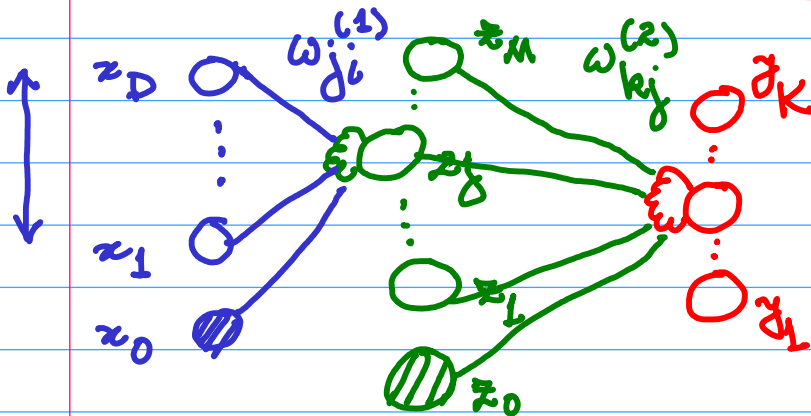


A representative example

## BACKPROPAGATION (contd.) [An example]



Assume

→  $h(\cdot)$ : activation function is  $\tanh(\cdot)$

→  $\sigma(\cdot)$ : activation function is a unit function

[2 layers of connections]

Why is this called 'Backpropagation'?



Training phase: error is at the output: this is fed back to update the weights

$n \in \{1, N\}$

FOR each pattern in the training set in turn, we follow the following operations.

① A 'forward propagation'

(j: hidden layer)

$$a_j^{(1)} = \underline{w}_j^{(1)T} \underline{x} = \sum_{i=0}^D w_{ji}^{(1)} x_i$$

$$z_j = \tanh(a_j)$$

$$y_k = \underline{w}_k^{(2)T} \underline{z} = \sum_{j=0}^M w_{kj}^{(2)} z_j$$

② Evaluate  $\delta_k$ 's for each output unit

$$\delta_k = y_k - t_k \quad \text{What is this, and how?}$$

$$E_n \triangleq \frac{1}{2} \sum_{k=1}^K (y_k - t_k)^2$$

$$\delta_k \triangleq \frac{\partial E_n}{\partial a_k} = \frac{1}{2} \cdot 2 (y_k - t_k) \left( \frac{\partial y_k}{\partial a_k} \right) = 1 \text{ as } y_k = a_k \text{ (} \sigma(\cdot) = \text{unit fn.} \text{)}$$

output layer activation

$\delta_k = y_k - t_k$  (Else, according to the specific activation function  $\sigma(\cdot)$  at the output layer)

③ Backpropagate these to obtain  $\delta_j$ 's for the hidden layer units

$$\delta_j = (1 - z_j^2) \sum_{k=1}^K w_{kj} \delta_k$$

What is this, and how?  
previous step (step #2)  
[output layer]

[hidden layer]

$$\delta_j \triangleq \frac{\partial E_n}{\partial a_j} = \sum_k \left( \frac{\partial E_n}{\partial a_k} \right) \left( \frac{\partial a_k}{\partial a_j} \right)$$

$\delta_k$  [step #2]

$a_k = \underline{w}^{(2)T} \underline{z} = \sum_{j=0}^M w_{kj}^{(2)} z_j = \sum_{j=0}^M w_{kj}^{(2)} h(a_j)$

$$\Rightarrow \frac{\partial a_k}{\partial a_j} = w_{kj}^{(2)} \frac{\partial h(a_j)}{\partial a_j} = w_{kj}^{(2)} h'(a_j) = w_{kj}^{(2)} (1 - z_j^2)$$

$$\delta_j = \sum_k \delta_k w_{kj}^{(2)} (1 - z_j^2) = (1 - z_j^2) \sum_k w_{kj}^{(2)} \delta_k$$

$$E_n \triangleq \frac{1}{2} \sum_{k=1}^K (y_k - t_k)^2$$

$$\delta_k \triangleq \frac{\partial E_n}{\partial a_k} = \frac{1}{2} \cdot 2 (y_k - t_k) \left( \frac{\partial y_k}{\partial a_k} \right) = 1 \text{ as } y_k = a_k \text{ (} \sigma(\cdot) = \text{unit fn.} \text{)}$$

output layer activation

$\delta_k = y_k - t_k$  (Else, according to the specific activation function  $\sigma(\cdot)$  at the output layer)

③ Backpropagate these to obtain  $\delta_j$ 's for the hidden layer units

$$\delta_j = (1 - z_j^2) \sum_{k=1}^K w_{kj} \delta_k$$

What is this, and how?  
previous step (step #2)  
[output layer]

[hidden layer]

$$\delta_j \triangleq \frac{\partial E_n}{\partial a_j} = \sum_k \left( \frac{\partial E_n}{\partial a_k} \right) \left( \frac{\partial a_k}{\partial a_j} \right)$$

$\delta_k$  [step #2]

$a_k = \underline{w}^{(2)T} \underline{z} = \sum_{j=0}^M w_{kj}^{(2)} z_j = \sum_{j=0}^M w_{kj}^{(2)} h(a_j)$

$$\Rightarrow \frac{\partial a_k}{\partial a_j} = w_{kj}^{(2)} \frac{\partial h(a_j)}{\partial a_j} = w_{kj}^{(2)} h'(a_j) = w_{kj}^{(2)} (1 - z_j^2)$$

$$\delta_j = \sum_k \delta_k w_{kj}^{(2)} (1 - z_j^2) = (1 - z_j^2) \sum_k w_{kj}^{(2)} \delta_k$$

# ④ Use the chain rule to evaluate the gradient

$$\frac{\partial E_n}{\partial \omega_{ji}^{(1)}} = \delta_j x_i, \quad \frac{\partial E_n}{\partial \omega_{kj}^{(2)}} = \delta_k z_j \quad \text{What is this, and how?}$$

Diagram: A pink box with a left-pointing arrow, containing a vertical column of four circles. To the right, an orange vector  $\bar{\nabla} E = \begin{bmatrix} \frac{\partial E_n}{\partial \omega} \end{bmatrix}$  is shown, with a double-headed arrow indicating it spans "all weights (1st & 2nd layers of connections)".

$$\frac{\partial E_n}{\partial \omega_{ji}^{(1)}} = \left( \frac{\partial E_n}{\partial a_j} \right) \frac{\partial a_j}{\partial \omega_{ji}^{(1)}}; \quad a_j = \underline{w}_j^{(1)T} \underline{x} = \sum_{i=1}^D \omega_{ji}^{(1)} x_i$$

$\delta_j$  (previous step: step #3)

$$\Rightarrow \frac{\partial a_j}{\partial \omega_{ji}^{(1)}} = x_i$$

$$\Rightarrow \frac{\partial E_n}{\partial \omega_{ji}^{(1)}} = \delta_j x_i$$

$$\frac{\partial E_n}{\partial \omega_{kj}^{(2)}} = \left( \frac{\partial E_n}{\partial a_k} \right) \frac{\partial a_k}{\partial \omega_{kj}^{(2)}}; \quad a_k = \underline{w}_k^{(2)T} \underline{z} = \sum_{j=1}^M \omega_{kj}^{(2)} z_j$$

$\delta_k$  (step #2)

$$\Rightarrow \frac{\partial a_k}{\partial \omega_{kj}^{(2)}} = z_j$$

$$\Rightarrow \frac{\partial E_n}{\partial \omega_{kj}^{(2)}} = \delta_k z_j$$

## Side topic: Numerical Evaluation of the gradient

Empirically, all of these alternative methods are numerically not as good as Backpropagation

$$(*) \quad \frac{\partial E_n}{\partial w_{ji}} = \frac{E_n(w_{ji} + \epsilon) - E_n(w_{ji})}{\epsilon} + o(\epsilon)$$

$\epsilon \longrightarrow (w_{ji} + \epsilon) - (w_{ji})$

OR: symmetrical central differences

$$\frac{\partial E_n}{\partial w_{ji}} = \frac{E_n(w_{ji} + \epsilon) - E_n(w_{ji} - \epsilon)}{2\epsilon} + o(\epsilon)$$

Why? The first direct formula?

$$E_n(w_{ji} + \epsilon) = E_n(w_{ji}) + \epsilon \frac{\partial E_n}{\partial w_{ji}} + \frac{1}{2} \epsilon H \epsilon + \text{higher order terms}$$

$$\Rightarrow \frac{E_n(w_{ji} + \epsilon) - E_n(w_{ji})}{\epsilon} = \frac{\partial E_n}{\partial w_{ji}} + \frac{1}{2} \epsilon H$$

Why? The symmetrical central differences formula:

$$E_n(w_{ji} + \epsilon) = E_n(w_{ji}) + \epsilon \frac{\partial E_n}{\partial w_{ji}} + \frac{1}{2} \epsilon H \epsilon + o(\epsilon^3)$$

$$E_n(w_{ji} - \epsilon) = E_n(w_{ji}) - \epsilon \frac{\partial E_n}{\partial w_{ji}} + \frac{1}{2} \epsilon H \epsilon - o(\epsilon^3)$$

$$E_n(w_{ji} + \epsilon) - E_n(w_{ji} - \epsilon) = 2\epsilon \frac{\partial E_n}{\partial w_{ji}} + 2o(\epsilon^3)$$

$$\frac{\partial E_n}{\partial w_{ji}} = \frac{E_n(w_{ji} + \epsilon) - E_n(w_{ji} - \epsilon)}{2\epsilon} + o(\epsilon^2)$$

Physical Significance: NN-based solution vis-a-vis the linear (or restricted linear) method done before.

(\*) NN: the weight parameters in the first layer are shared between the outputs

linear: each classification is performed independently.

(\*) The first layer of the network can be viewed as performing a non-linear feature extraction, and sharing the features between different outputs can lead to savings in computation and also lead to improved generalisation.

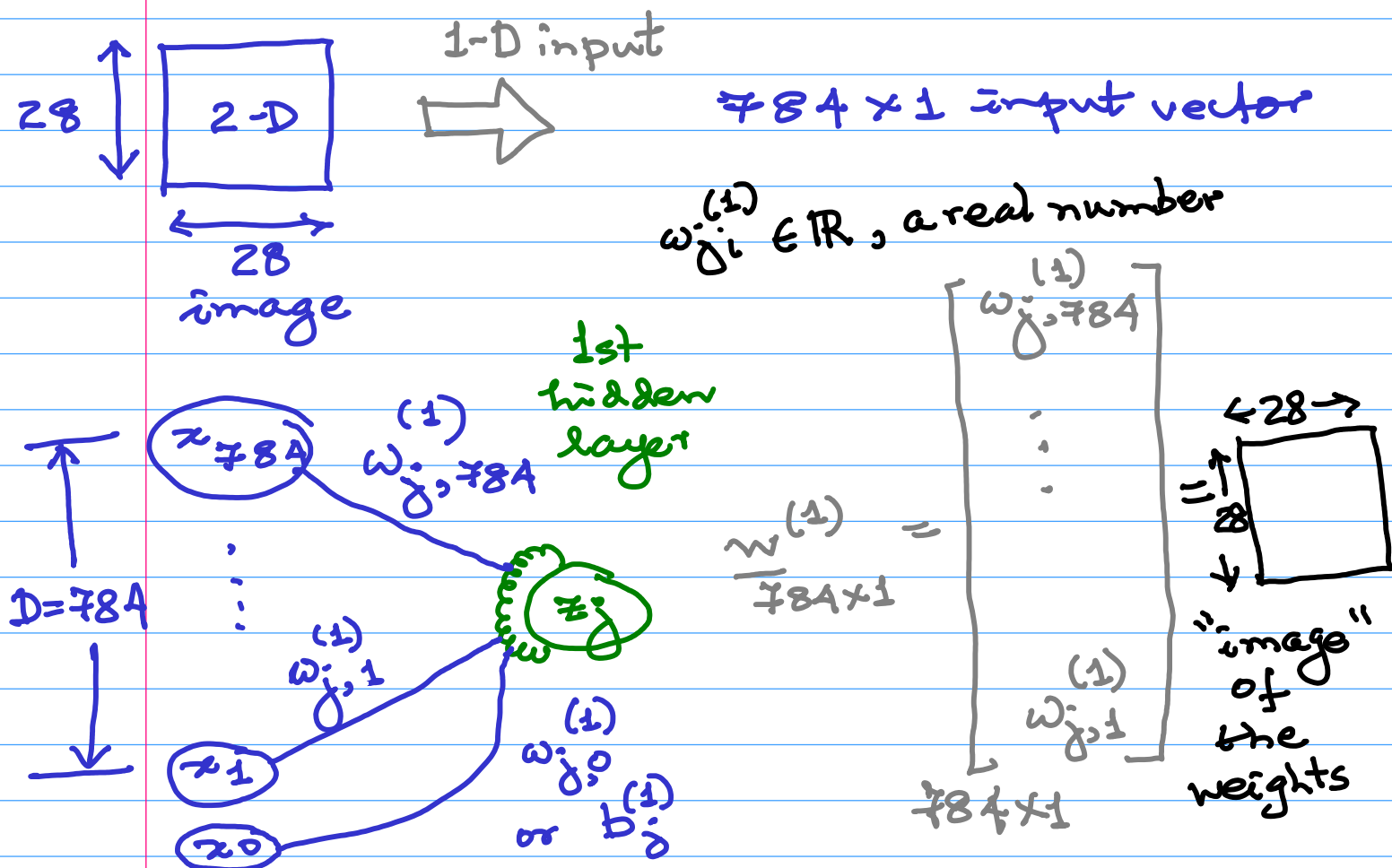
## Building Block (Input: 2-D: Image)

MNIST Numeral database:  $28 \times 28$  images  
(grayscale: not binary (not 0 and 255, or normalised 0 and 1))  
~ smooth

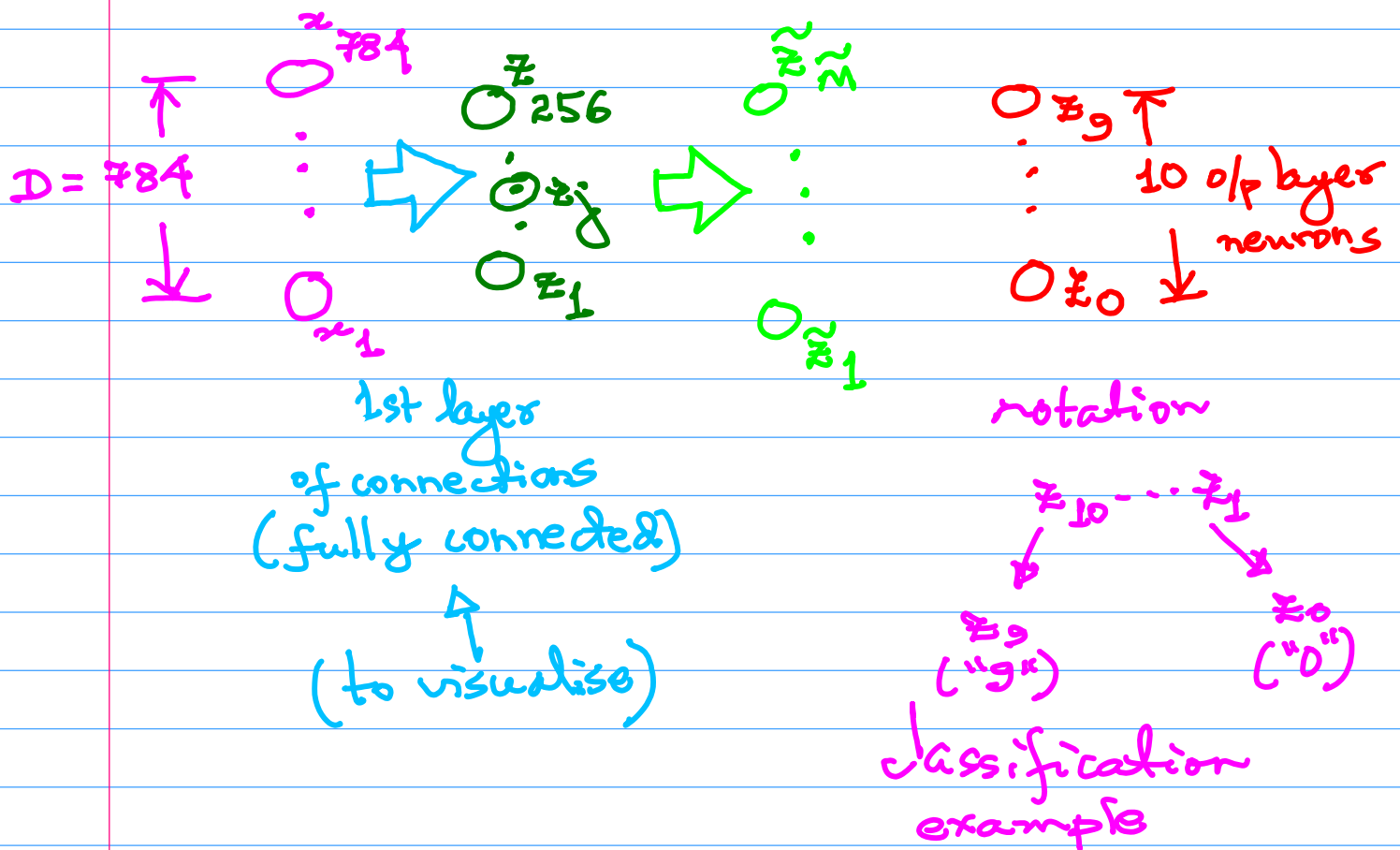
→ shades of grey as well, though most of the image is black or white.  
(0) (255, or 1, normalised)

Images of the 10 numerals: 0 to 9

Basic structure: an MLP with 2 hidden layers  
→ of specific interest is the first layer of connections

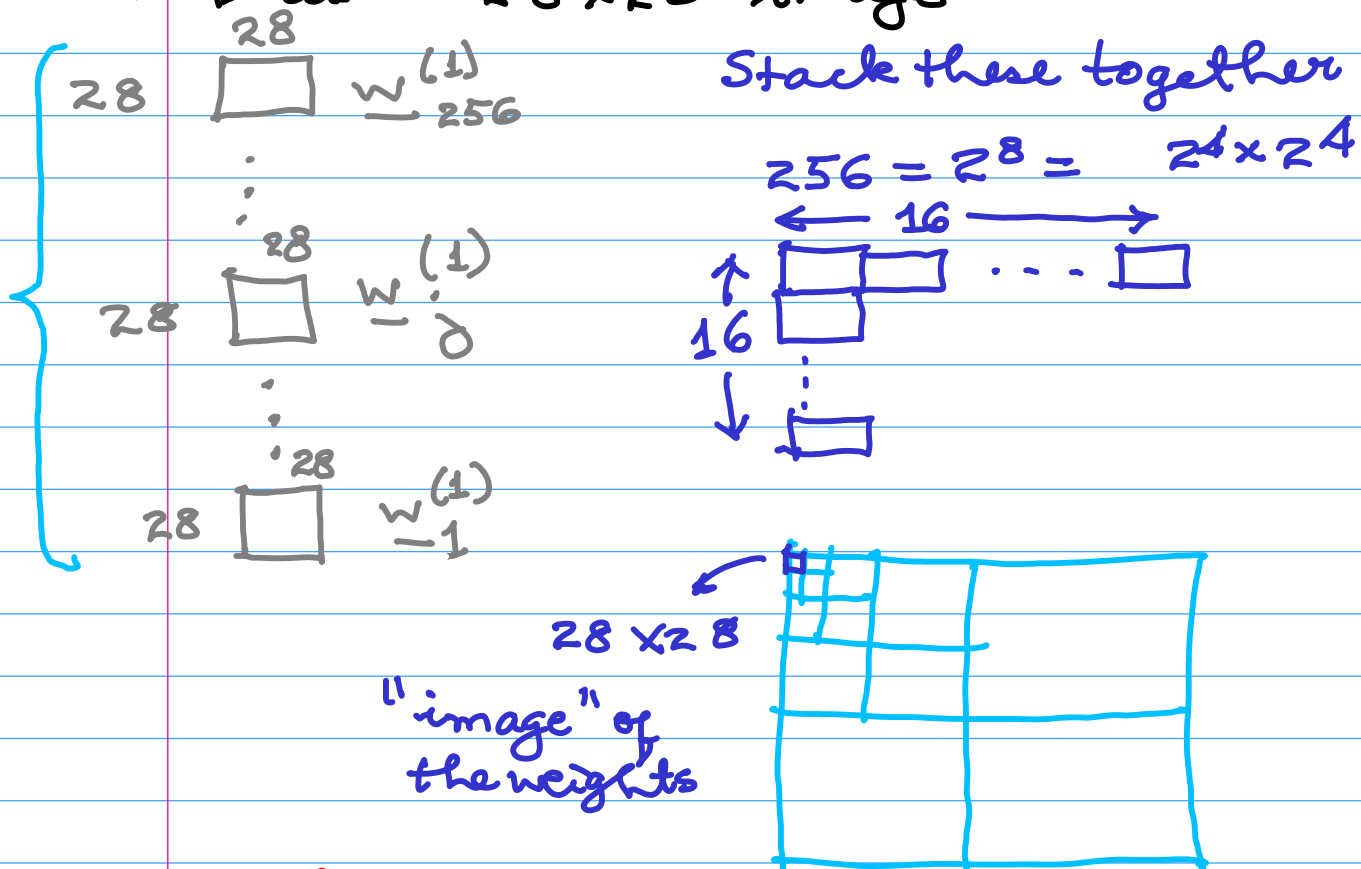


Take-home point: The input is not an ordinary 1-D vector, but a 2-D image  
 $\Rightarrow$  we can visualise the weights not as an ordinary 1-D vector, but a 2-D image with a similar spatial configuration / arrangement as the input itself.



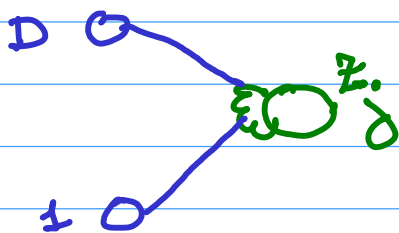


In this example, we are given 256 neurons in the first hidden layer. We will visualise the  $784 \times 1$  weight vector associated with each of the first hidden layer neurons e.g.  $\#j$ ,  $j \in \{1, 256\}$   $\rightarrow$  as a  $28 \times 28$  "image"



Empirical observation these "images" correspond to directional edge detectors, or directional derivatives, for a network which is subjected to weight learning.

1-D example



$$a_j^{(1)} = \underline{w}_j^{(1)} \underline{x} + b_j^{(1)}$$

or  $w_{0,j}^{(1)}$

$$= \sum_{i=1}^p w_{ji}^{(1)} x_i + b_j^{(1)}$$

Empirical observation  $\rightarrow$  most of the learnt weights are zero, except for a select few "close to" the neuron in question

e.g.,  $a_j^{(1)} = w_{j,l}^{(1)} x_l + w_{j,l+1}^{(1)} x_{l+1}$

and the other  $w_{ji}^{(1)}$ 's are  $\approx$  zero

and  $w_{j,l}^{(1)} = 1$ ,  $w_{j,l+1}^{(1)} = -1$

$\Rightarrow a_j^{(1)} = \underbrace{x_l - x_{l+1}}_{\text{difference}}$

$y = f(x)$

$$\frac{\partial f}{\partial x} = \lim_{\delta x \rightarrow 0} \frac{f(x + \delta x) - f(x)}{(x + \delta x) - (x)} = \lim_{\delta x \rightarrow 0} \frac{f(x + \delta x) - f(x)}{\delta x}$$

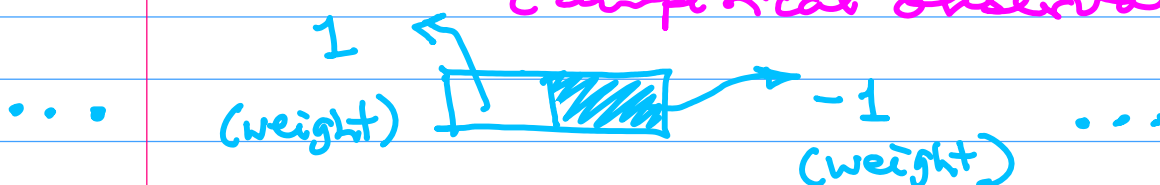
discrete case  $\delta x$ : 1 sampling point

$\delta x = 1$

derivate  $\rightarrow f[x+1] - f[x]$

discrete approximation of the derivative: in 1-D.

The first layer: computes a "local" derivative. (Empirical observation)



# CONVOLUTION

Neural Networks

$$\text{Activation } a = \underbrace{\underline{w}^T \underline{x}}_{\text{sum of pairwise products}} + b$$

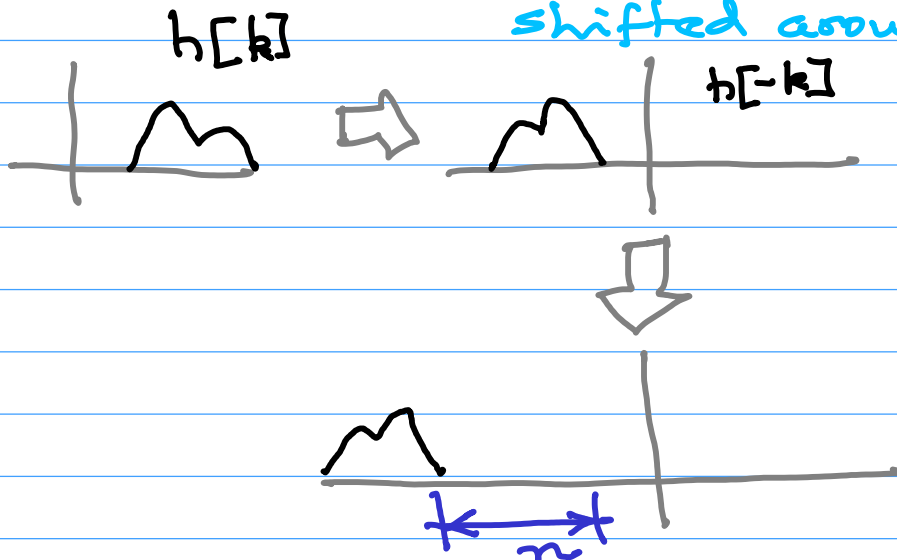
or  $w_0$

sum of pairwise products

[ Discrete domain : easier to understand ]  
Formula for Convolution:

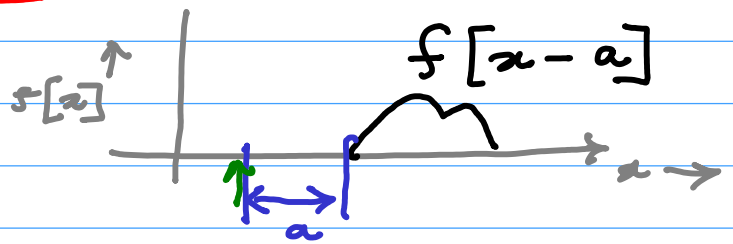
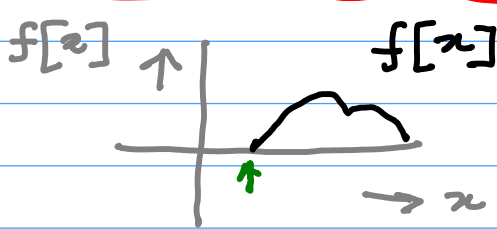
$$x[n] * h[n] \triangleq \sum_k \underbrace{x[k]}_{\text{input}} \underbrace{h[n-k]}_{\substack{\text{"mask" or} \\ \text{system} \\ \text{response}}}$$

'flipped' and then shifted around

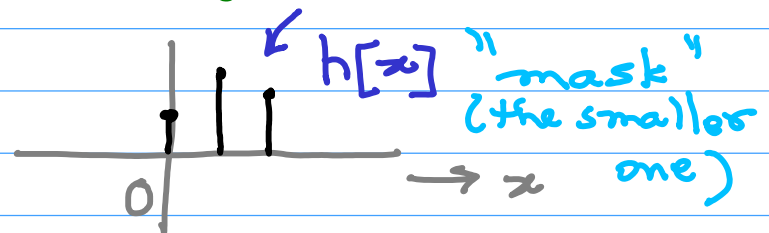
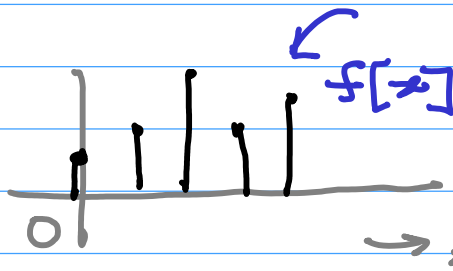


# Notes on Convolution:

1-D discrete functions

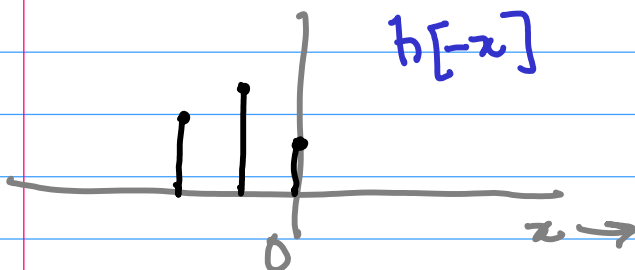


"delayed function"



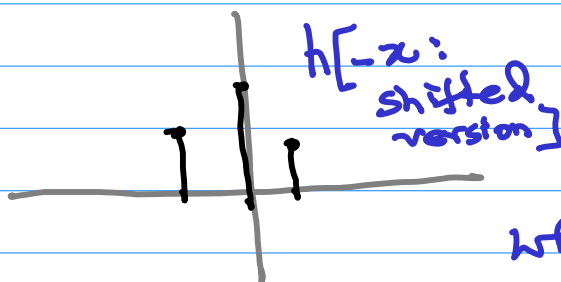
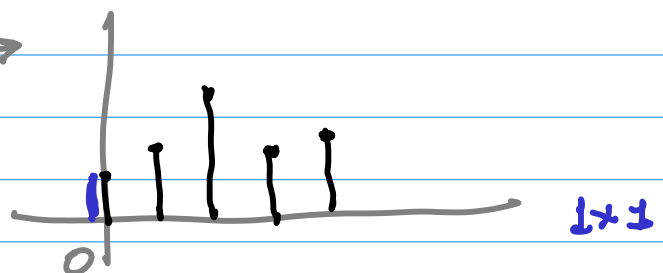
$$f[x] * h[x] = \sum_k f[k] \underbrace{h[x-k]}$$

$h[-k]$ : flipped mask



$$h[-k+x]$$

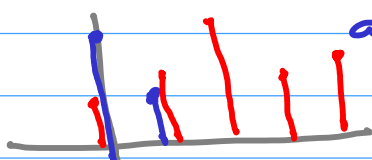
shift it all around



$h[-x]$ : shifted version

shifting all around & looking for overlaps.

Wherever there is an overlap, we perform pairwise mult & sum all of them up



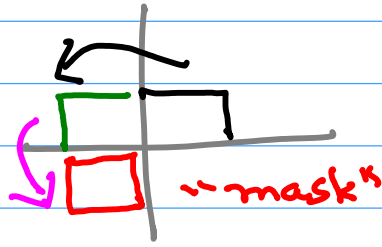
How many overlaps

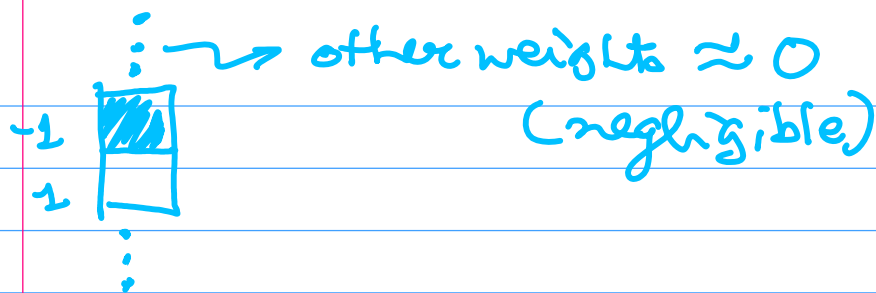
$$A + B - 1 \text{ (Brigham's Result)}$$

Efficient manner: using arrays

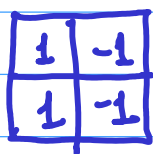
2-D: Conv.

Correlation:  
no flipping!

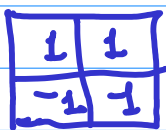
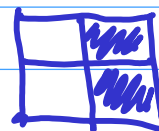




2-D:

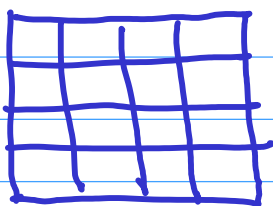


Captures a "vertical" edge



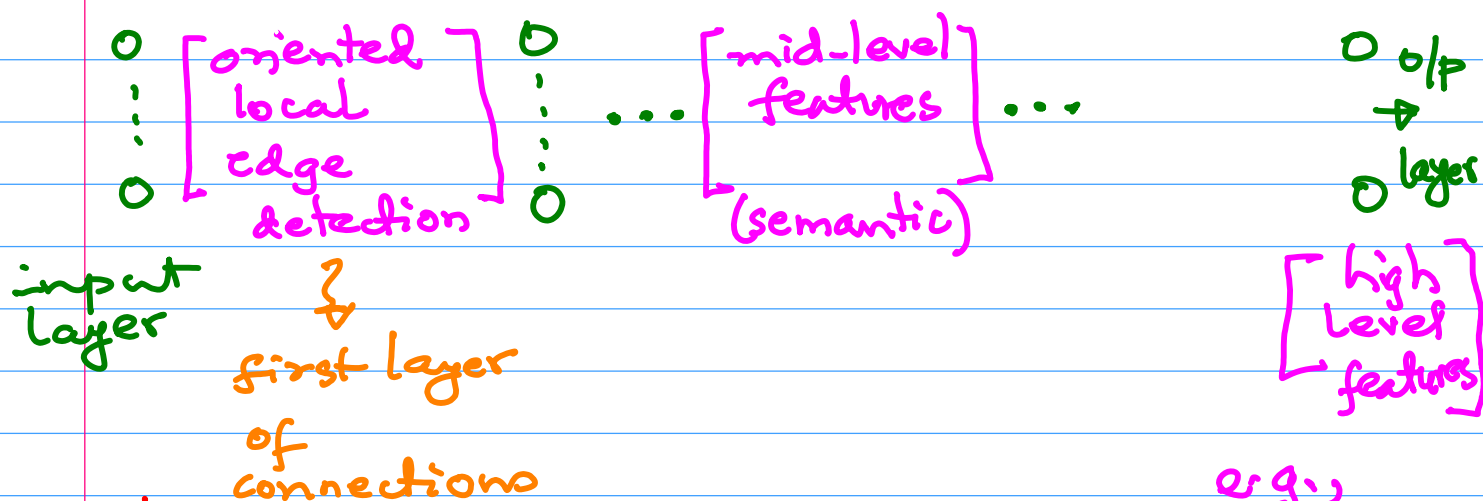
Captures a "horizontal" edge

larger neighbourhood



capture edges in different orientations.

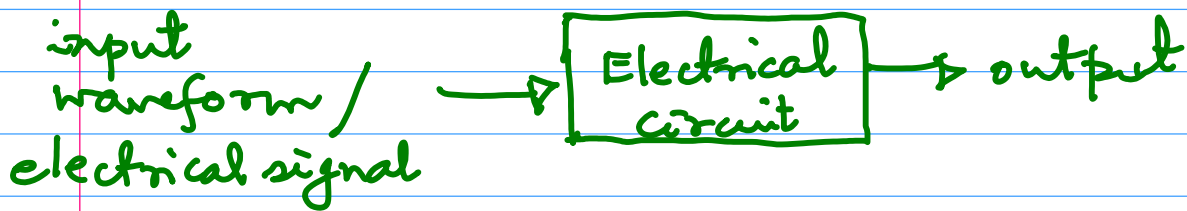
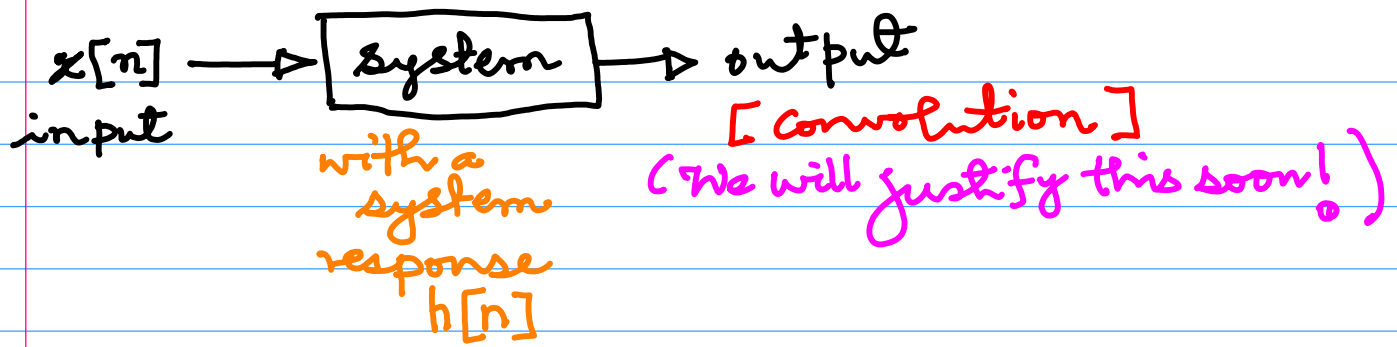
## Interpretation: (deep network)



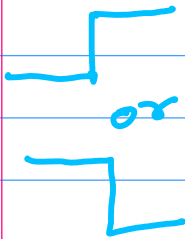
also corresponds  
 to the action in V1/V2  
 of the visual cortex

- $\Rightarrow$
- 1) simple derivative features
  - 2) only "local", not global

e.g.,  
 image of  
 a cat/dog



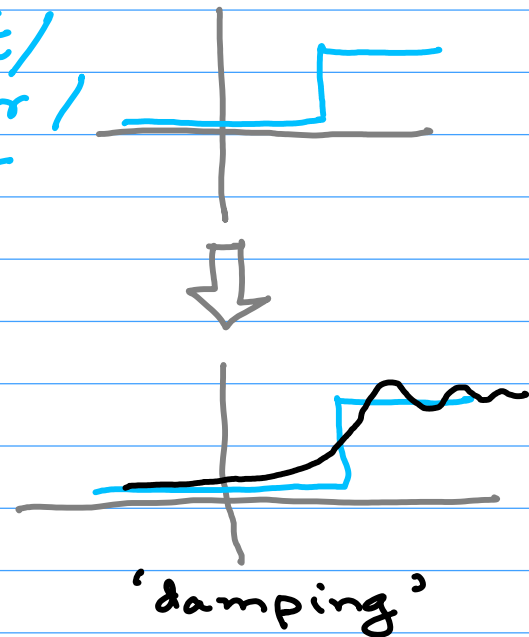
(pothole)



or

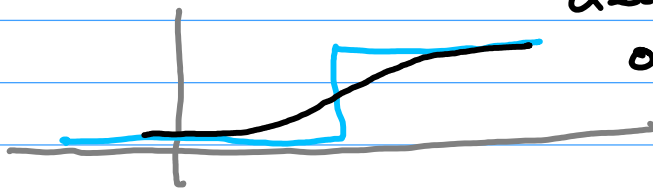
smoothing agent /  
smoothing filter /  
Low-pass filter

(unit step function /  
extreme case of a  
sigmoid)



'damping'

good suspension



dead-beat /  
overdamped  
system

(On the other hand, for a chronometer: a time-keeping device, we want an underdamped system)

# LSI SYSTEMS (LINEAR SHIFT INVARIANT) SYSTEMS

Approximation

Electrical: diode/BJT linear region

Mechanical: Mass-spring system

$$F = -kx$$

↳ the nature of the output is the same, just shifted in time/space

e.g., weighing scales with a tray to keep objects

Why are people obsessed with LSI systems?

\* Many practical systems can be approximated by LSI systems

(EE specific:

- LSI permit an equivalent freq domain <sup>analysis</sup>
- complex exponentials / sinusoids are eigenfunctions of LSI systems)

What is linearity? [ — Additivity  
— Homogeneity

$$x_1[n] \rightarrow y_1[n]$$

$$x_2[n] \rightarrow y_2[n]$$

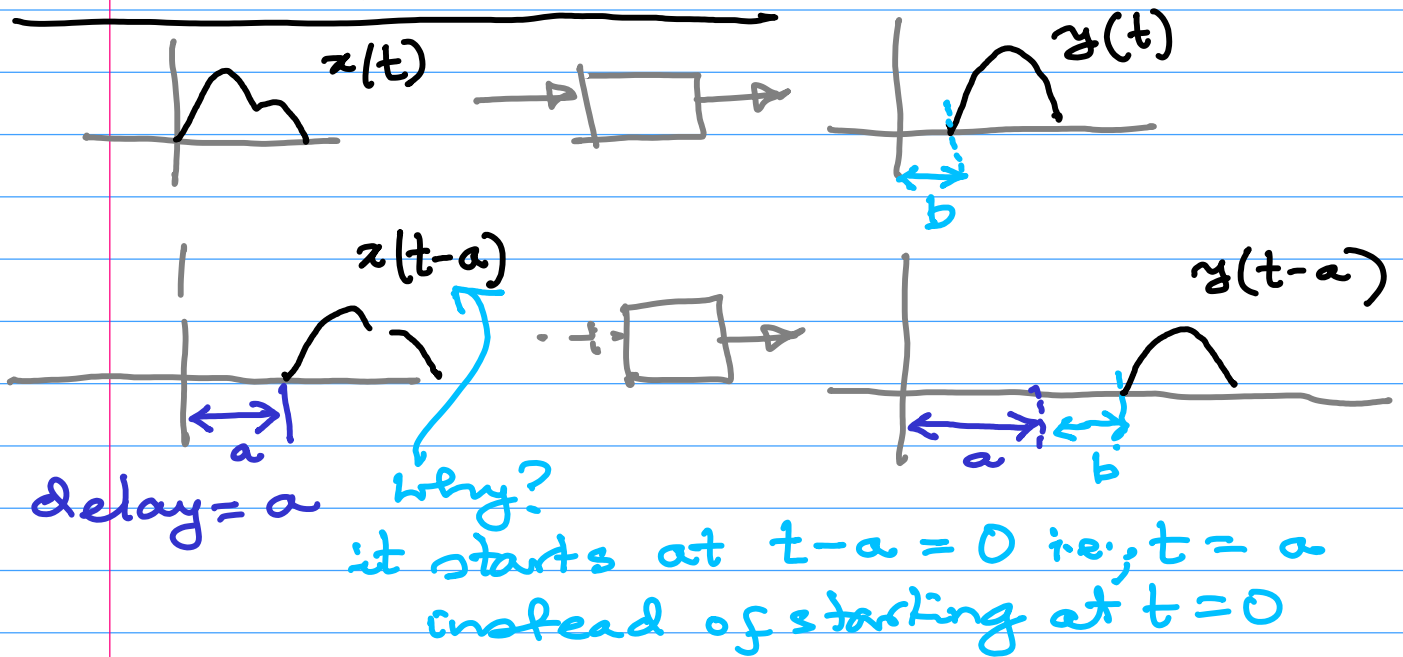
Additivity  $\Rightarrow x_1[n] + x_2[n] \rightarrow y_1[n] + y_2[n]$

Homogeneity  $\Rightarrow \alpha x_1[n] \rightarrow \alpha y_1[n]$

"Superposition Principle"



## What is Shift Invariance?



## How do we 'characterise' a system?

can be constructed out of different classes of

- e.g.: electronic component <
  - vacuum tubes
  - transistors
- e.g.: mechanical components (time measurement)
  - pendulum
  - flywheel

we need a  
"standardisation  
of the equivalence"

"standard input" → system → output

e.g., monitor (CRT/LCD/LED/Plasma)  
input [switching on]

characterisation: how quickly does it respond

This output is a waveform which "characterises" the system irrespective of the exact internal hardware construction (as a black box)

→ This allows us to treat a system (electrical/mechanical) as a waveform itself

→ The name of the EE course "Signals and Systems"



Impulse  $\delta$

(Kronecker Delta / 'bar')  
This is a standard input

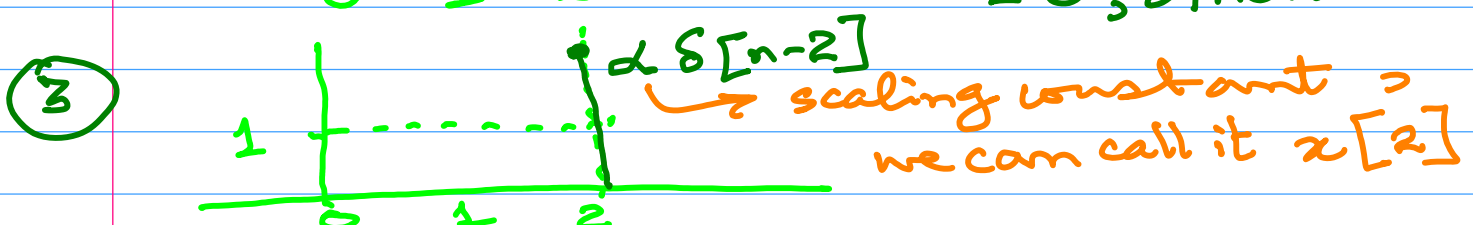
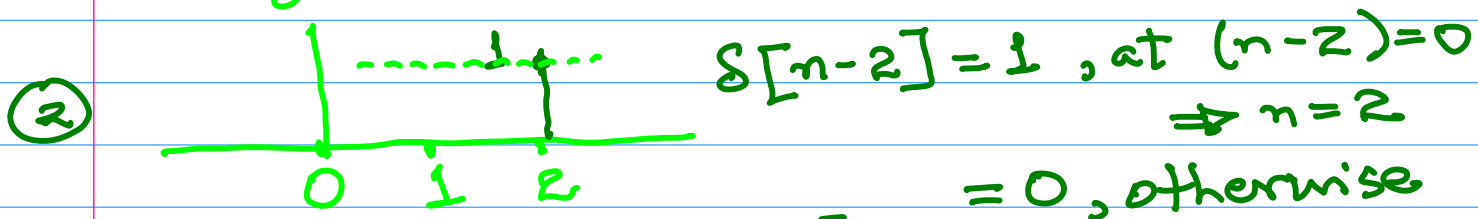
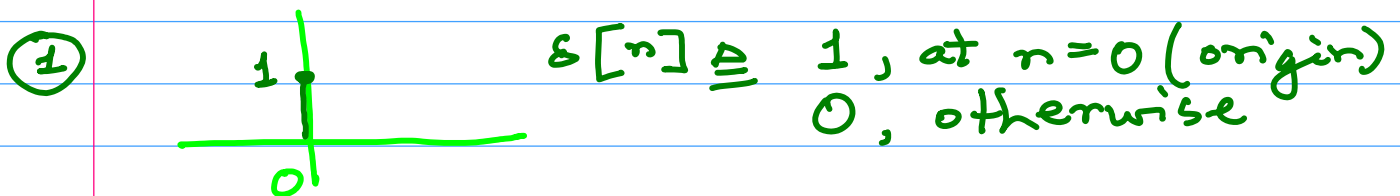
'Impulse Response'

is the characterisation of the system

Civil Engg example: Bridge: (Suspension bridge / ... / )

input: perturbation

output: vibration response



**Step I**  $\delta[n] \rightarrow \boxed{\text{LSI}} \rightarrow h[n]$  "Impulse Response"

**Step II**  $\delta[n-k] \rightarrow \boxed{\text{LSI}} \rightarrow h[n-k]$

delayed the impulse by 'k' units

(Shift Invariance)

i.e., impulse is not at  $n=0$ , but at  $k$

**Step III**  $x[k] \delta[n-k] \rightarrow \boxed{\text{LSI}} \rightarrow x[k] h[n-k]$

The impulse at  $k$  is no longer a unit impulse, but has a magnitude  $x[k]$

(Homogeneity)

**Step IV**  $\sum_k x[k] \delta[n-k] \rightarrow \boxed{\text{LSI}} \rightarrow \sum_k x[k] h[n-k]$

This is the input waveform itself! A collection of scaled unit impulses which are appropriately delayed.

(Additivity)

→ this is the definition of Convolution itself!

