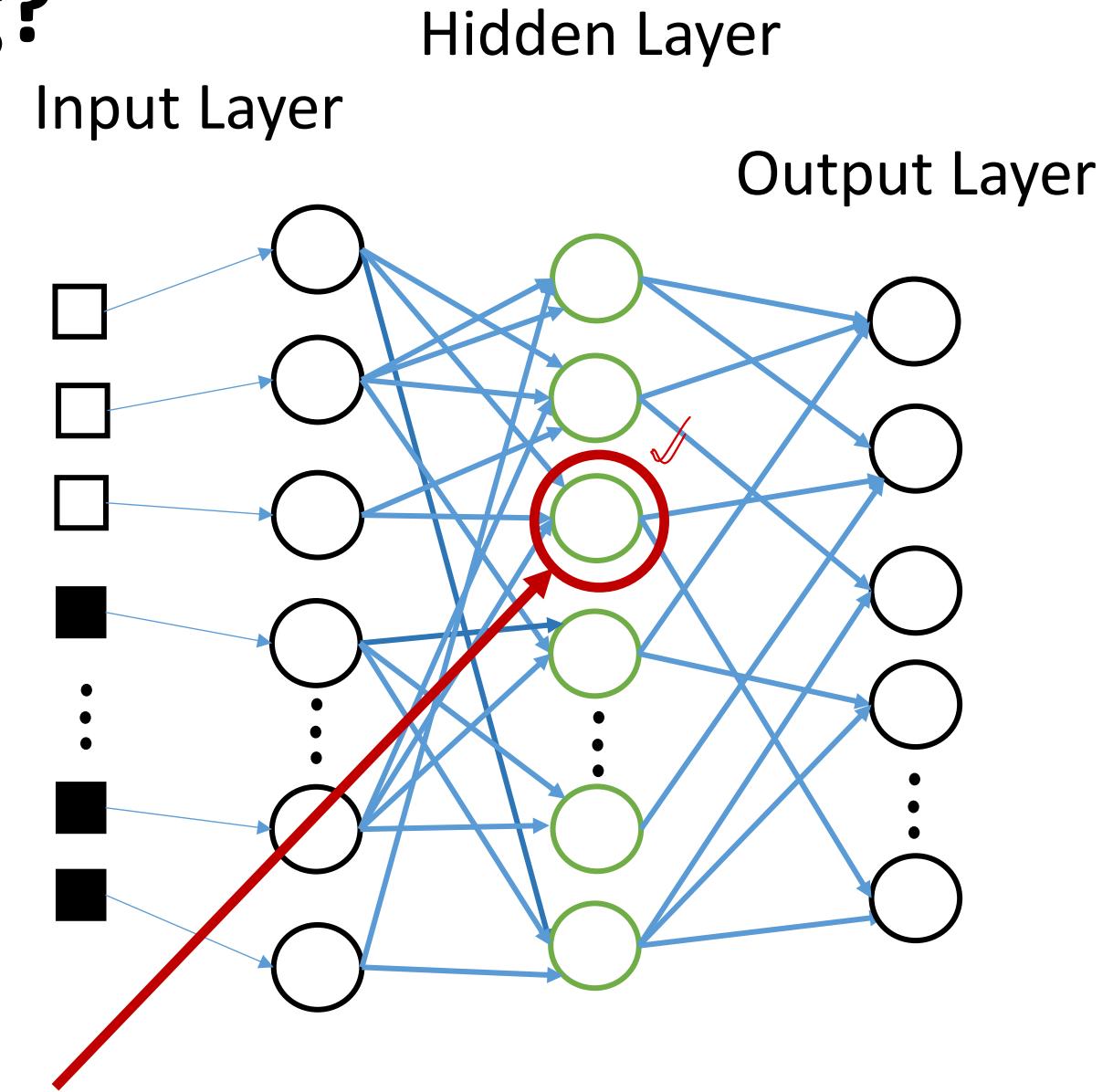


# What is this unit doing?

0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9

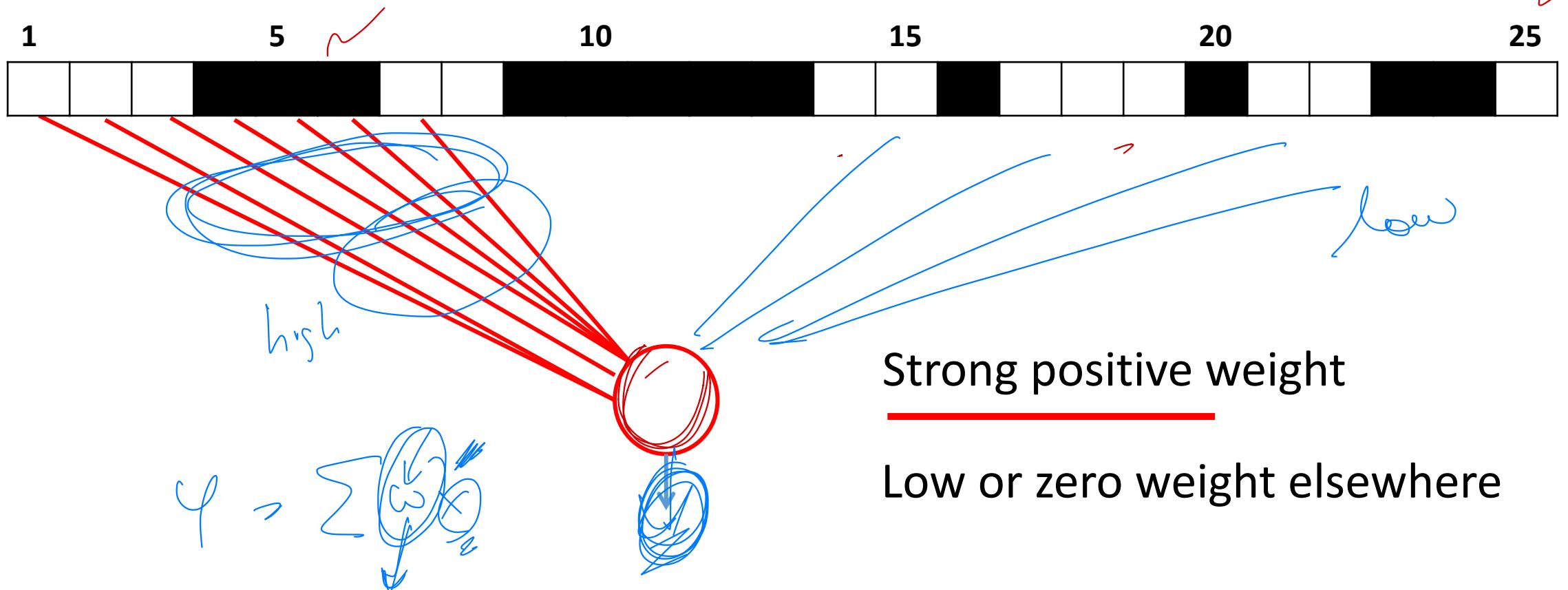


Handwritten images



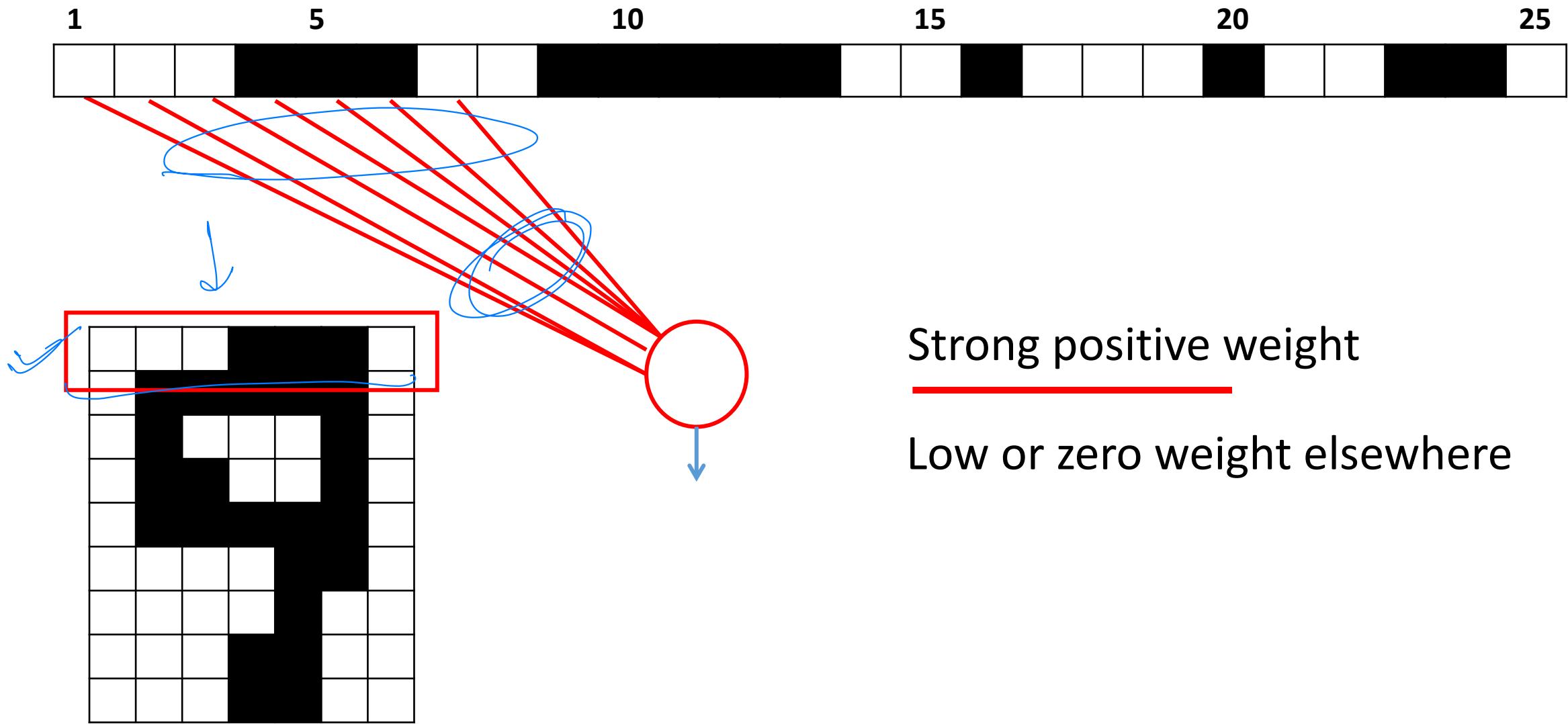


# What does this unit detect?



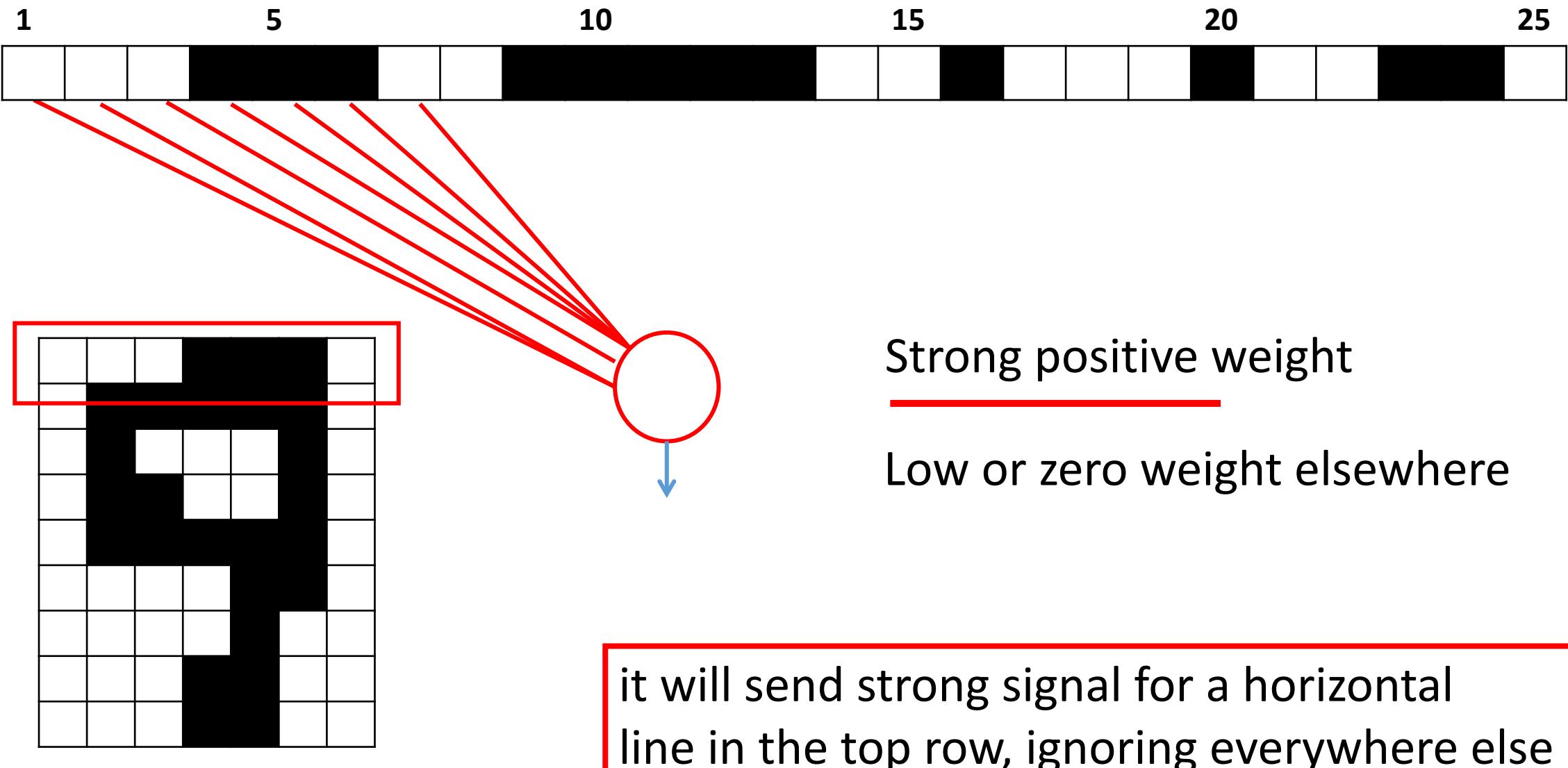


# What does this unit detect?



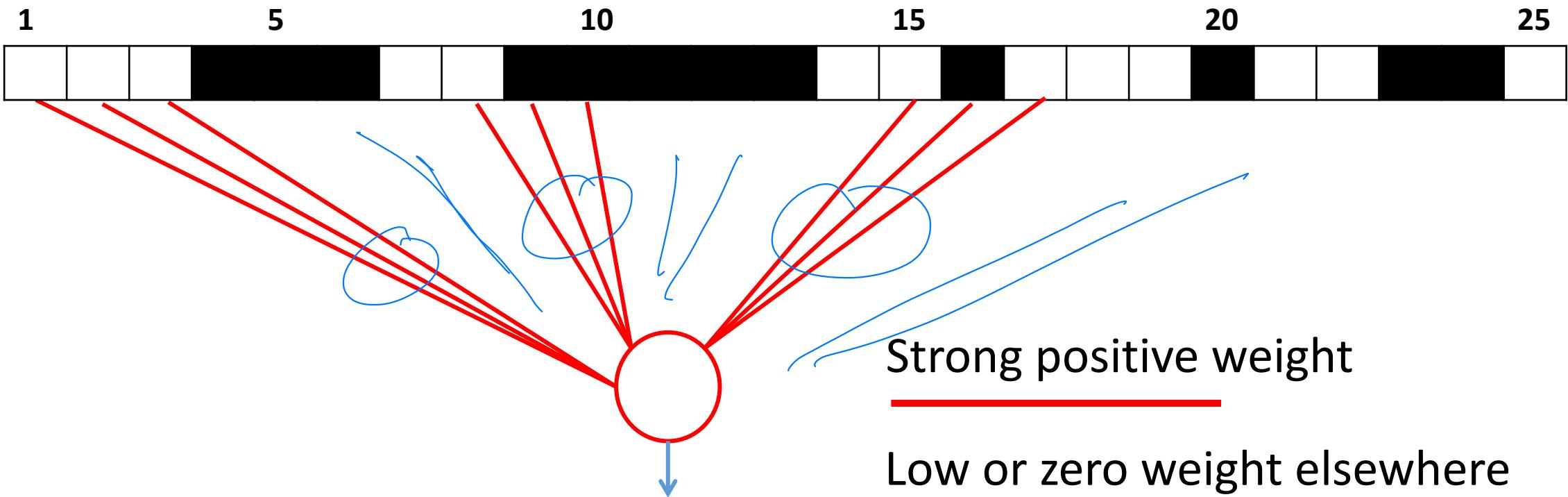


# What does this unit detect?



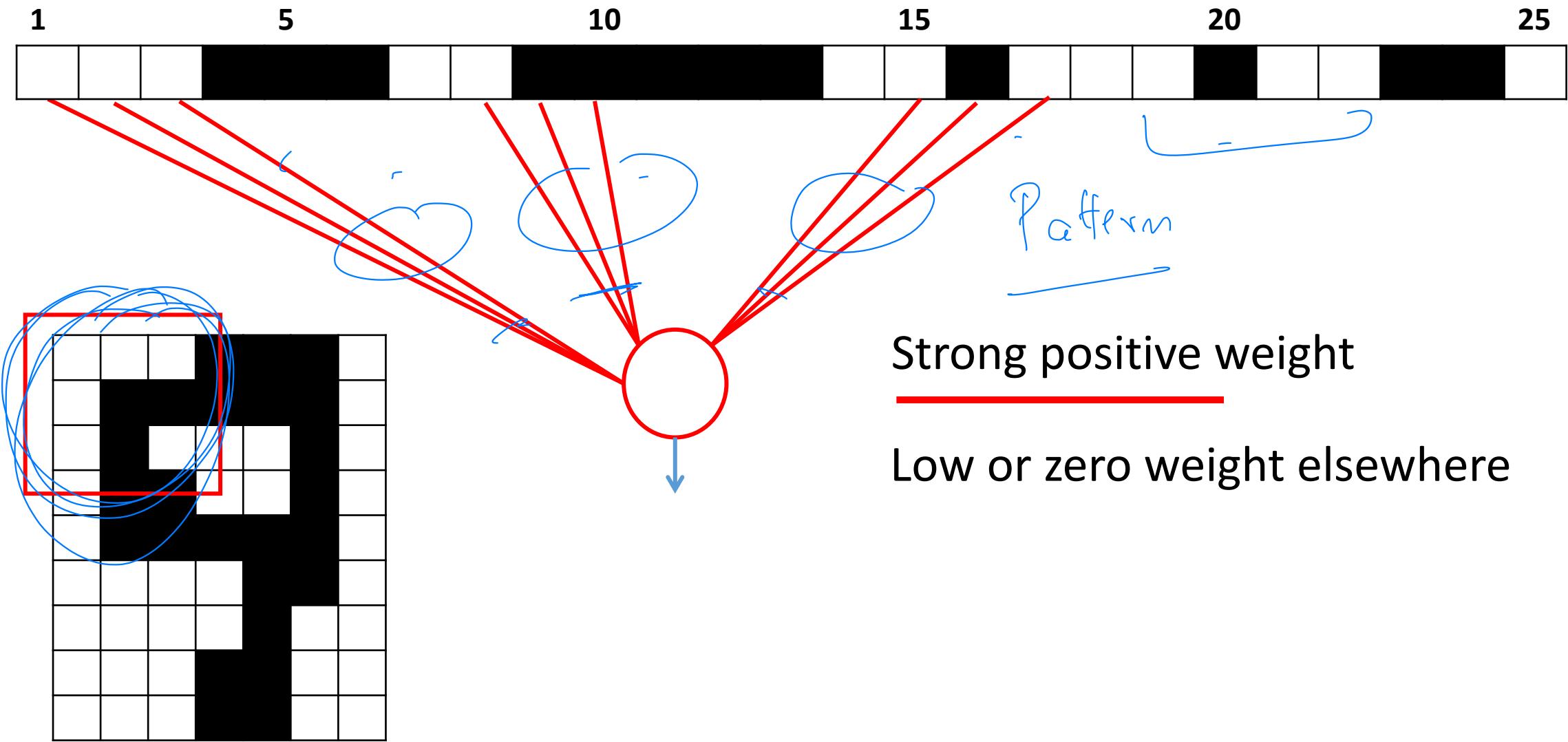


# What does this unit detect?



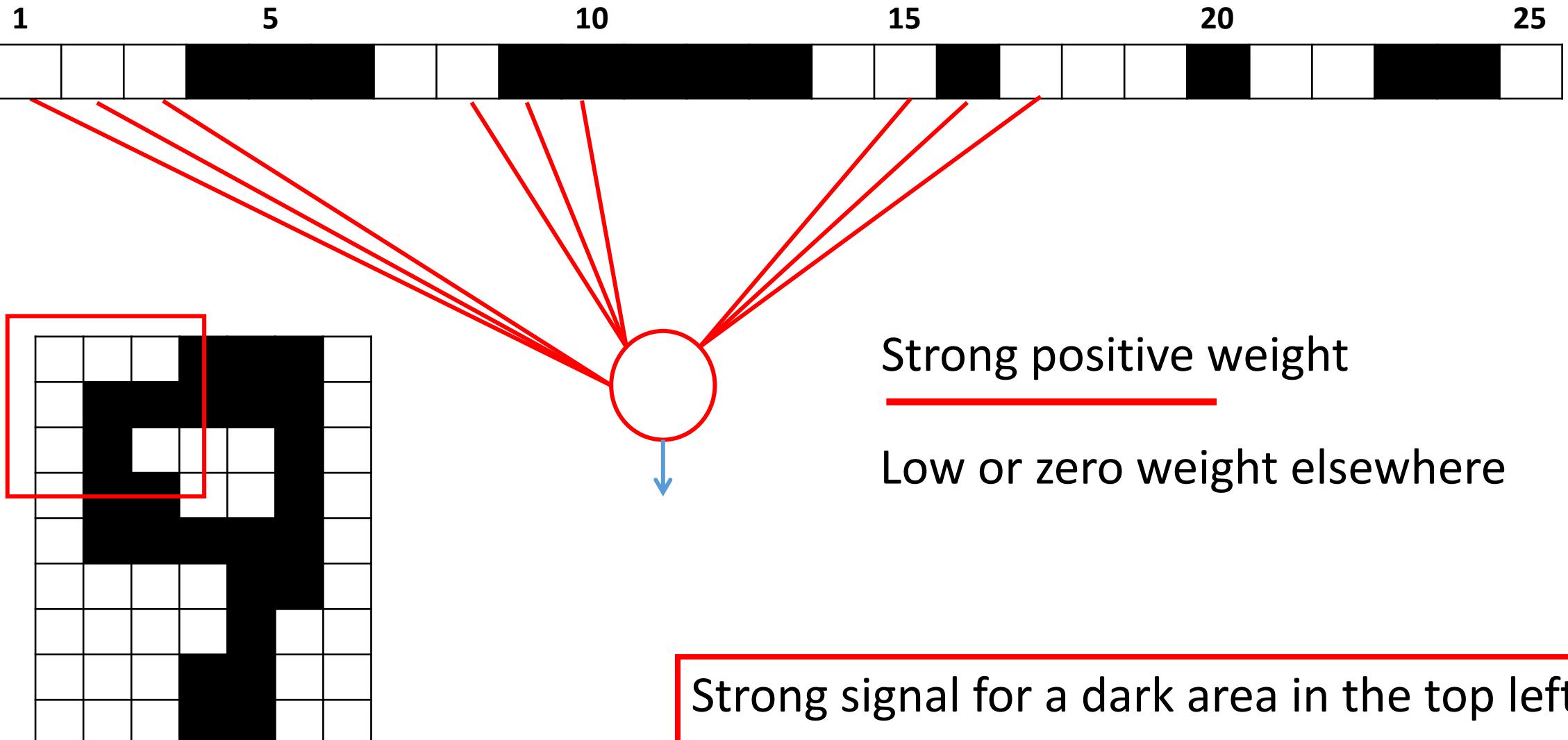


# What does this unit detect?



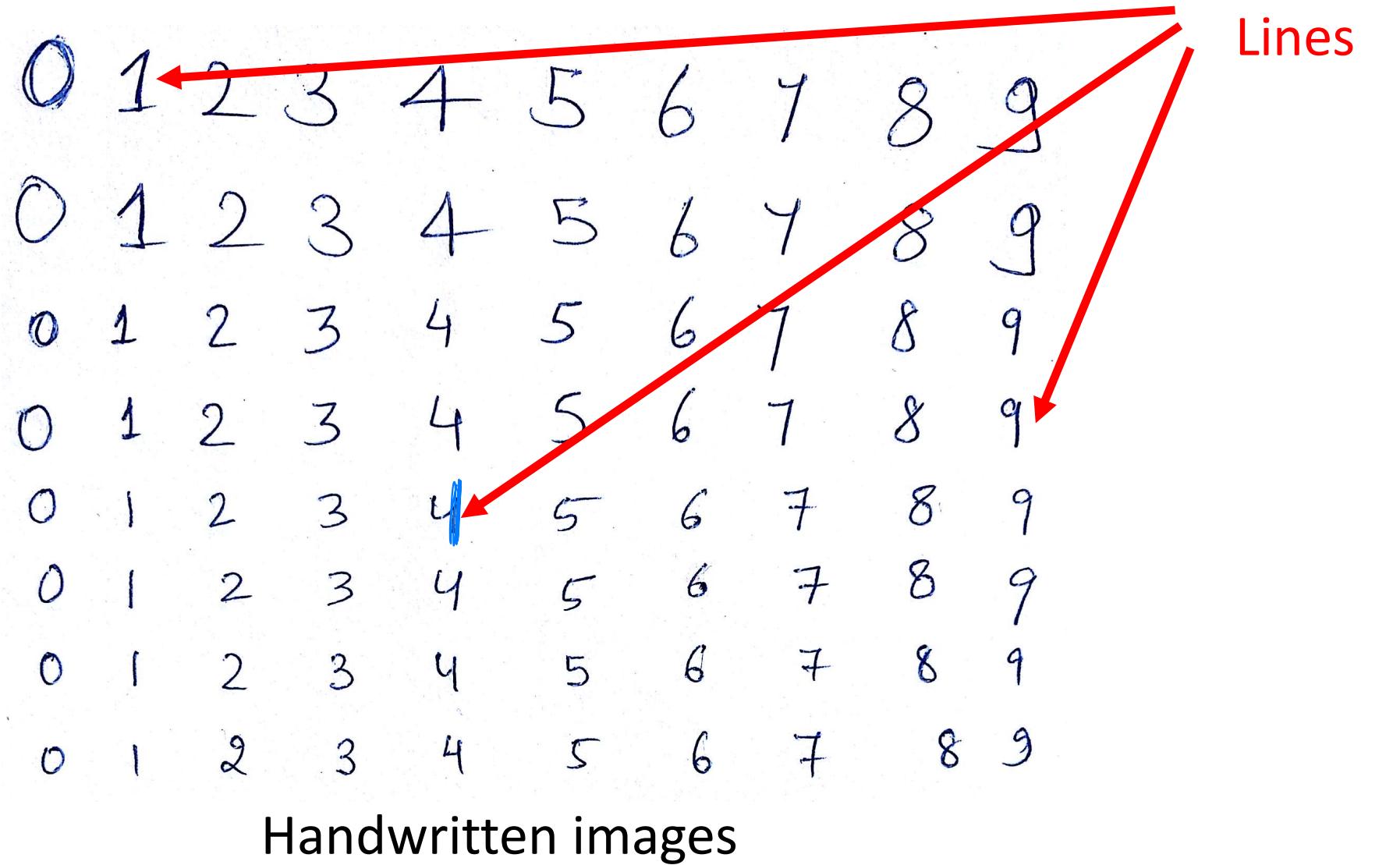


# What does this unit detect?



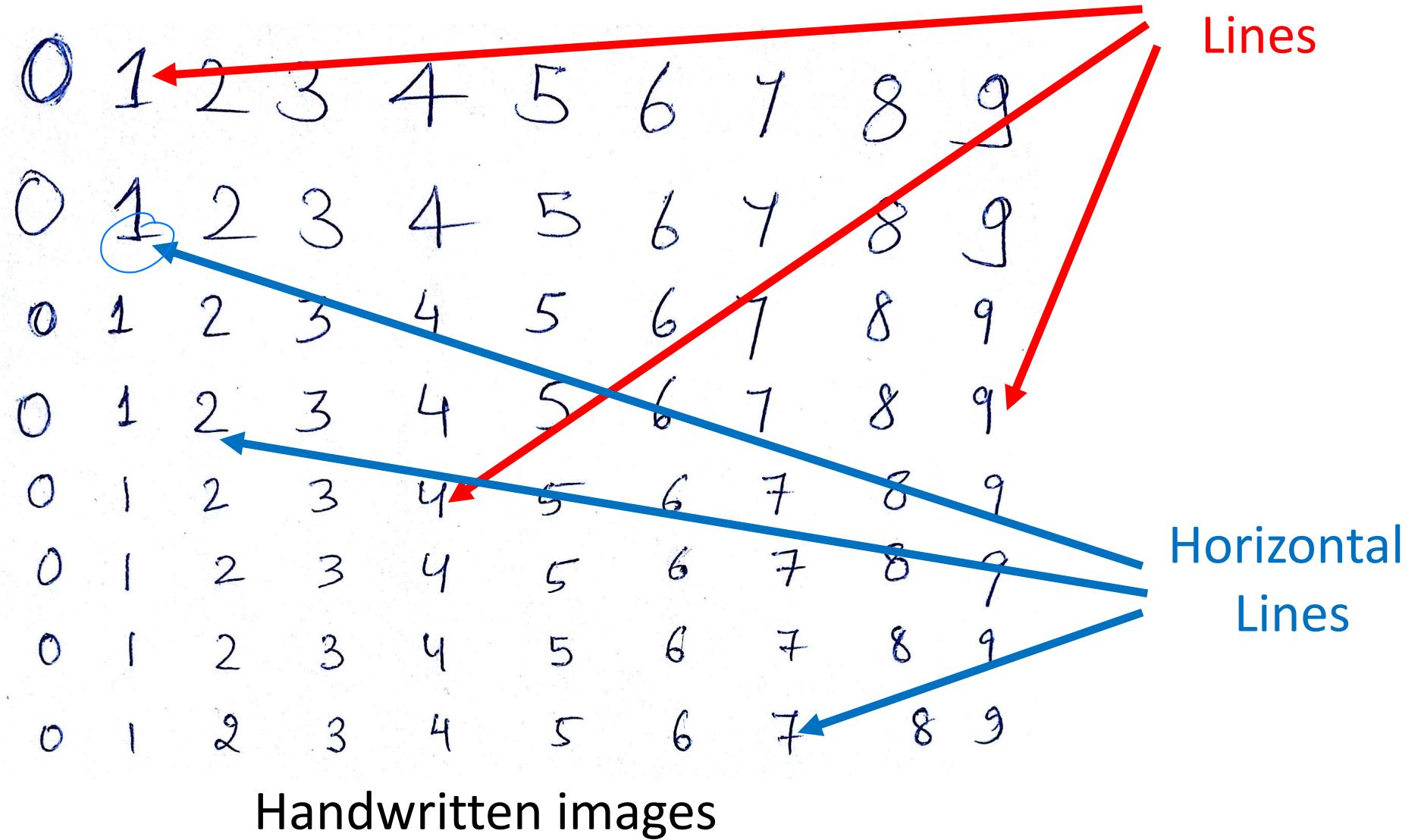


# Perceptron as a Feature Detector





# Perceptron as a Feature Detector





# Perceptron as a Feature Detector

Two blue outlines of eggs are shown. The top egg is whole, and the bottom egg is cut open, revealing a small yellow yolk inside.

# Small Loops

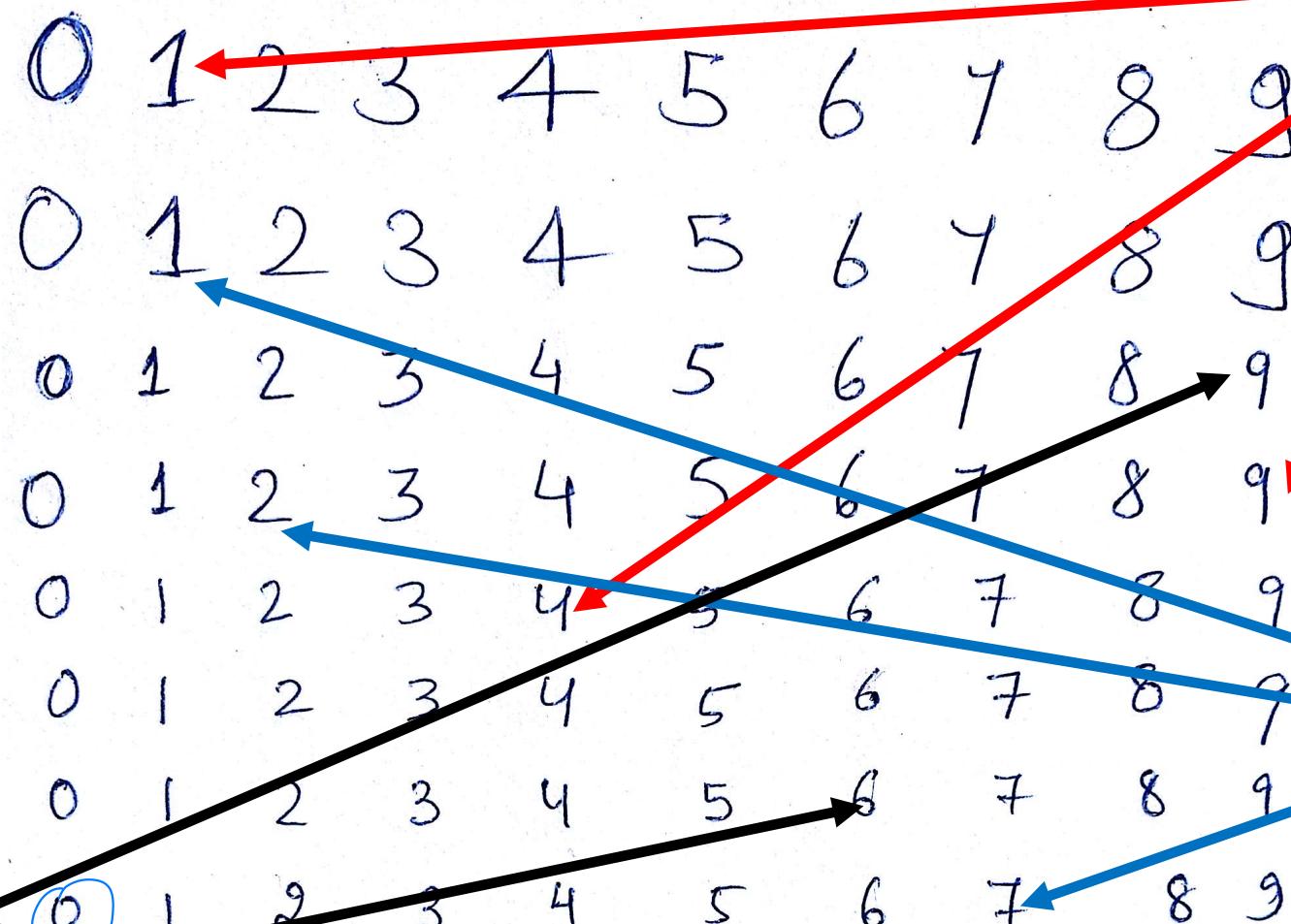
# Handwritten images

# Vertical Lines

7

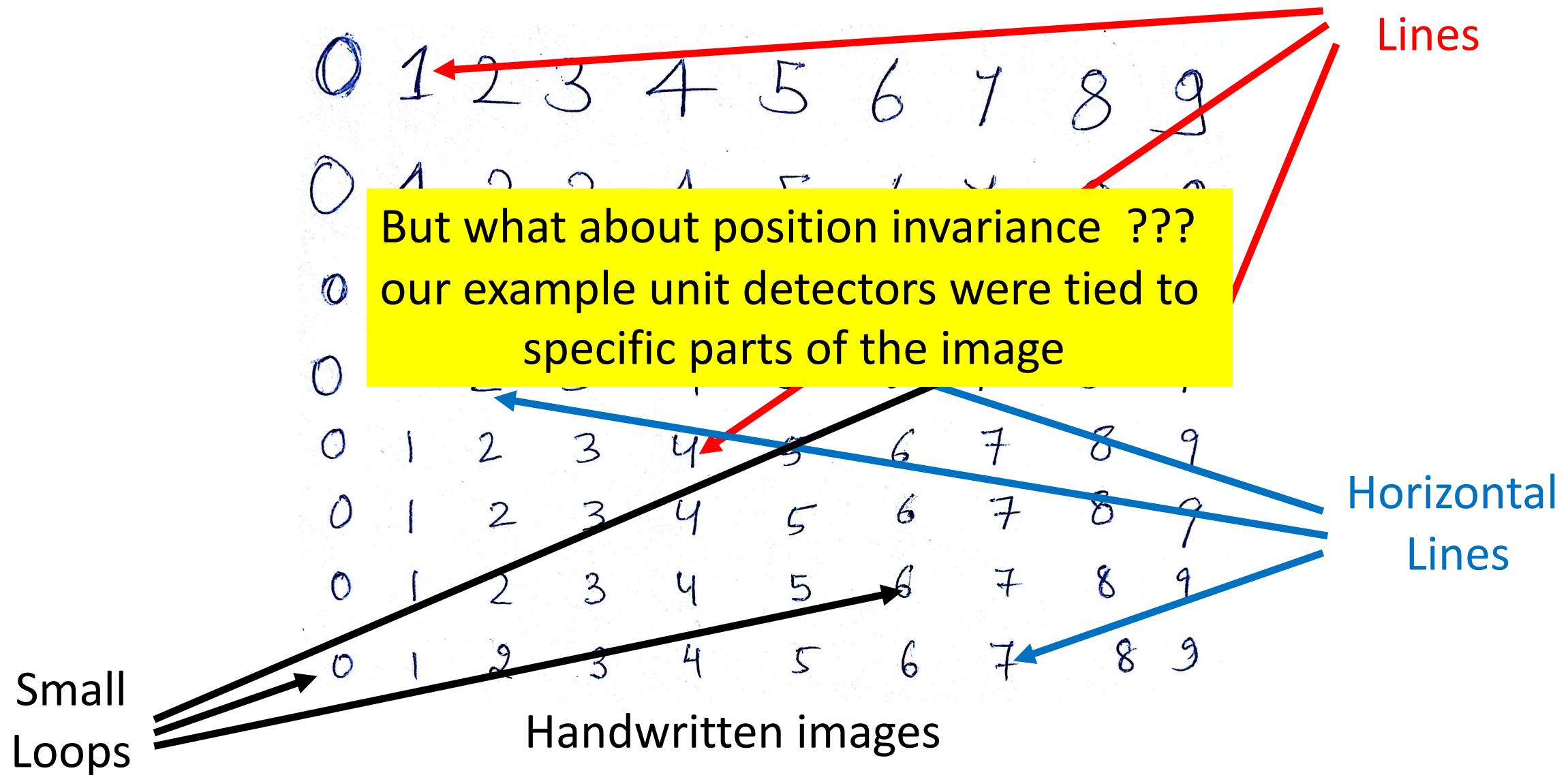
# Horizontal Lines

1



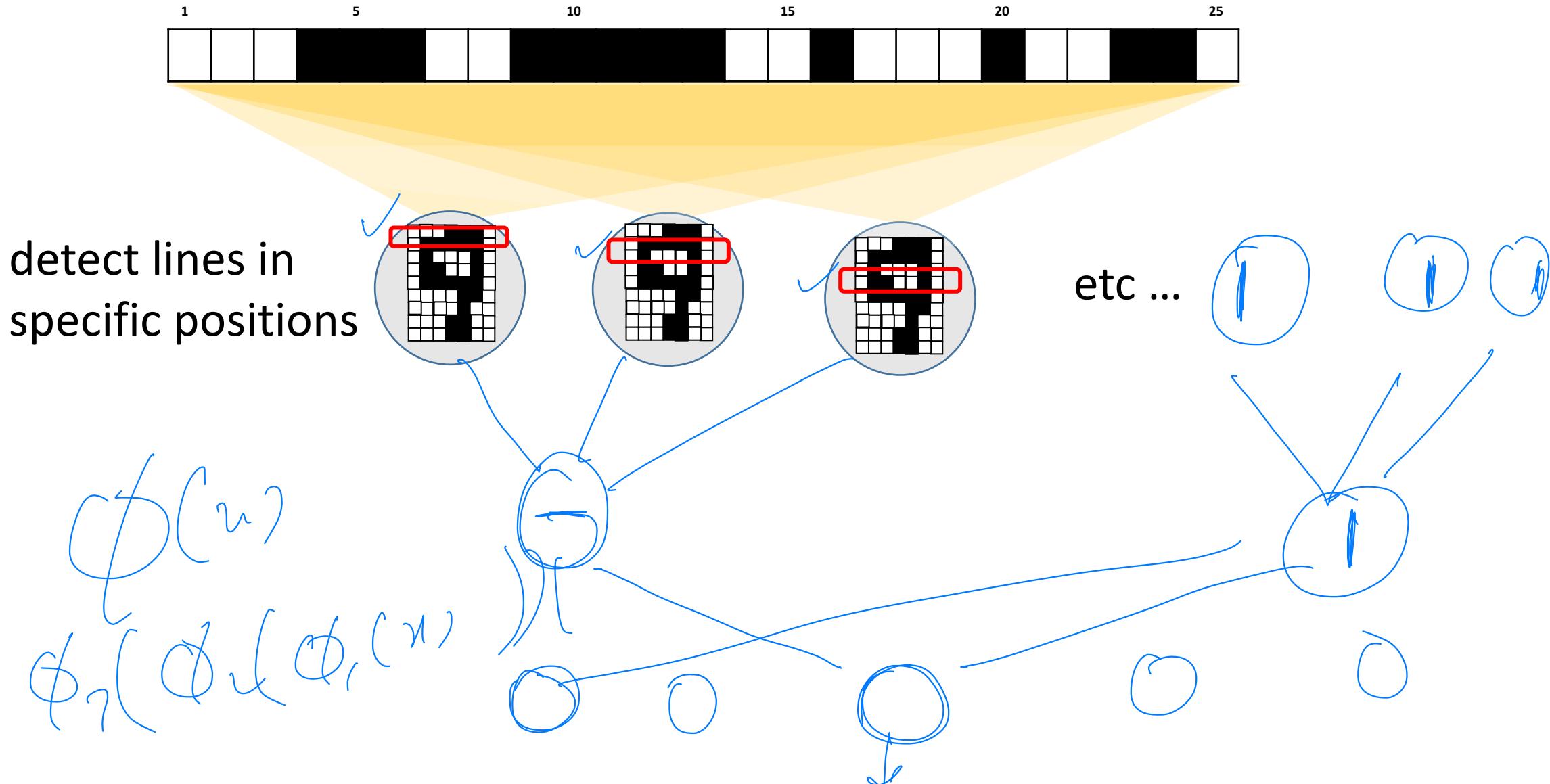


# What does this unit detect?



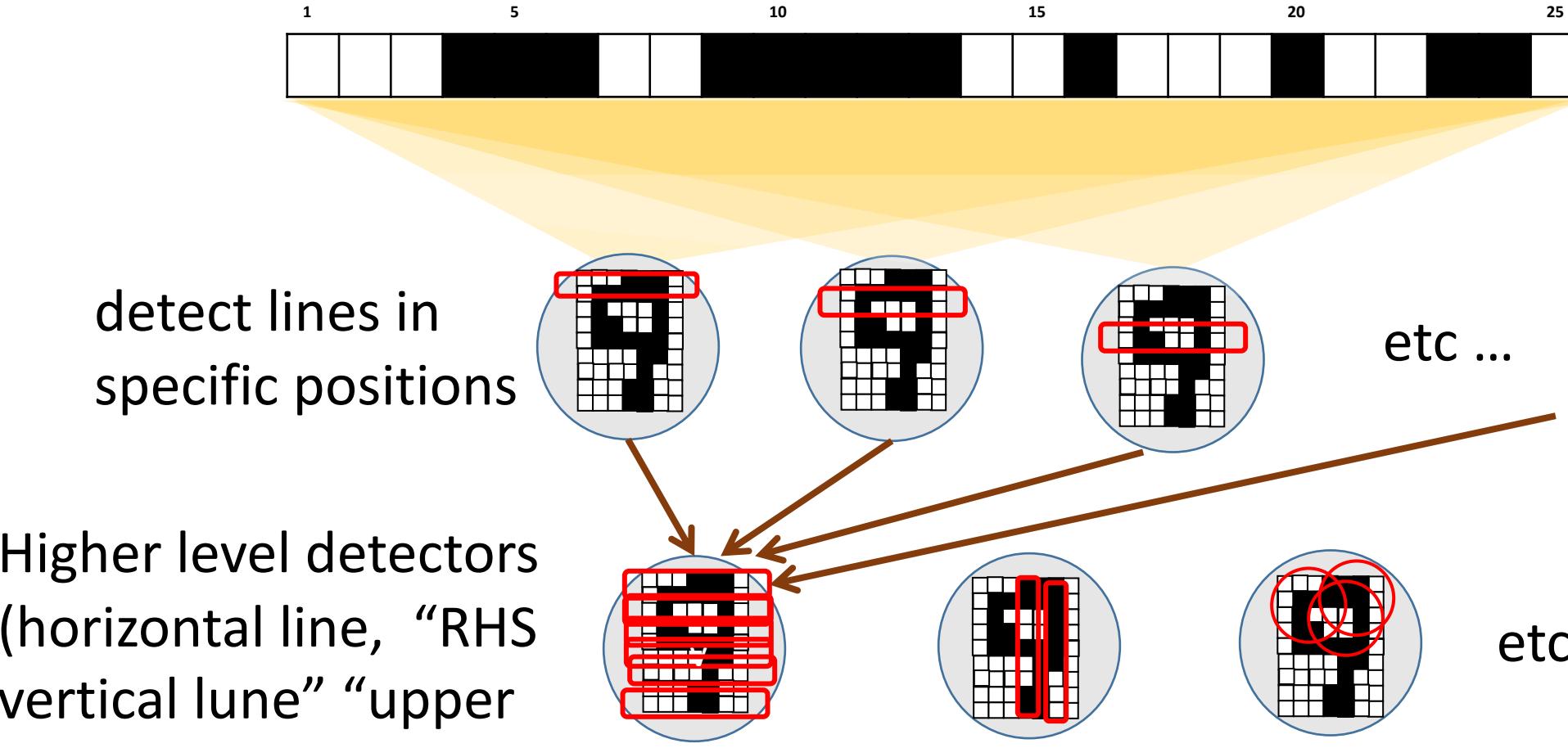


# Successive layers can learn higher-level features



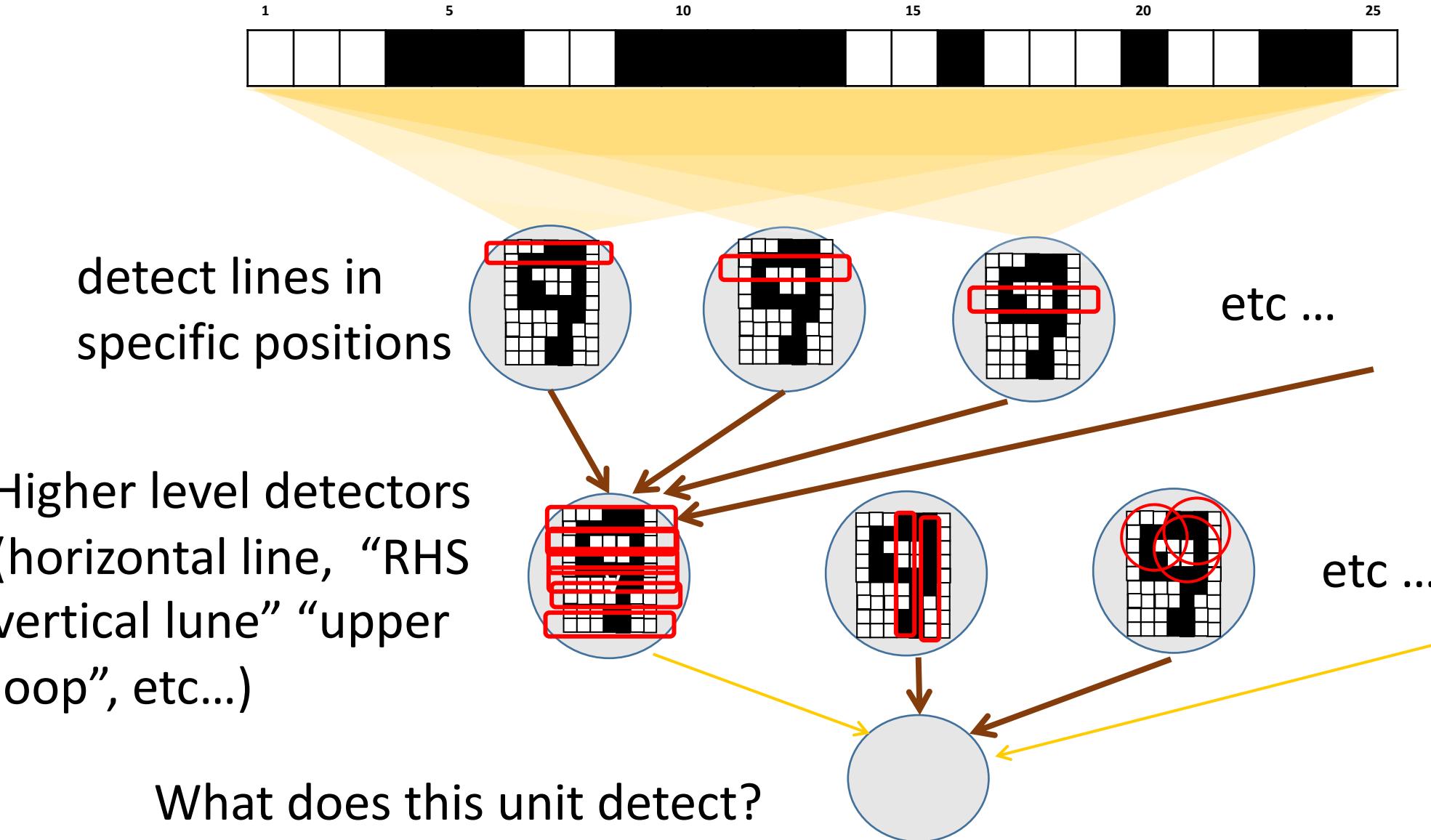


# Successive layers can learn higher-level features





# Successive layers can learn higher-level features

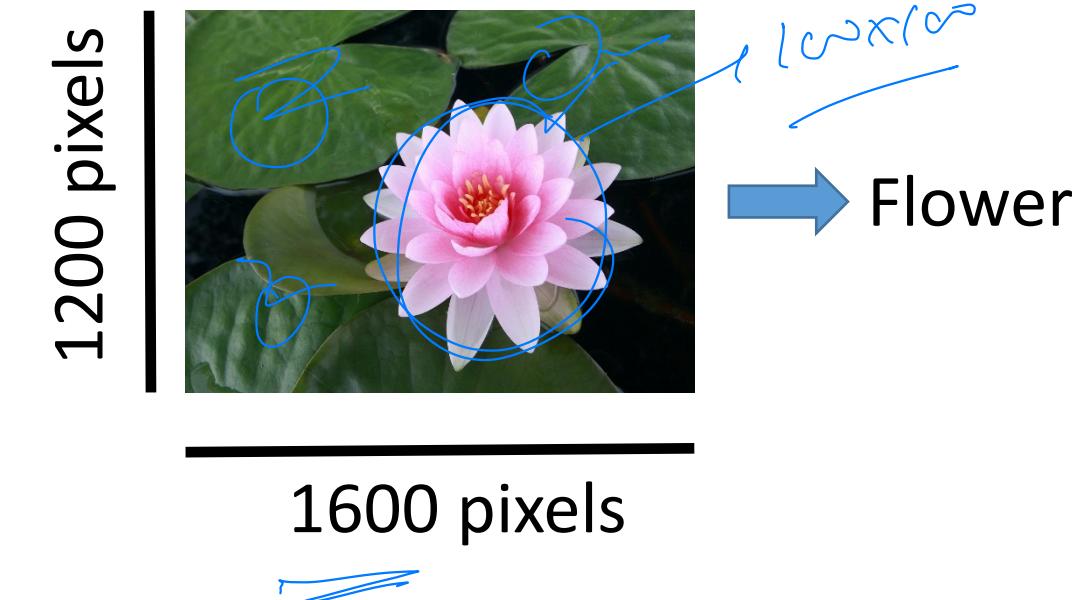


# **Convolutional Neural Networks**



# Motivation: Object Detection in Vision

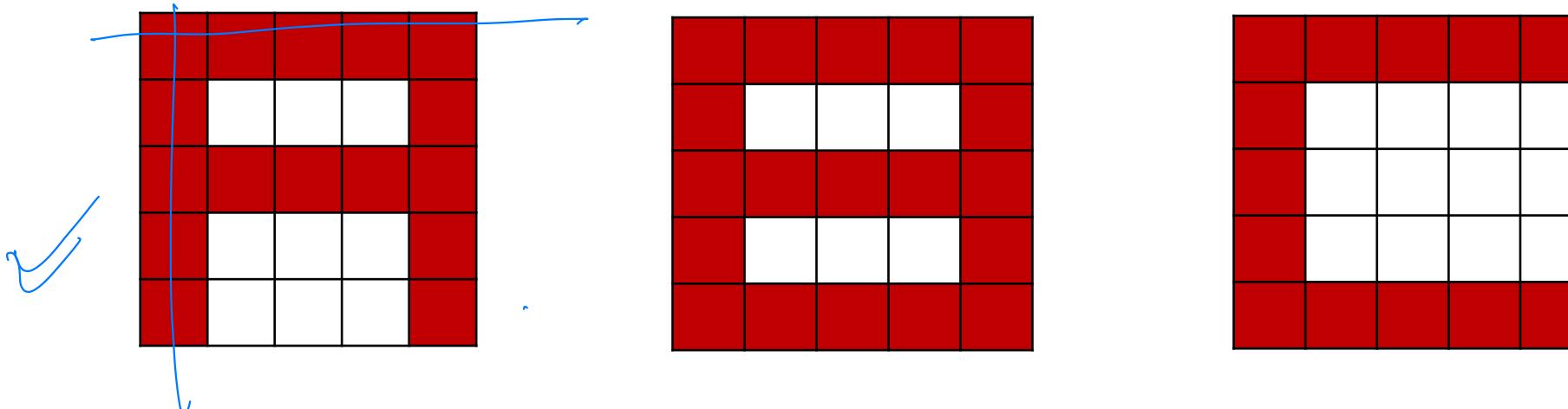
- Challenges
  - Must deal with very high-dimensional inputs
  - $1600 \times 1200 \text{ pixels} = 1.9\text{m inputs, or } 3 \times 1.9\text{m if RGB pixels}$
  - Can we exploit the 2D topology of pixels (or 3D for video data)
  - Can we build invariance to certain variations we can expect
    - Translations, illumination, etc



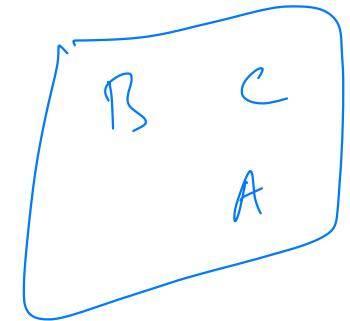
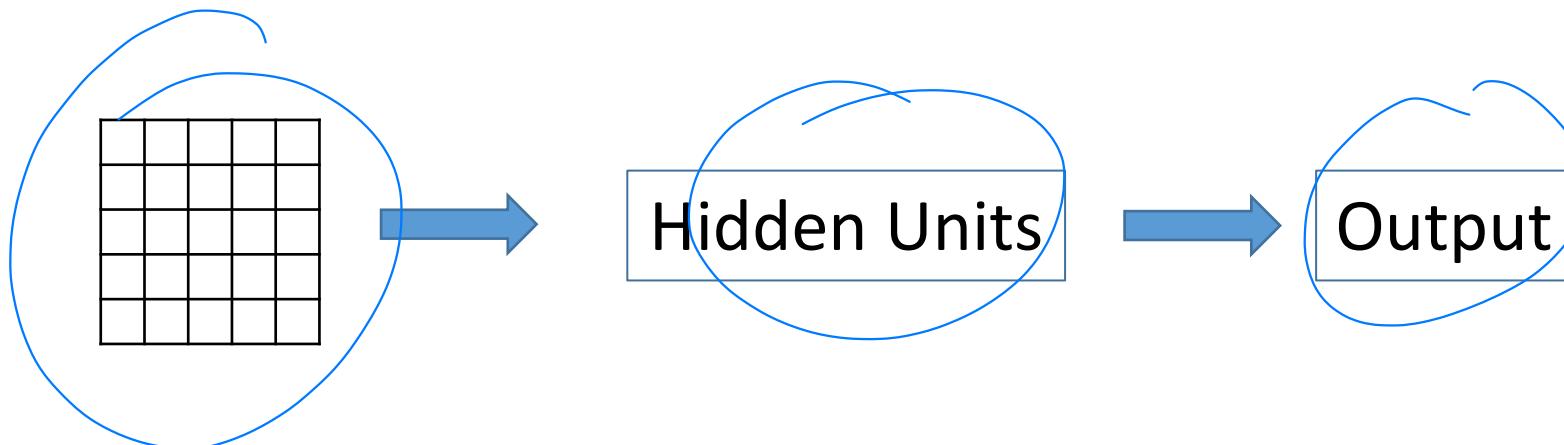


# Motivation: Recognizing an Image

- Input is 5x5 pixel array



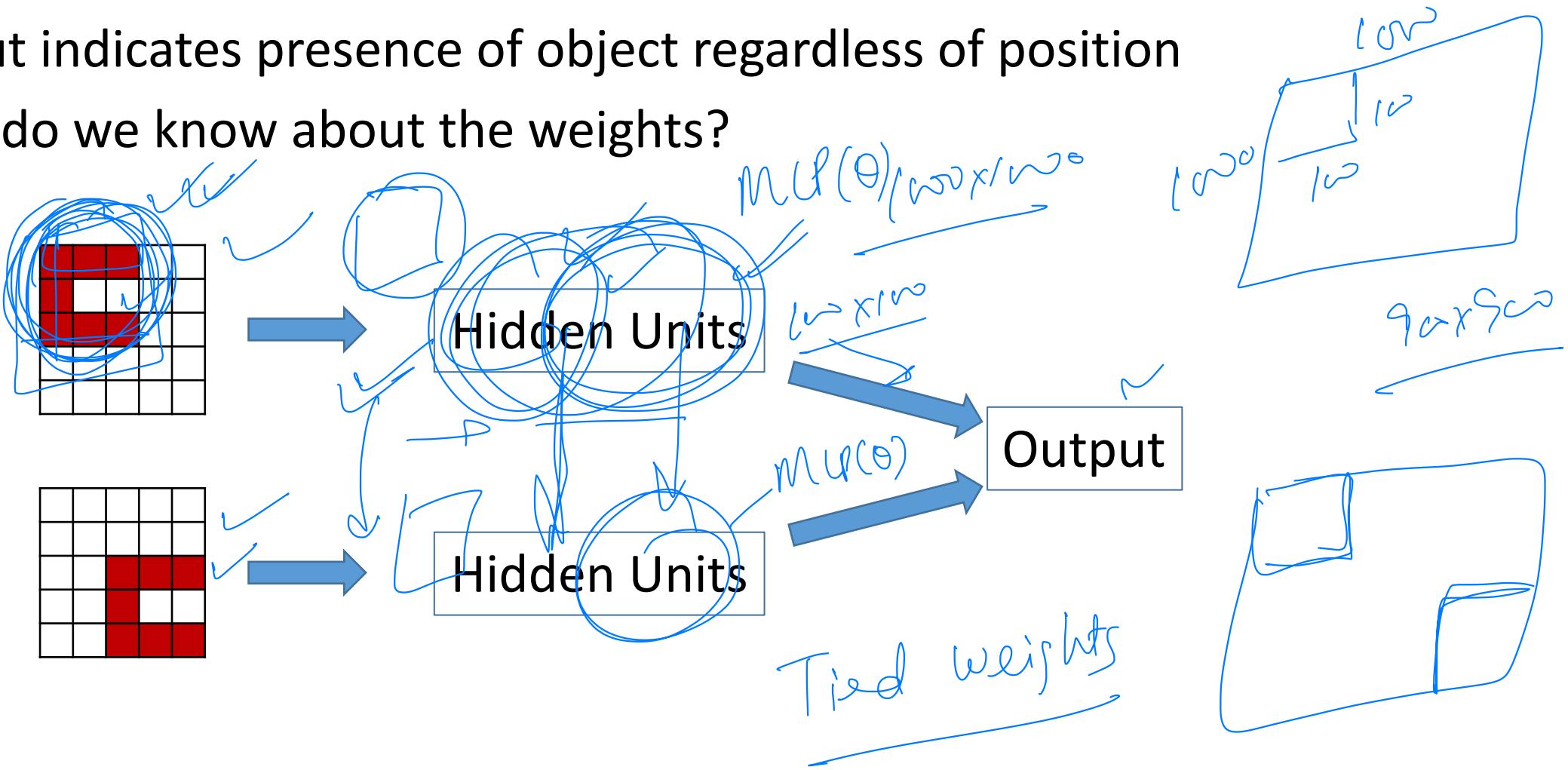
- Simple back propagation net





# Recognizing an Image with unknown location

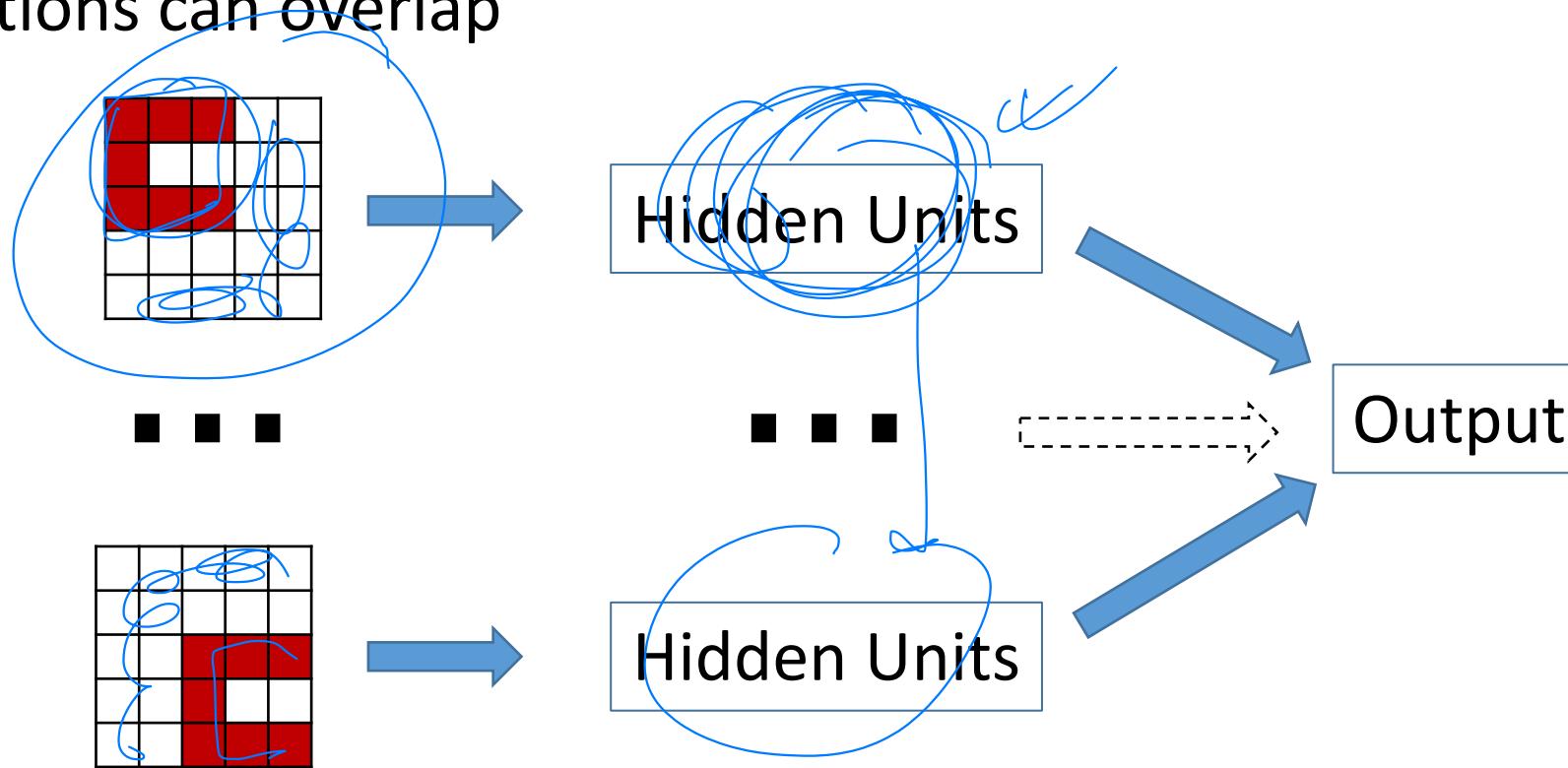
- Object can appear either in the top image or in the bottom image location
- Output indicates presence of object regardless of position
- What do we know about the weights?





# Recognizing an Image with **any** location

- Each possible location the object can appear in has its own set of hidden units
- Each set detects the same features except in a different location
- Locations can overlap





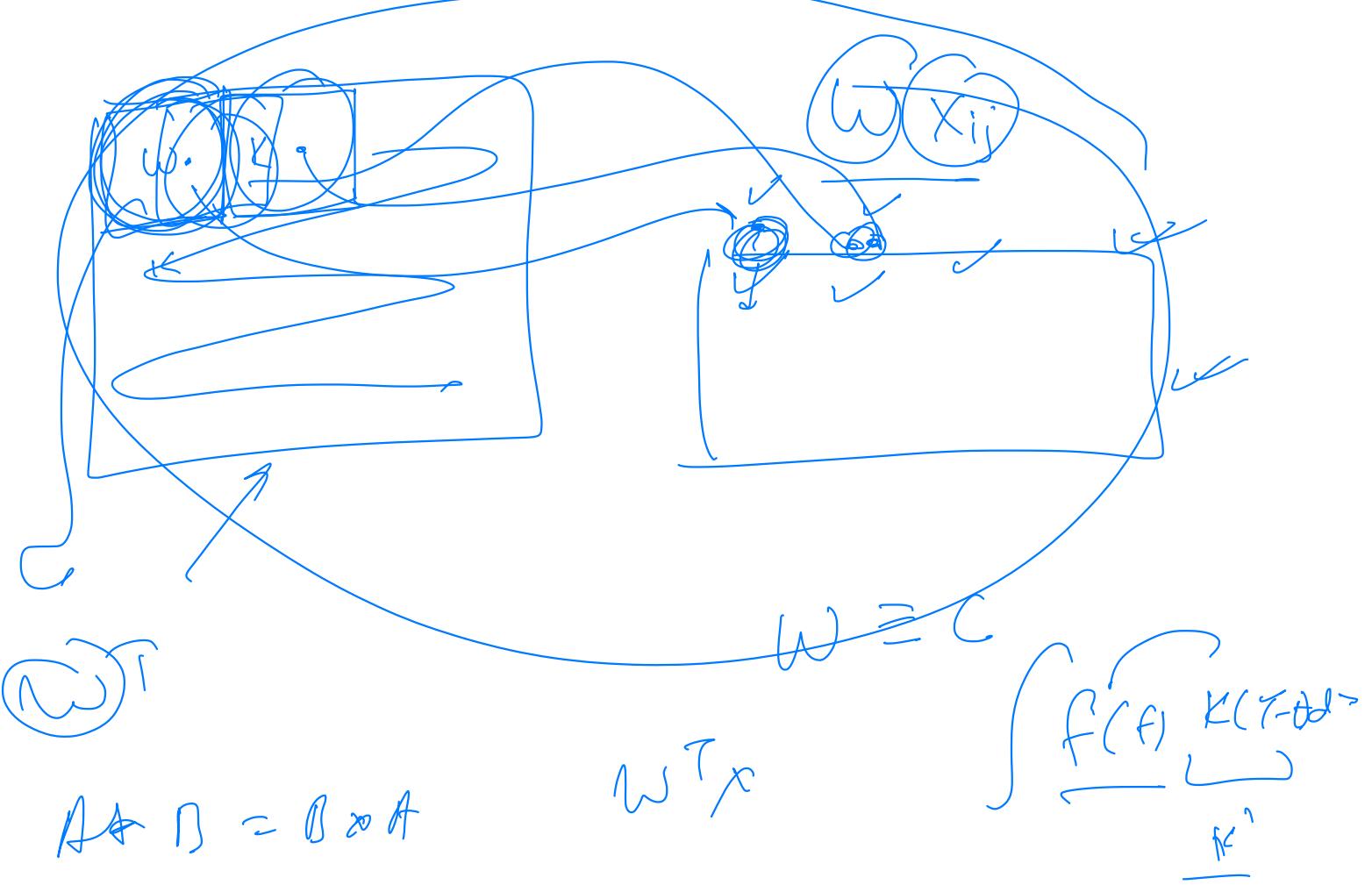
# Convolutional Neural Network (CNN): Key Idea

## Exploit

1. Structure ✓
2. Local Connectivity
3. Share Parameter

## To Give

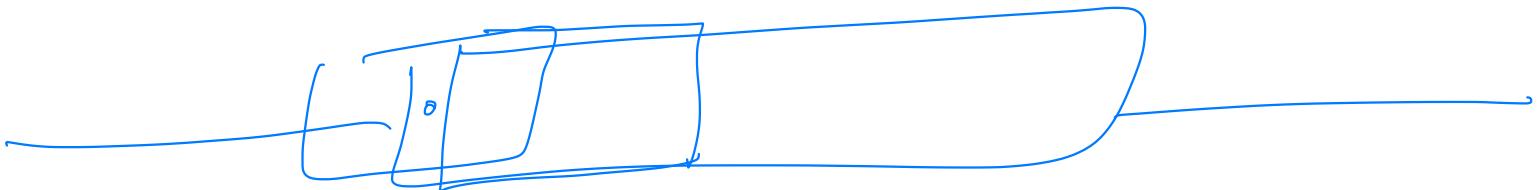
1. Translation Invariance



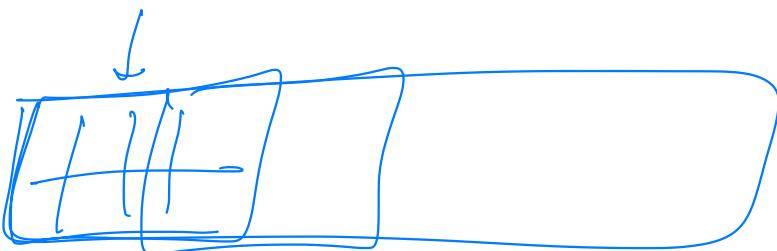
$$s(t) = \int f(\cancel{t}) K(t-\tau) d\tau$$

$(L(t))$

$\int f(\cancel{t}) K(t-\tau) d\tau$   $(L(t))$



$f(t)$

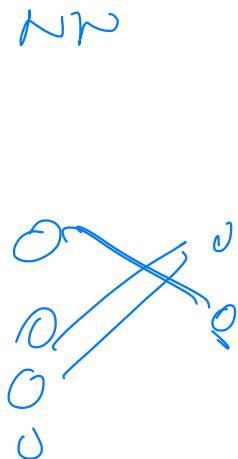




Fully Connected

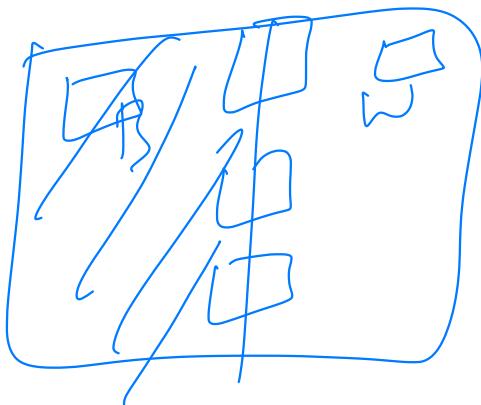
FC

Locally connected  
+ Tied weights



$$k = 100$$
$$k = 7 \times 3$$

$-1$	$0$	$+1$
$-1$	$0$	$+1$
$-1$	$0$	$+1$



$\varphi_1$	$\varphi_1$	$\varphi_1$
$0$	$0$	$0$
$-1$	$-1$	$-1$

$+1$	$+1$	$+1$
$+1$	$+1$	$+1$
$+1$	$+1$	$+1$

$$\begin{bmatrix} 2 & -1 & -1 & -1 \\ -1 & 2 & 2 & 2 \\ 2 & 2 & 2 & 2 \\ +1 & +1 & +1 & +1 \end{bmatrix}$$

Blur + Edge

$K_C \otimes K_I \otimes I$

Architecture

$$\begin{bmatrix} 1 & 1 & 1 \\ 2 & 2 & 2 \\ 1 & 1 & 1 \end{bmatrix}$$



$$\begin{bmatrix} -1 & -1 & 1 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \end{bmatrix}$$

$K \otimes I$

$\downarrow$   
 $(K_2 \otimes K_3)$

1	1	1
2	2	2
1	1	1

=

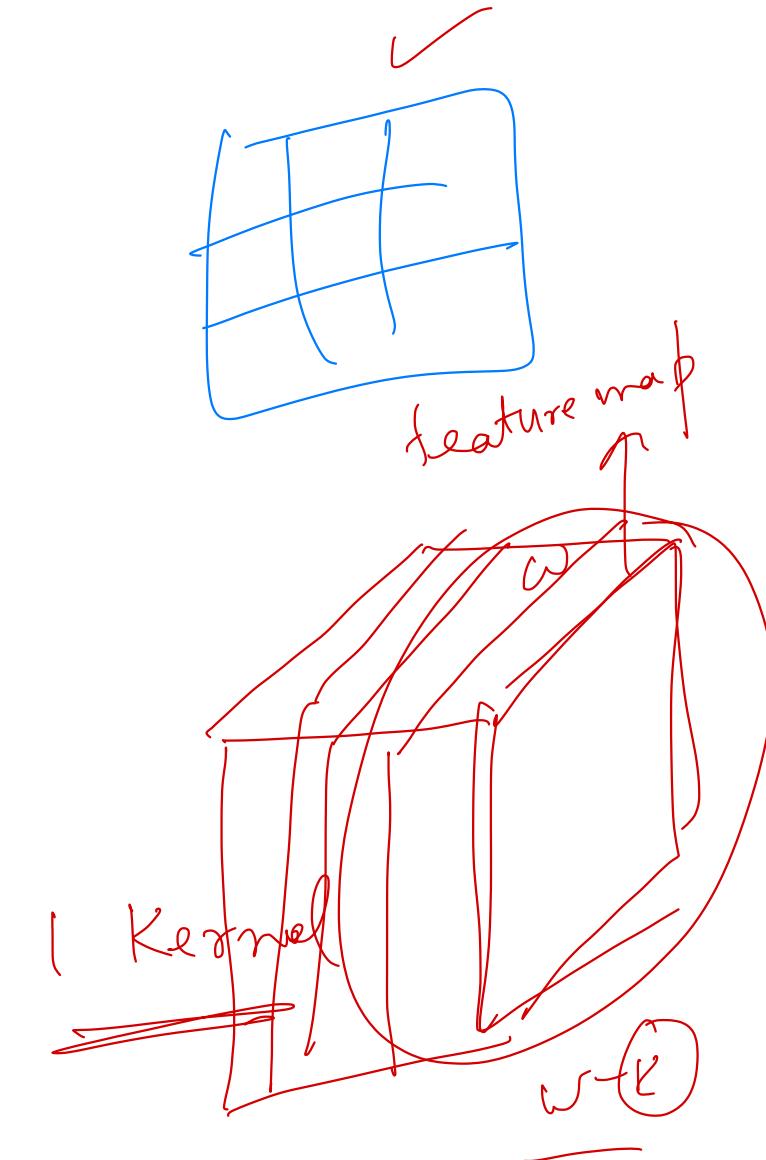
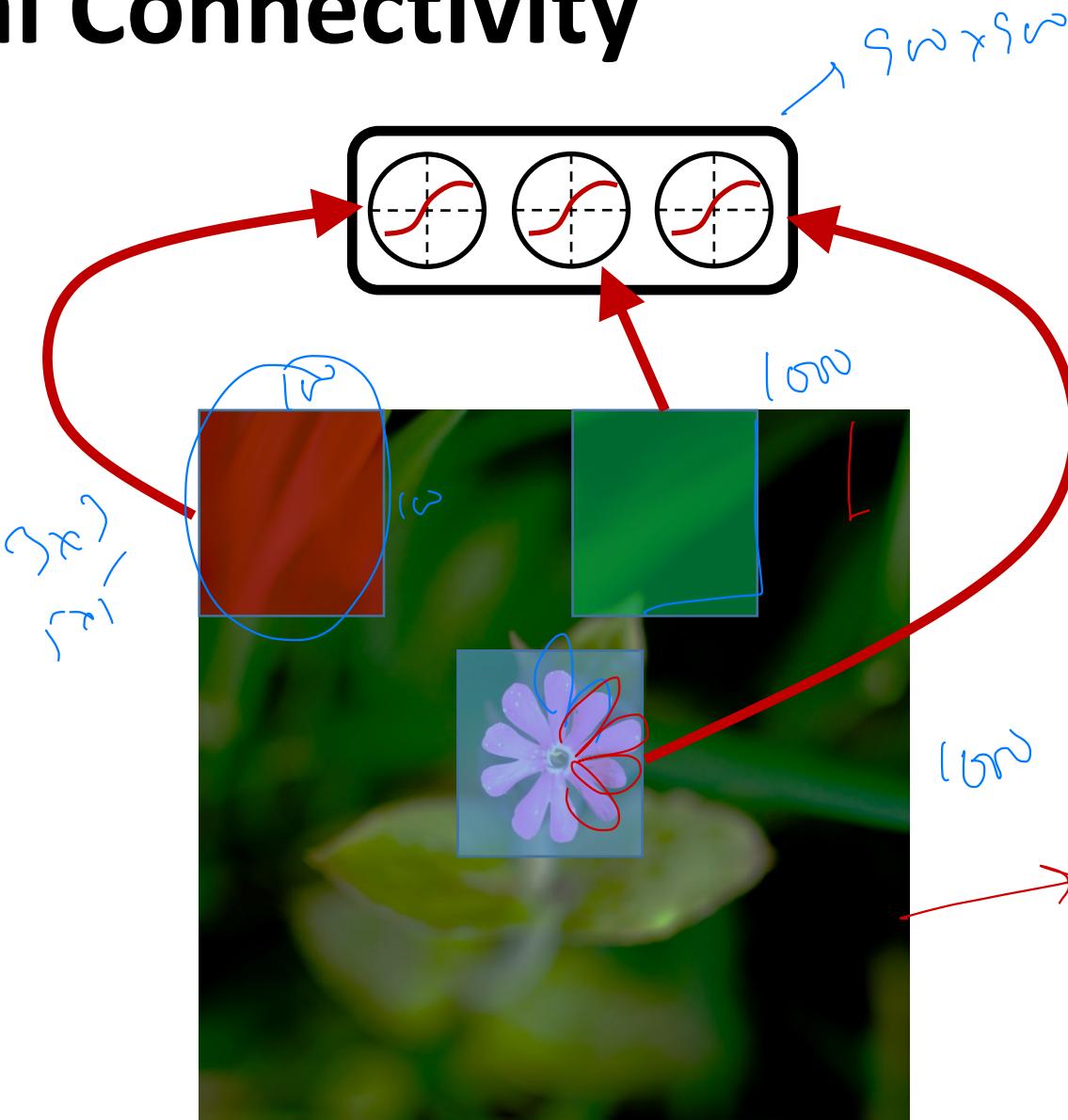
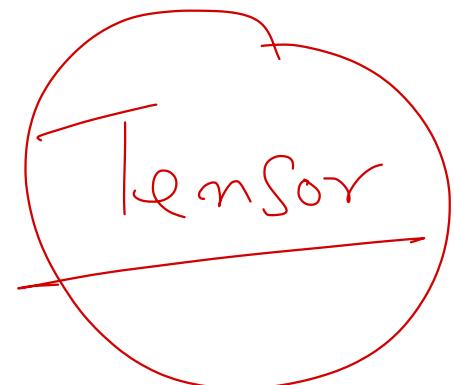
$$\begin{pmatrix} 1 \\ 2 \\ 1 \end{pmatrix}$$

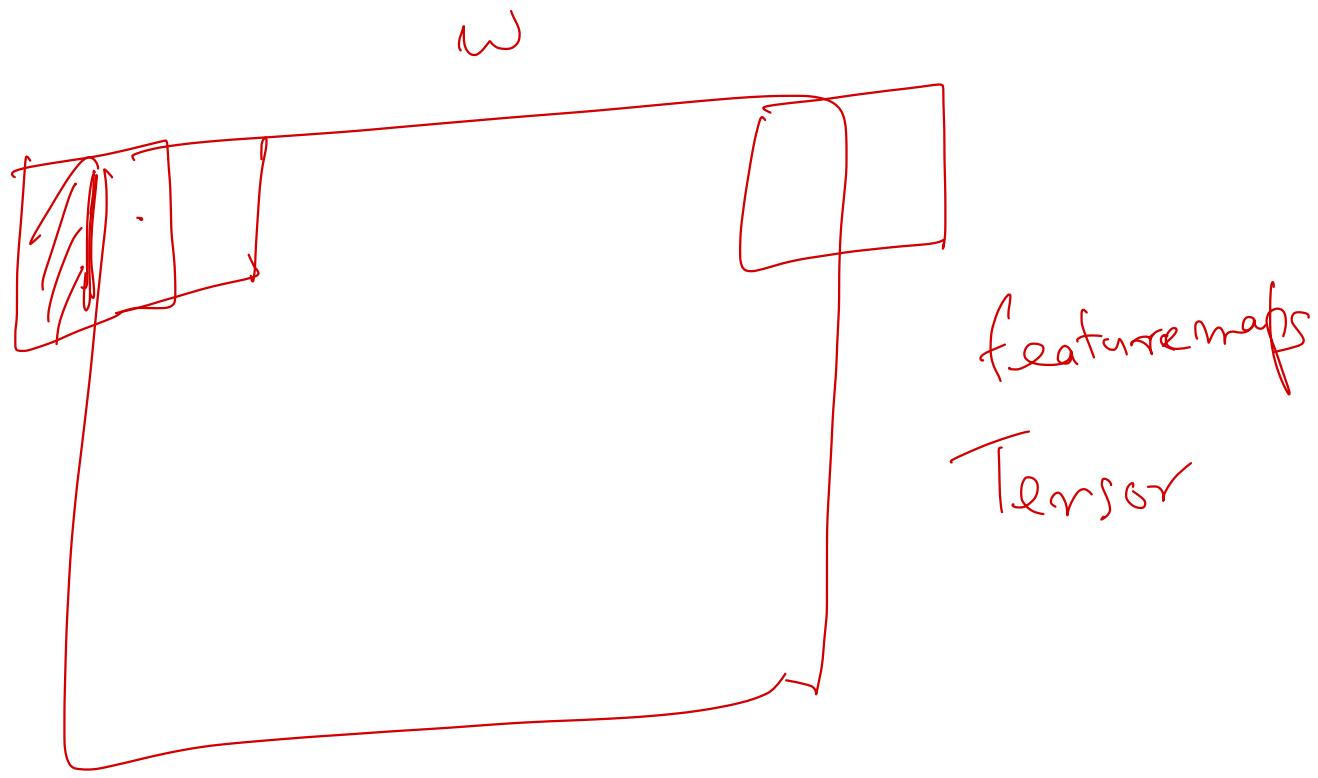
$\Rightarrow$

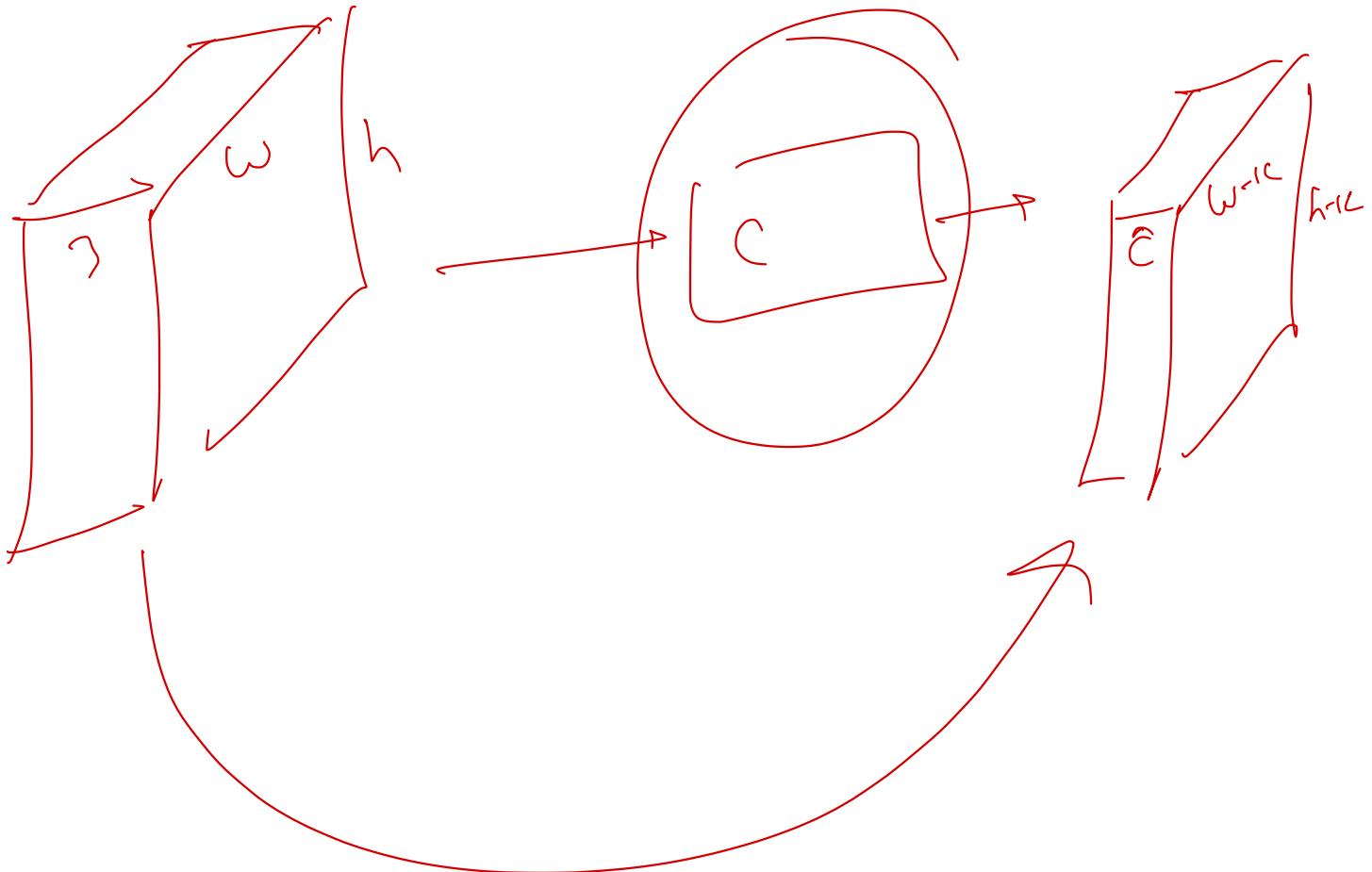
1	1	1
---	---	---

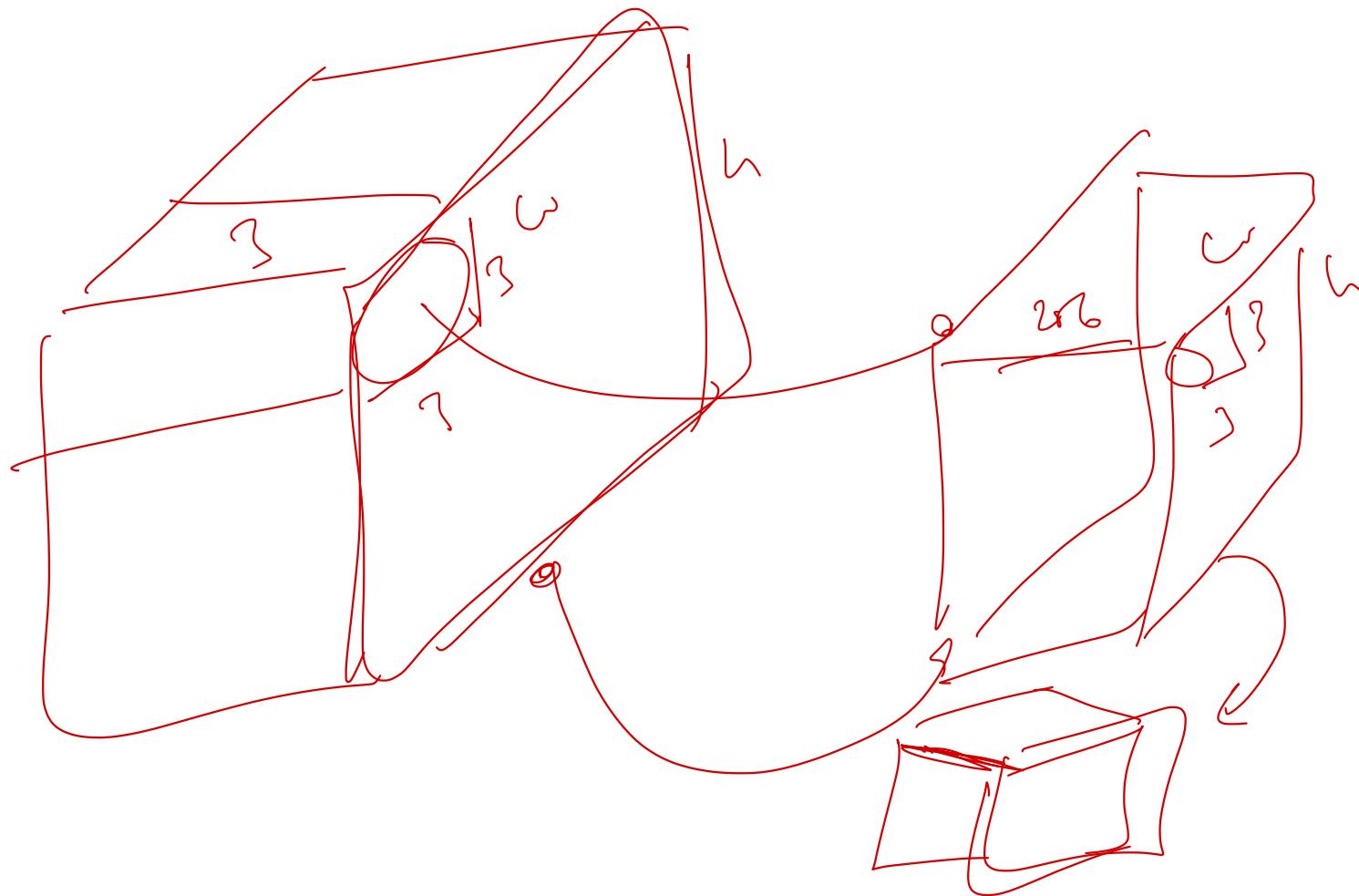


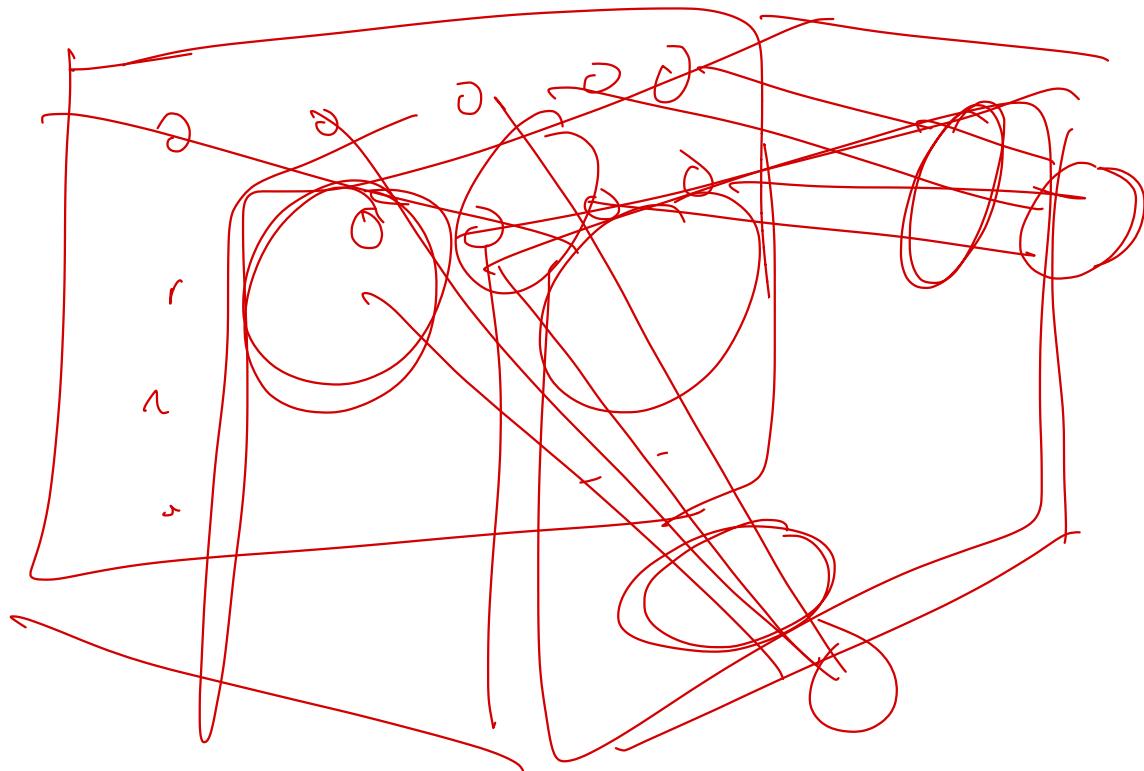
# CNN: Local Connectivity

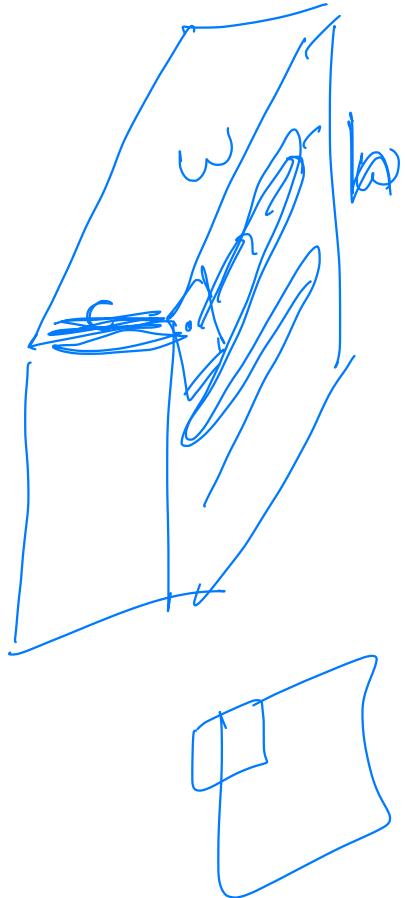




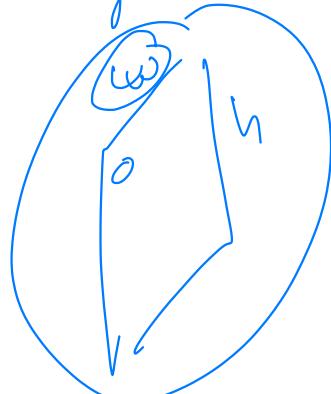






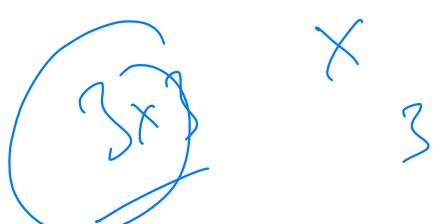


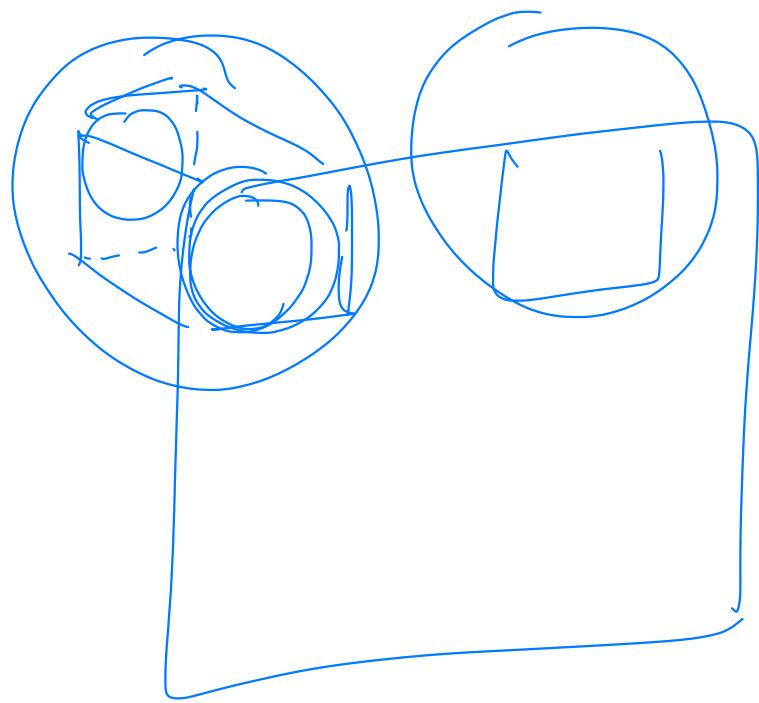
Convolutional  
Layer



$$\sum w_i(x_i) \rightarrow [ ]$$

27







# CNN: Local Connectivity

- Each hidden unit is connected only to a sub-region (patch) of the input image
- It is connected to all channels
  - 1 if greyscale image
  - 3 (R, G, B) for color image

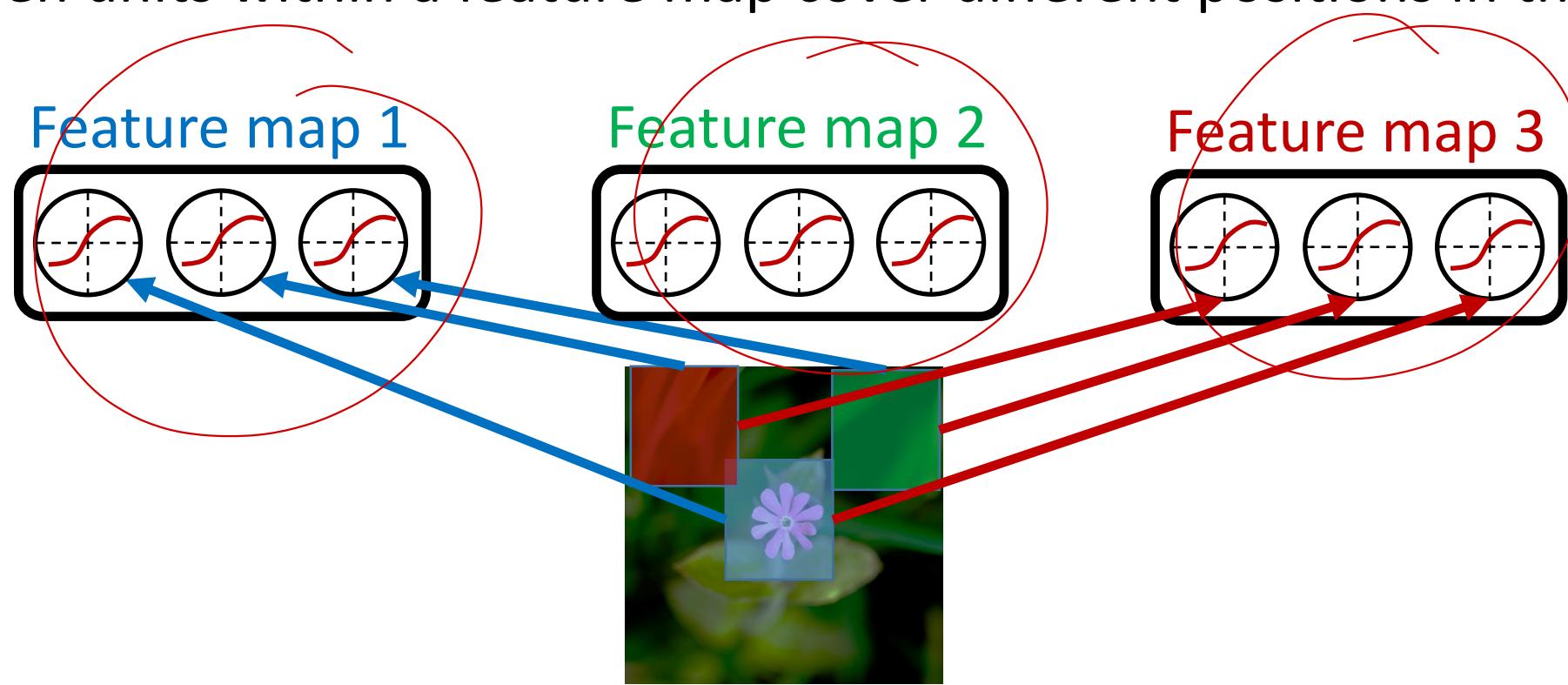
Solves the following problems:

- Fully connected hidden layer would have an unmanageable number of parameters
- Computing the linear activations of the hidden units would be very expensive



# CNN: Parameter Sharing

- Units organized into the same “feature map/template” share parameters
- In this way all neurons detect the same feature/template at different positions in the input image
- Hidden units within a feature map cover different positions in the image





# CNN: Parameter Sharing

## How does it help?

- Reduces even more number of parameters (compared to local connectivity)
- Will extract the same features at every position
  - features are “equi-variant”

$$\phi_3(\phi_2(\phi_1(x)))$$

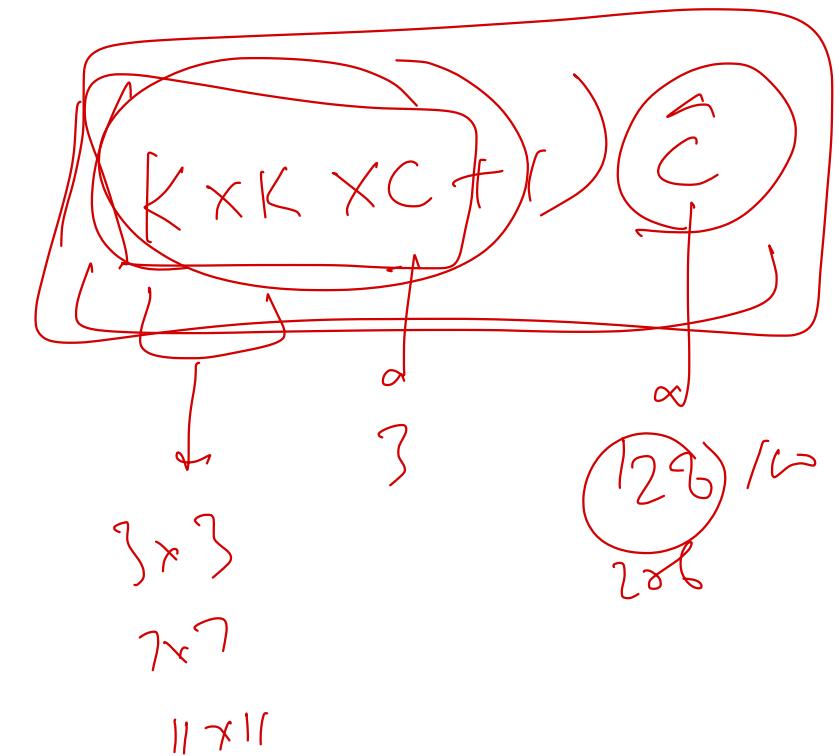
$3w \times 3x \Rightarrow 27 \approx$

independent of image size

$\approx$

$$y = \phi(x)$$

$\downarrow$



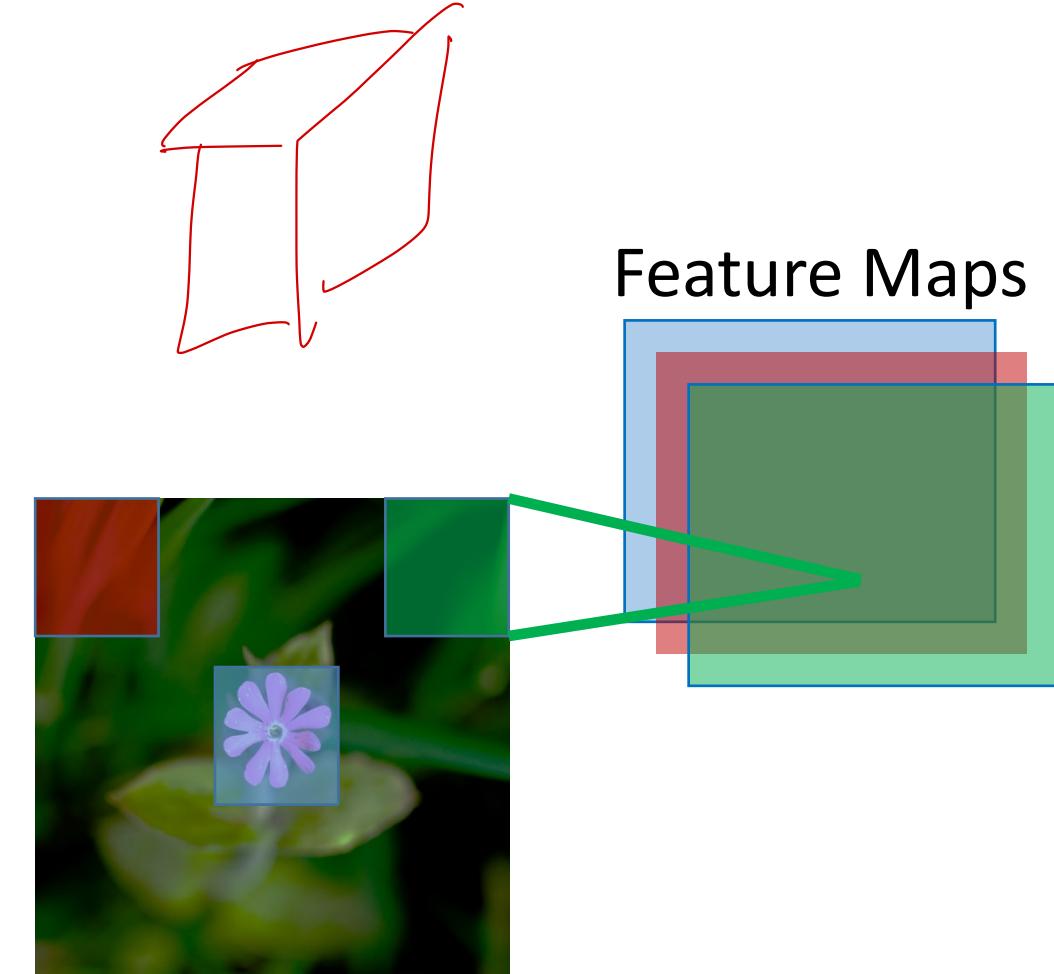


# CNN: Parameter Sharing

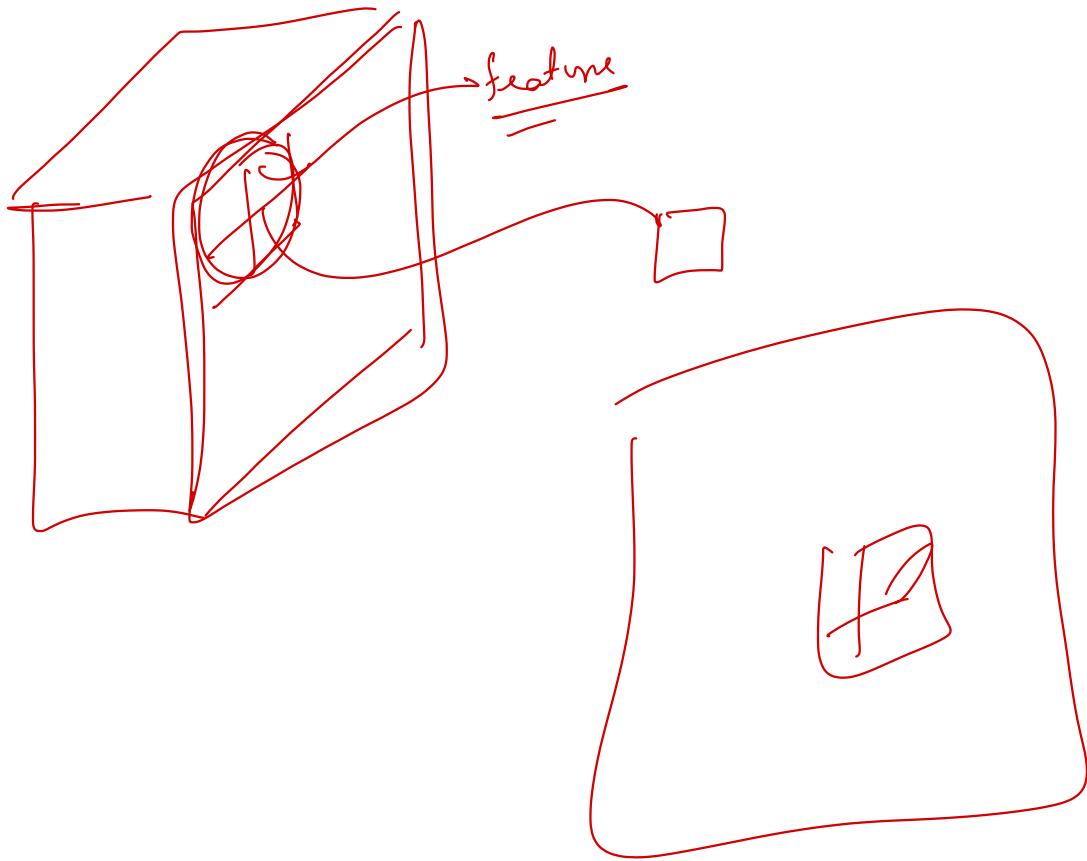
- Each feature map forms a 2D grid of features
- Can be computed with a discrete convolution (\*) of a kernel matrix  $k_{ij}$  which is the hidden weights matrix  $W_{ij}$  with its rows and columns flipped:

$$y_j = \tanh \sum_i k_{ij} x_i$$

- $x_i$  is the  $i^{th}$  channel of input
- $k_{ij}$  is the convolution kernel
- $y_j$  is the hidden layer



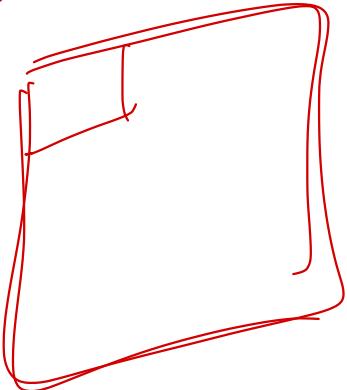
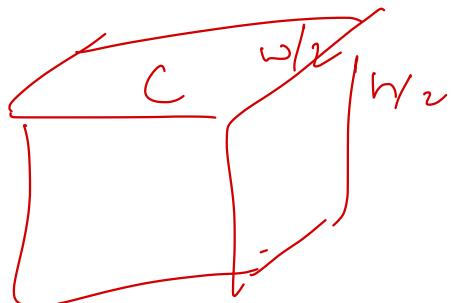
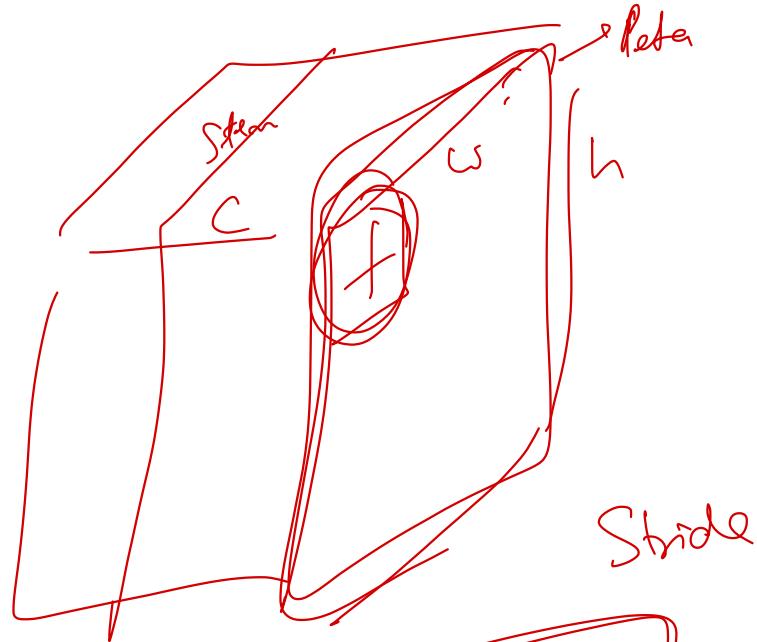
Pooling

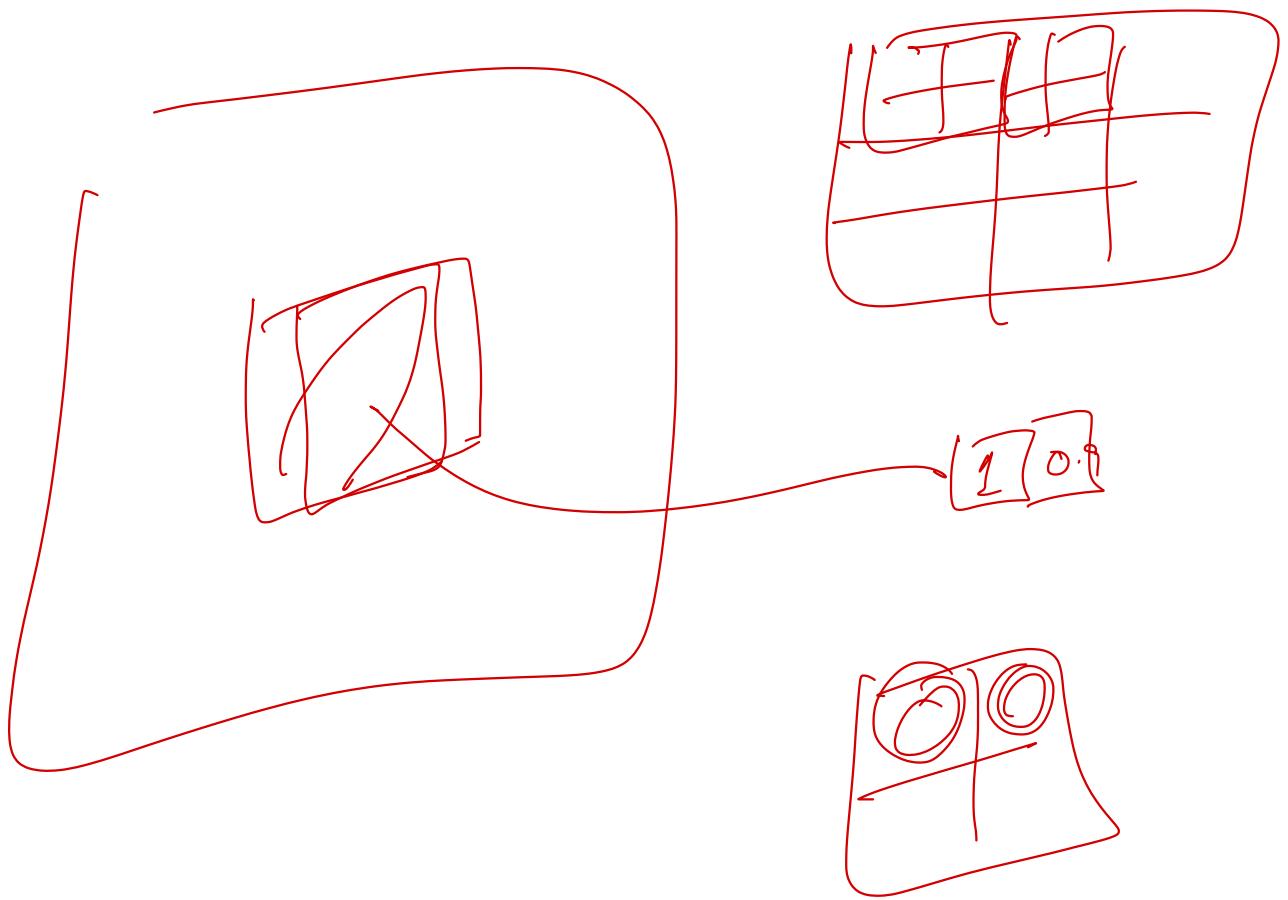


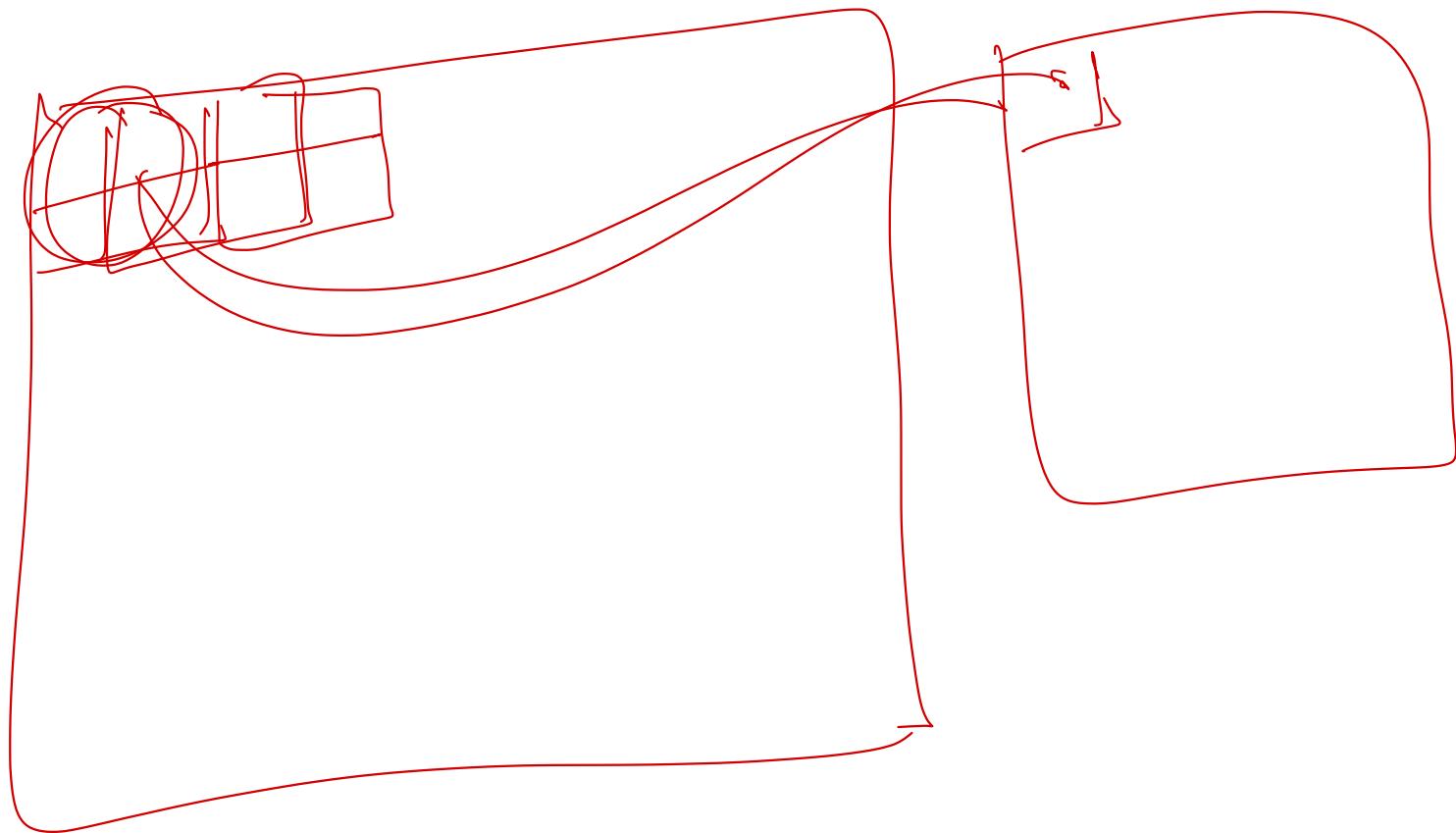
Pooling

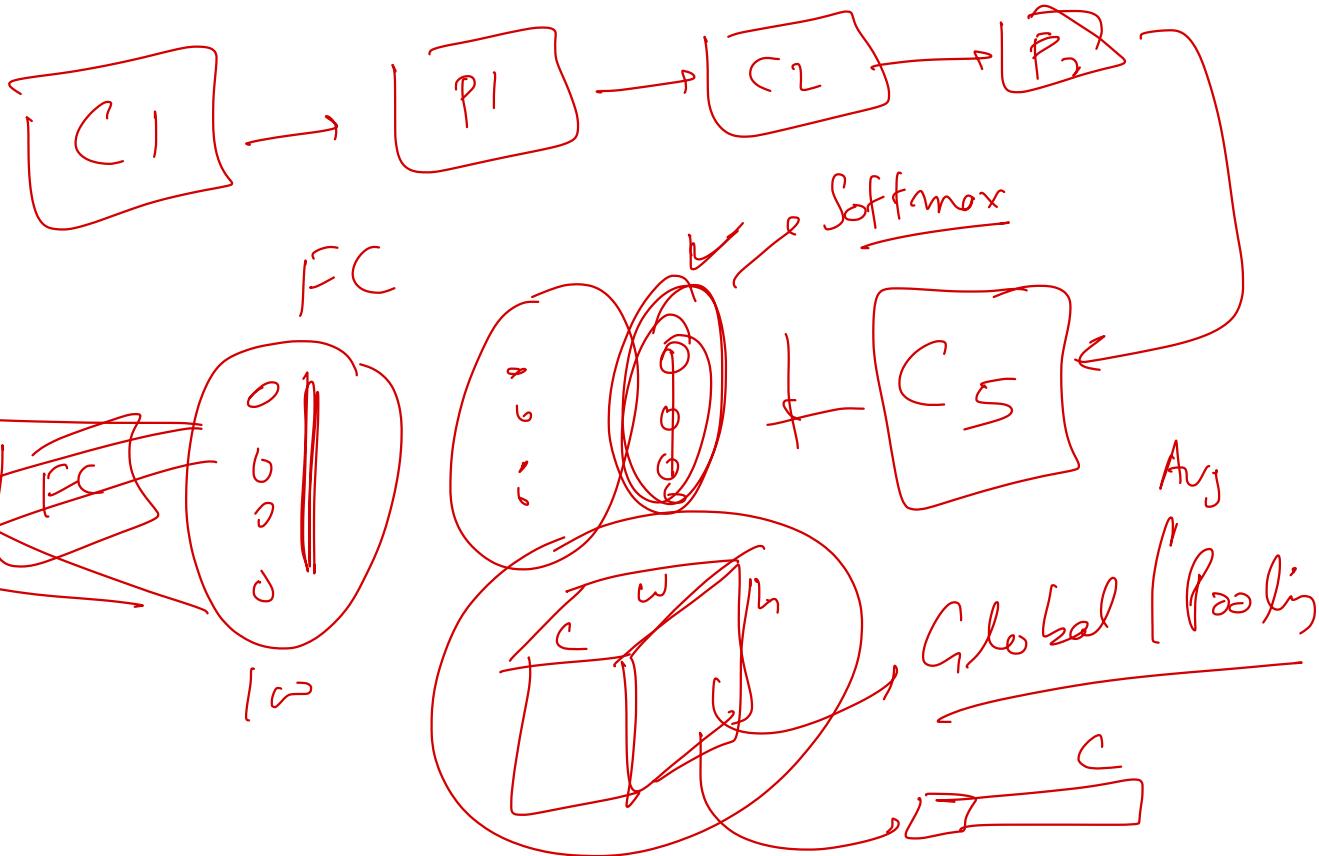
max  
average

Pooling  
Convolution + Stride











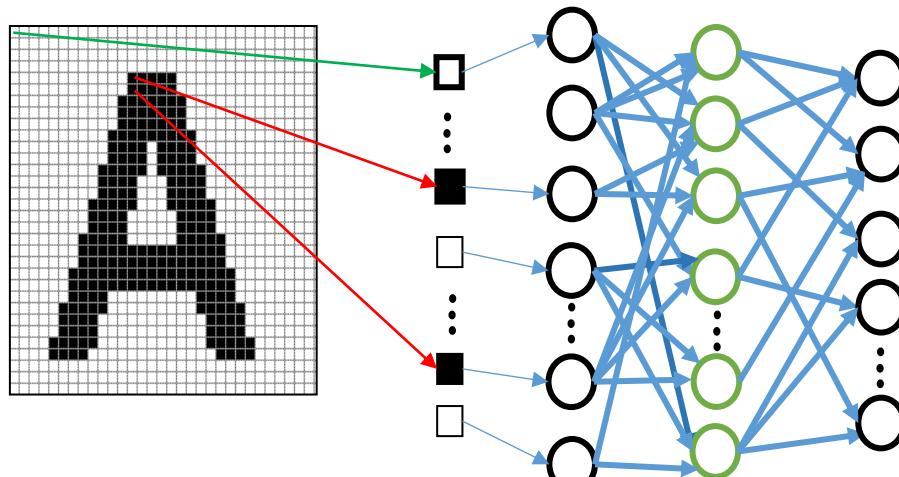
# CNN: Parameter Sharing

- **Stride** (typically denoted by **s** in CNNs) decides the spacing between overlapping convolutions
- Helps in reducing parameters further

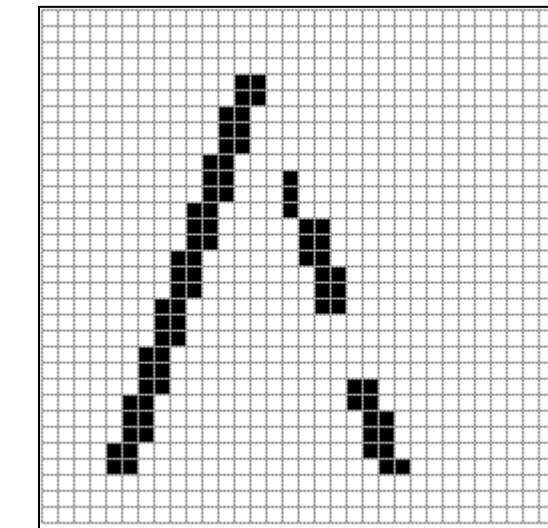
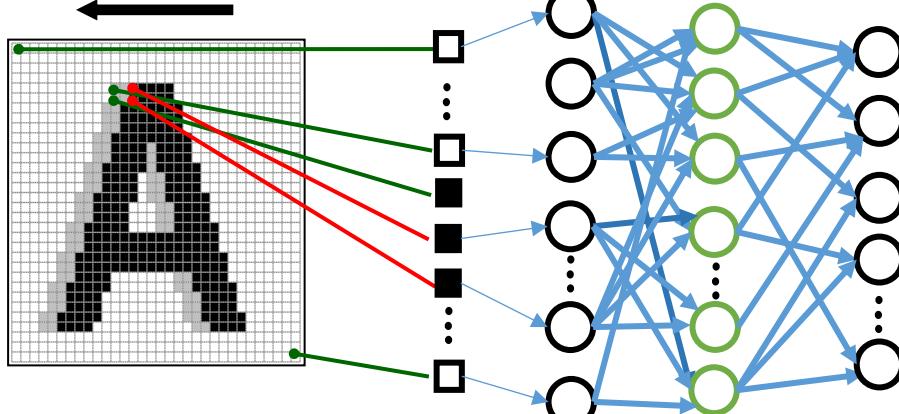


# Drawbacks of Neural Networks

- Little or no invariance to shifting



Shift left



154 input change  
from 2 shift left  
77 : black to white  
77 : white to black



# Drawbacks of Neural Networks

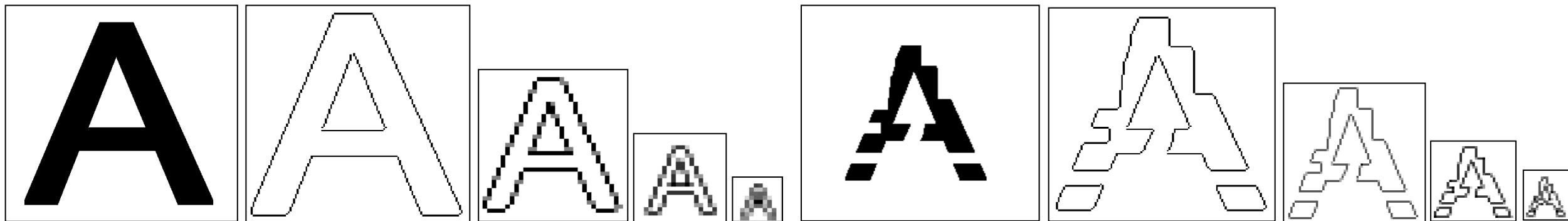
- Scaling, or other forms of distortion





# Pooling or Subsampling Layer

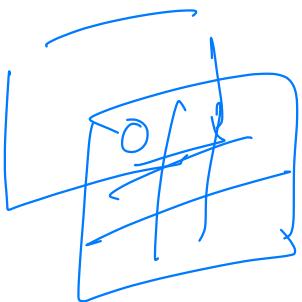
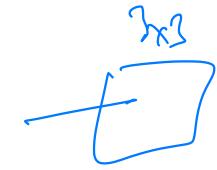
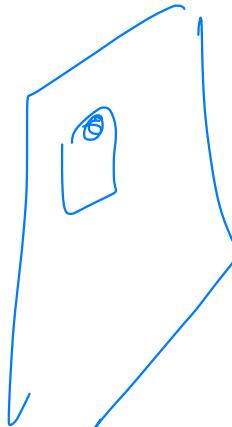
- The pooling layers reduce the spatial resolution of each feature map. Helps to collate the features in a region
- By reducing the spatial resolution of the feature map, a certain degree of noise, shift and distortion invariance is also achieved.



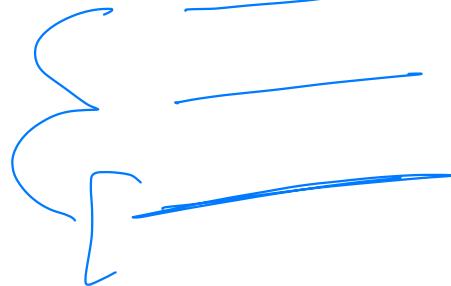


# Jargon

- Convolutional Neural Networks
  - also called CNNs, Conv Nets etc.
- Each hidden unit channel
  - also called **map, feature, feature type, dimension**
- Weights for each channel
  - also called **kernels**
- Input patch to a hidden unit at  $(x,y)$ 
  - also called **receptive field**



→ higher layer



lower layers



# Typical CNN

- Alternates convolutional and pooling layers
- Output layer is a regular, fully connected layer with softmax non-linearity
  - Output provides an estimate of the conditional probability of each class
- The network is trained by stochastic gradient descent

# **Backpropagation**



# Backpropagation in General Cases

1. Decompose operations in layers of a neural network into function elements whose derivatives w.r.t. inputs are known by symbolic computation.

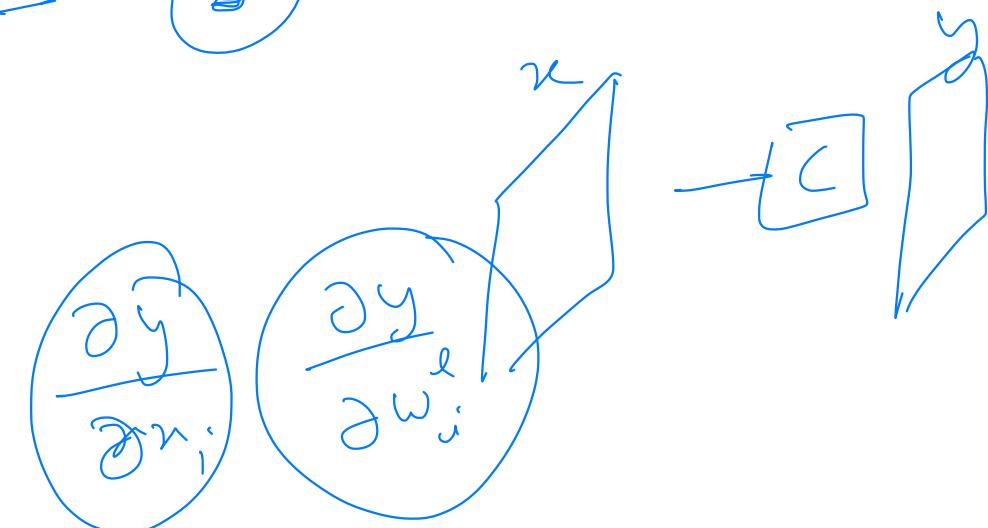
$$h_{\theta}(x) = \left( f^{(l_{max})} \circ \dots \circ f_{\theta^{(l)}}^{(l)} \circ \dots \circ f_{\theta^{(2)}}^{(2)} \circ f^{(1)} \right)(x)$$

where  $f^{(1)} = x$

$f^{(l_{max})} = h_{\theta}$  cost function

Sigmoid  
Softmax

and  $\forall l: \frac{\partial f^{(l)}}{\partial f^{(l-1)}}$  is known



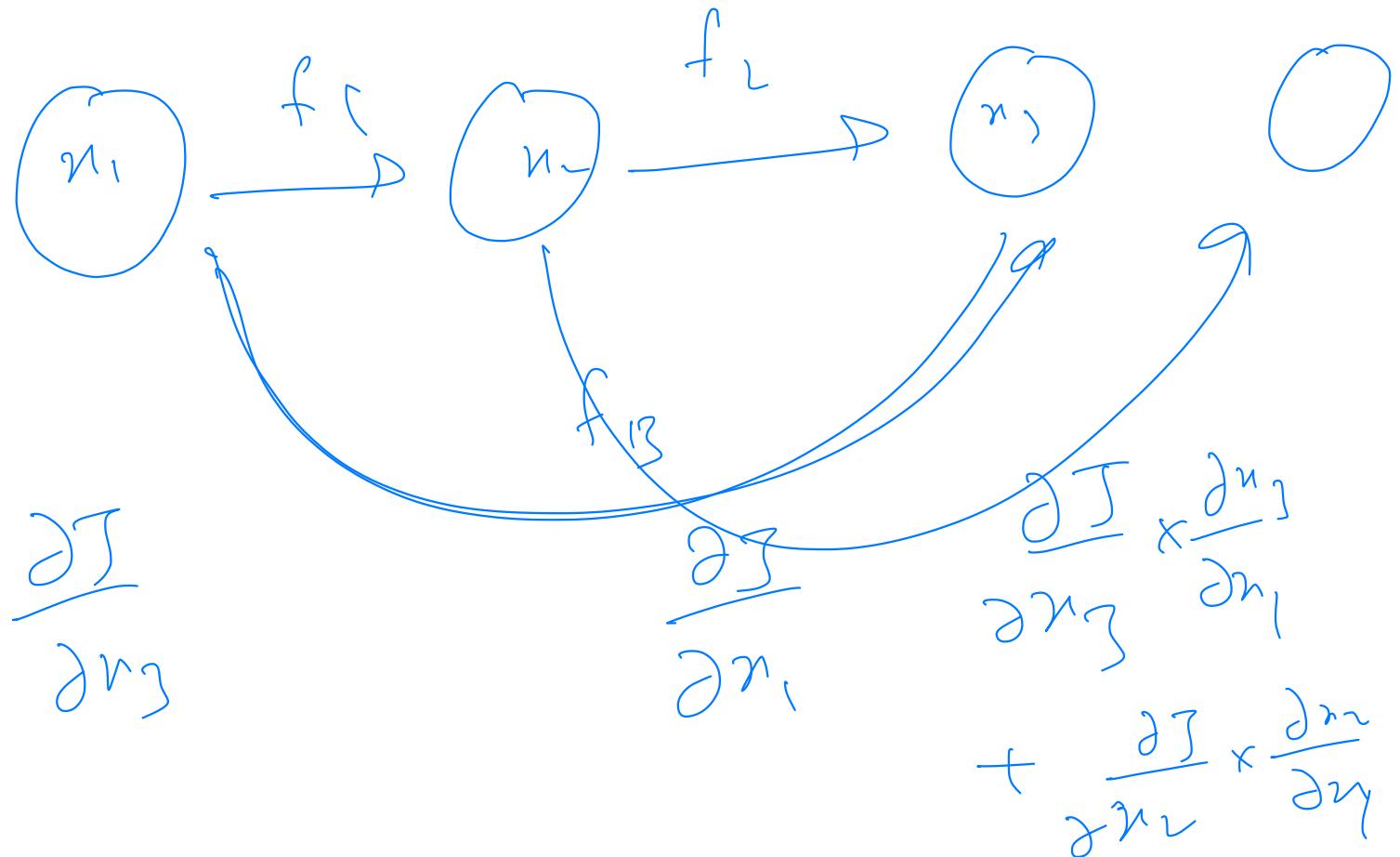


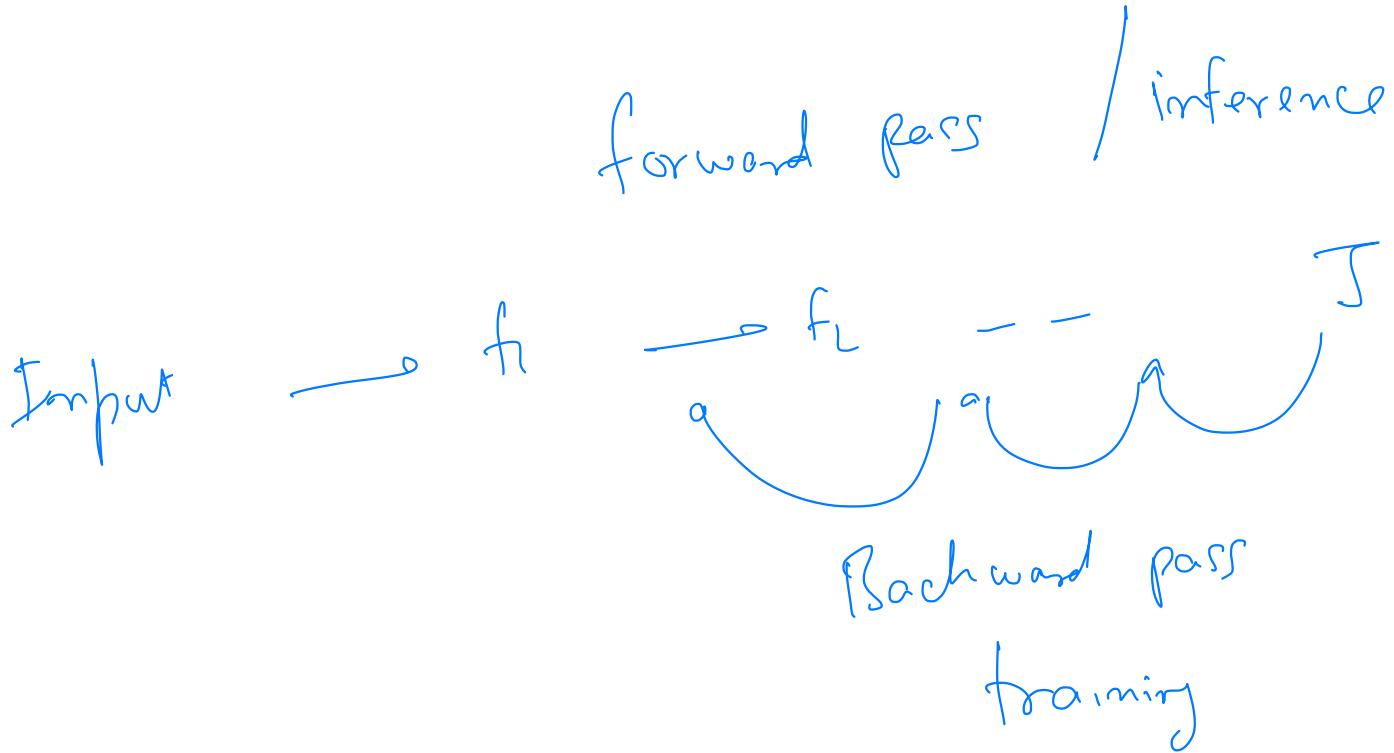
# Backpropagation in General Cases

2. Backpropagate error signals corresponding to a differentiable cost function by numerical computation (Starting from cost function, plug in error signals backward).

$$\delta^{(l)} = \frac{\partial J}{\partial f^{(l)}} J(\theta; x, y) = \frac{\partial J}{\partial f^{(l+1)}} \frac{\partial f^{(l+1)}}{\partial f^{(l)}} = \delta^{(l+1)} \frac{\partial f^{(l+1)}}{\partial f^{(l)}}$$

where  $\frac{\partial J}{\partial f^{(l_{max})}}$  is known





$$\sum_k \left( \frac{\partial J}{\partial w_i^l} \right)^k$$

mini-batches

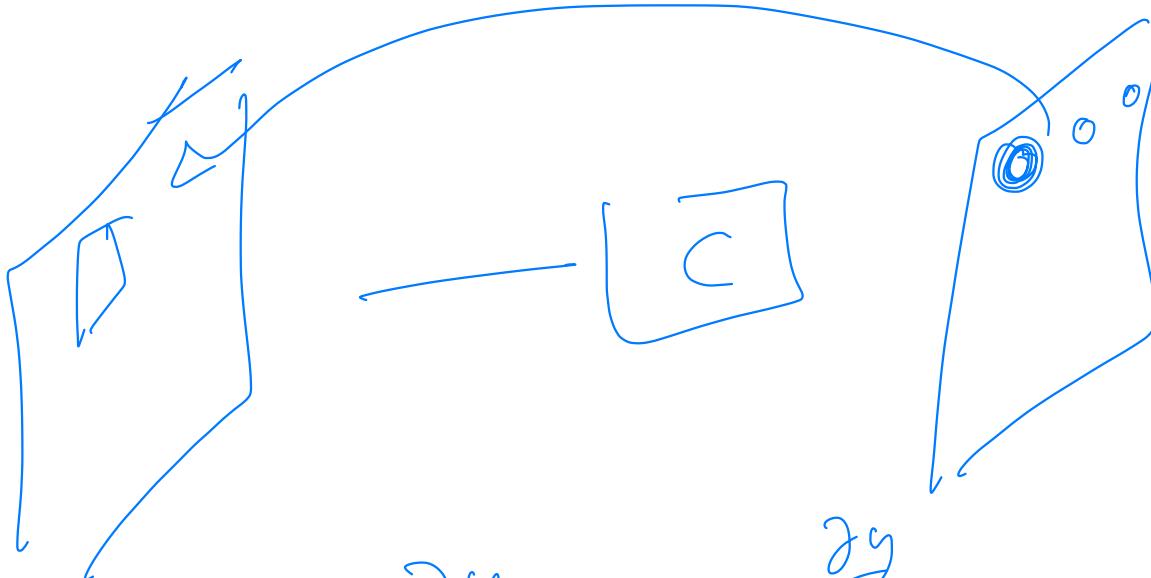
Batch mode

GD

$$w_i^l$$

$w$	$s_1$	$s_2$	$s_3$
$w_1$	#	-	+
$w_2$	-	+	+
$w_3$	+	-	-

SAD

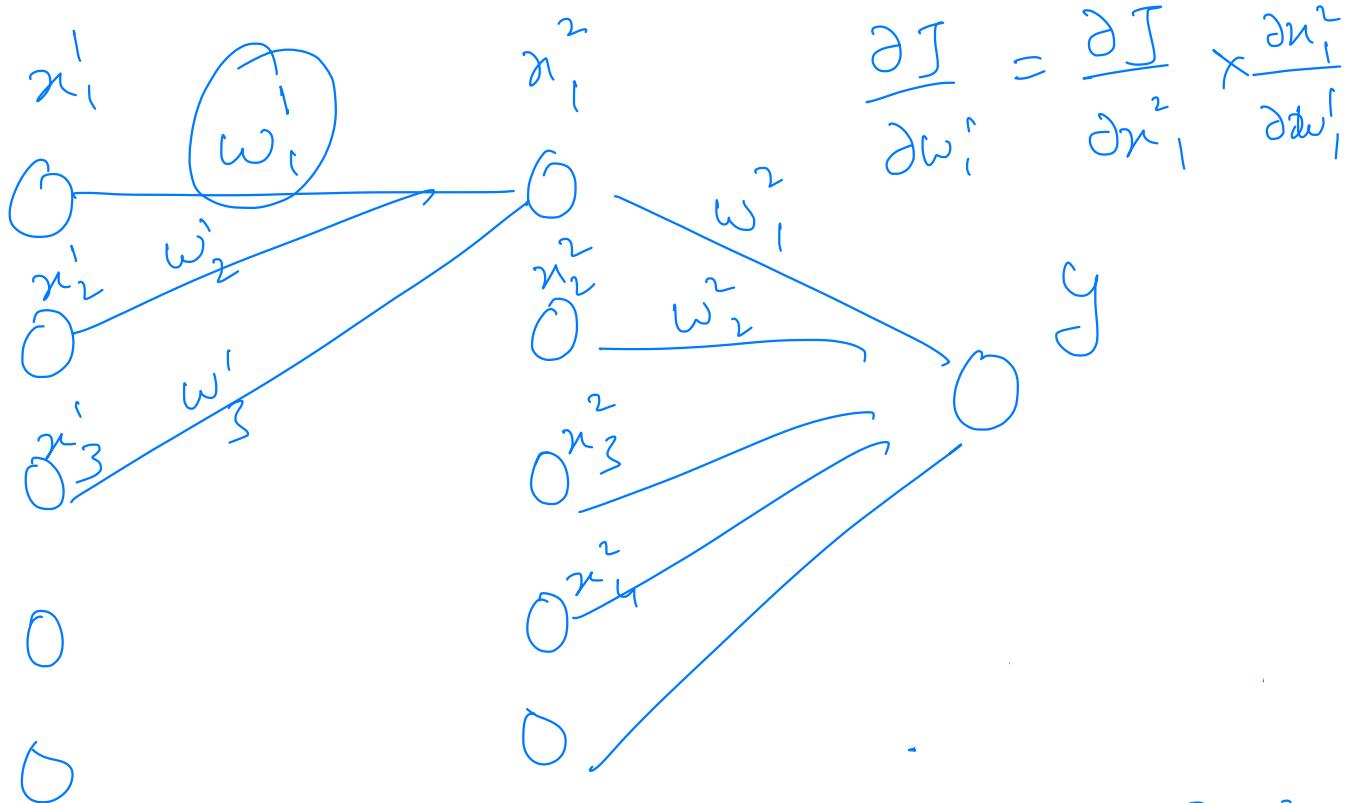


$$\frac{\partial y}{\partial w_i}$$

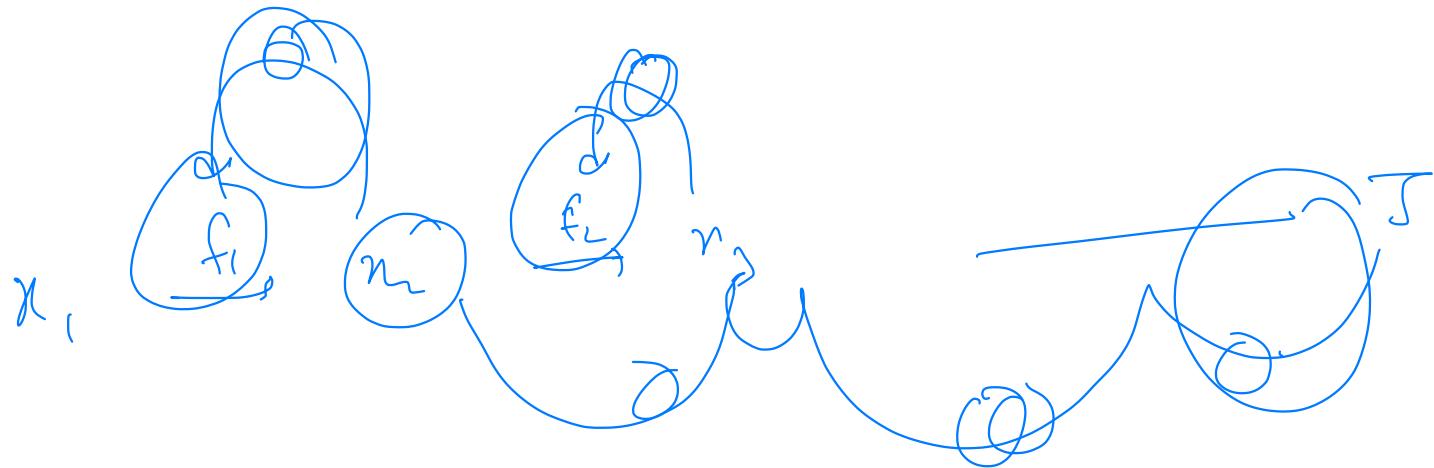
$$\frac{\partial y}{\partial w_j}$$

$$y = \text{softmax}(\mathbf{x})$$

$$S \left( \sum_{j \neq w} x_j \right)$$



$$\frac{\partial y}{\partial x_i^1} = \frac{\partial y}{\partial x_1^2} \times \frac{\partial x_1^2}{\partial x_i^1} + \frac{\partial y}{\partial x_2^2} \frac{\partial x_2^2}{\partial x_i^1}$$





# Backpropagation in General Cases

3. Use back-propagated error signals to compute gradients w.r.t. parameters only for the function elements with parameters where their derivatives w.r.t. parameters are known by symbolic computation.

$$\nabla_{\theta^{(l)}} J(\theta; x, y) = \frac{\partial}{\partial \theta^{(l)}} J(\theta; x, y) = \frac{\partial J}{\partial f^{(l)}} \frac{\partial f_{\theta^{(l)}}^{(l)}}{\partial \theta^{(l)}} = \delta^{(l)} \frac{\partial f_{\theta^{(l)}}^{(l)}}{\partial \theta^{(l)}}$$

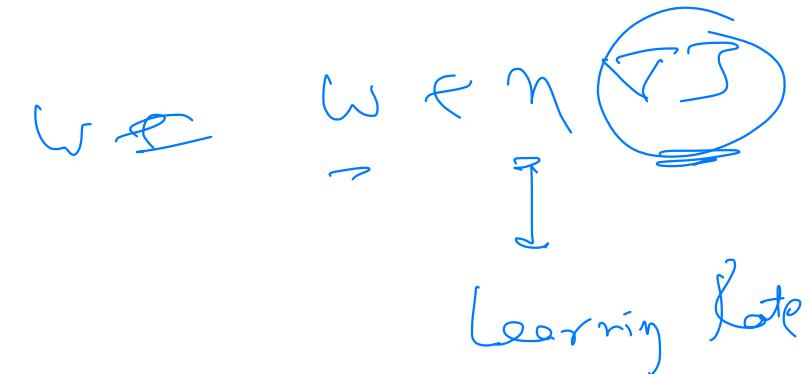
where  $\frac{\partial f_{\theta^{(l)}}^{(l)}}{\partial \theta^{(l)}}$  is known



# Backpropagation in General Cases

4. Sum gradients over all example to get overall gradient.

$$\nabla_{\theta^{(l)}} J(\theta) = \sum_{i=1}^m \nabla_{\theta^{(l)}} J(\theta; x^{(i)}, y^{(i)})$$



Output of penultimate layer =  $[a_1, \dots, a_k, \dots, a_k]$   $k = \# \text{ classes}$   
 Also called "logits" (we will soon describe why)

Softmax  $\Pr[y = i | a] =$

$$\frac{e^{a_i}}{\sum_{j=0}^k e^{a_j}}$$

denoted as " $\Sigma$ "

Sigmoid / Logistic function =  $\frac{1}{1 + e^{-z}}$

$$\Pr[y = o | a] = \frac{e^{a_o}}{e^{a_o} + e^{a_1}} = \text{Softmax for 2 classes.}$$

Softmax

$$= \frac{1}{1 + e^{-2}}$$

Output of penultimate layer =  $[a_1, \dots, a_k, \dots, a_k]$   $k = \# \text{ classes}$   
 Also called "logits" (we will soon describe why)

Softmax  $\Pr[g = i | a] = \frac{e^{a_i}}{\sum_{j=0}^k e^{a_j}}$  denoted as " $\Sigma$ "

Sigmoid / Logistic function =  $\frac{1}{1 + e^{-z}}$

= Softmax for 2 classes.

$$\Pr[g = o | a] = \frac{e^{a_o}}{e^{a_o} + e^{a_1}} = \frac{1}{1 + e^{-(a_o - a_1)}} = \frac{1}{1 + e^{-z}}$$

Softmax

Alternate form of softmax

$$\Pr[y=i|a] = \frac{e^{a_i}}{\sum_j e^{a_j}}$$

If we add a constant  $c$  to all the elements of "a" vector

$$\Pr[y=i|a] = \frac{e^{(a_i+c)}}{\sum_j e^{(a_j+c)}} = \frac{e^c \times e^{a_i}}{e^c \sum_j e^{a_j}} = \frac{e^{a_i}}{\sum_j e^{a_j}}$$

Hence we may choose the  $k^{th}$  logit as anchor and choose  $c$  such that  $e^{a_k+c} = e^0 = 1$  which makes the other  $k-1$  logits as independent.

$$\Pr[y=i|a] = \frac{e^{a_i}}{1 + \sum_{j=1}^{k-1} e^{a_j}}$$

Logits = Log odds

→ Imagine running K-1 independent binary logistic regression

models

→ Permute layer output =  $[a_1, \dots, a_x, \dots, a_K]$

→ Choose  $K^{\text{th}}$  output as the pivot

$$\text{log odds} = \log \left( \frac{\Pr(y=i)}{\Pr(y=K)} \right) = a_i$$

$$\Pr(y=i) = e^{a_i} \Pr(y=K)$$

We know that

$$P_n(y \geq k) = 1 - \sum_{i=1}^{k-1} P_n(y=i) = 1 - \sum_{i=1}^{k-1} P_n(y=k) e^{a_i}$$

$$\Rightarrow P_n(y=k) = \frac{1}{1 + \sum_{i=1}^{k-1} e^{a_i}}$$

$$\Rightarrow P_n(y=j) = \frac{e^{a_j}}{1 + \sum_{i=1}^{k-1} e^{a_i}}$$

Gradient of Softmax with respect to logits

$$\frac{\partial \xi_i}{\partial a_j} = \frac{\partial \frac{e^{a_i}}{\sum e^{a_i}}}{\partial a_j} = \frac{\partial \frac{e^{a_i}}{\sum}}{\partial a_j}$$

Let  $g = e^{a^i}$ ,  $h = \sum$

if  $f = \frac{g}{h} \Rightarrow f' = \frac{g'h - h'g}{h^2}$

$$\frac{\partial \xi_i}{\partial a_i} = \frac{e^{a_i} \sum - e^{a_i} e^{a_i}}{\sum \times \sum} = \frac{e^{a_i}}{\sum} - \frac{e^{a_i}}{\sum} \times \frac{e^{a_i}}{\sum}$$

$$\frac{\partial S_i}{\partial a_i} = \frac{e^{a_i}}{\sum} \left( 1 - \frac{e^{a_i}}{\sum} \right) = S_i (1 - S_i)$$

for  $\frac{\partial S_i}{\partial a_j}$

hence  $\frac{\partial S_i}{\partial a_j} = \frac{G - e^{a_j} e^{a_i}}{\sum^2} = \frac{-e^{a_j}}{\sum} \times \frac{e^{a_i}}{\sum}$

$$= -S_j S_i$$

## Cross Entropy

If  $q$  is the target probability &  $P$  is the predicted probability

$$CE(P, q) = - \sum_{i=1}^k q_{i,i} \log (P_i)$$

Recall that  $\lim_{x \rightarrow 0} x \log x = \lim_{n \rightarrow \infty} \frac{\log n}{n}$

Using L'Hospital Rule  $\lim_{n \rightarrow \infty} \frac{y_n}{-1/n^2} = 0$



# Convolution as Matrix Multiplication

$$\begin{matrix} 4 & 5 & 8 & 7 \\ 1 & 8 & 8 & 8 \\ 3 & 6 & 6 & 4 \\ 6 & 5 & 7 & 8 \end{matrix}$$

 $*$ 

$$\begin{matrix} 1 & 4 & 1 \\ 1 & 4 & 3 \\ 3 & 3 & 1 \end{matrix}$$

 $=$ 

$$\begin{matrix} 122 & 148 \\ 126 & 134 \end{matrix}$$



$$\begin{matrix} 4 & 5 & 8 & 7 & 1 & 8 & 8 & 8 & 3 & 6 & 6 & 4 & 6 & 5 & 7 & 8 \end{matrix}$$

 $T$



# Convolution as Matrix Multiplication

4	5	8	7
1	8	8	8
3	6	6	4
6	5	7	8

 $*$ 

1	4	1
1	4	3
3	3	1



122	148
126	134



1	4	1	0	1	4	3	0	3	3	1	0	0	0	0	0
	1	4	1	0	1	4	3	0	3	3	1	0			
				1	4	1	0	1	4	3	0	3	3	1	0
					1	4	1	0	1	4	3	0	3	3	1



# Convolution as Matrix Multiplication

1	4	1	0	1	4	3	0	3	3	1	0	0	0	0	0	0
	1	4	1	0	1	4	3	0	3	3	1	0				
				1	4	1	0	1	4	3	0	3	3	1	0	
				1	4	1	0	1	4	3	0	3	3	1	0	

 $*$ 

4
5
8
7
1
8
8
8
3
6
6
4
6
5
7
8

 $=$ 

122
148
126
134



# Backpropagation in Convolution Layer

Forward propagation

$$x * w = y$$

Backward propagation

$$\frac{\partial J}{\partial x} = \text{flip}(w) * \frac{\partial J}{\partial y}$$

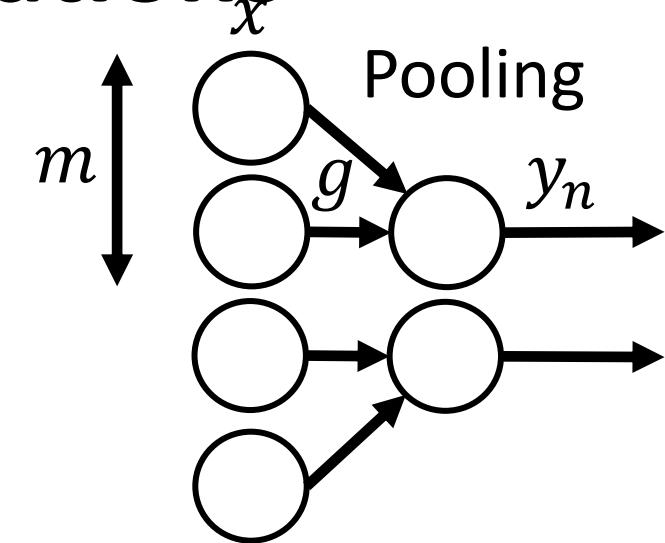
Gradient computation

$$x * \frac{\partial J}{\partial y} = \frac{\partial J}{\partial w}$$



# Pooling Layer: Subsampling Equations

Mean Pooling     $g(x) = \frac{\sum_{k=1}^m x_k}{m}$      $\frac{\partial g}{\partial x} = \frac{1}{m}$

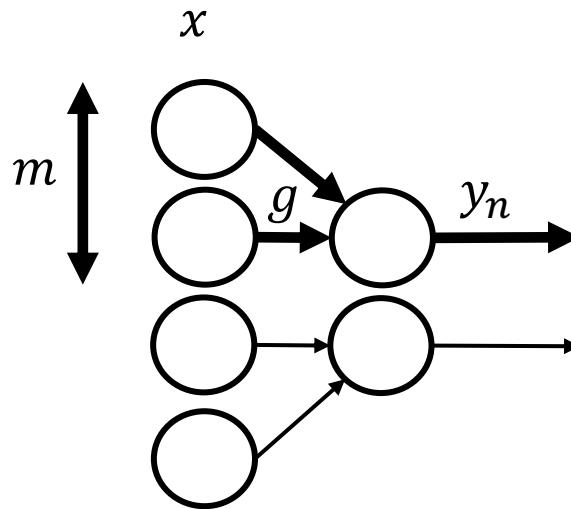


Max Pooling     $g(x) = \max(x)$      $\frac{\partial g}{\partial x_i} = \begin{cases} 1 & \text{if } x_i = \max(x) \\ 0 & \text{otherwise} \end{cases}$

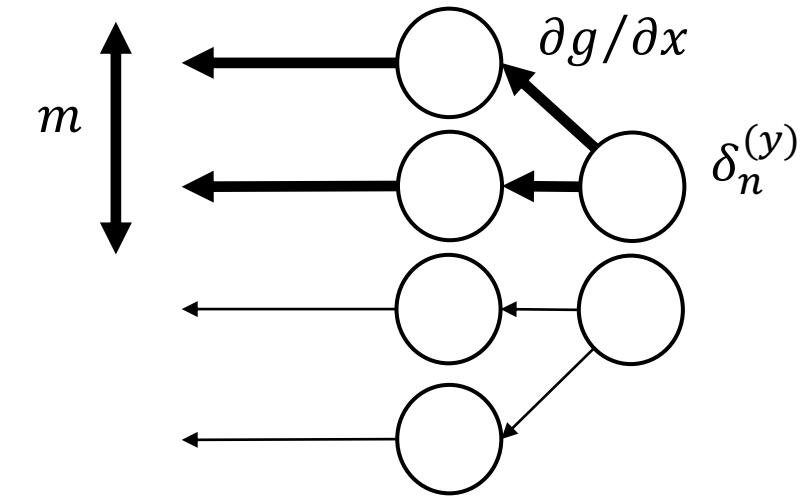
$L^p$  Pooling     $g(x) = \|x\|_p = \left( \sum_{k=1}^m |x_k|^P \right)^{\frac{1}{P}}$      $\frac{\partial g}{\partial x_i} = \left( \sum_{k=1}^m |x_k|^P \right)^{\frac{1}{P}-1} \cdot |x_i|^{p-1}$



# Backpropagation in the Pooling Layer



Forward propagation (subsampling)



Backward propagation (upsampling)