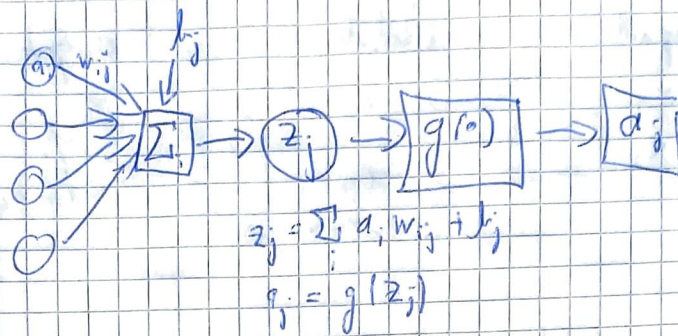


# Deep Learning 1 (Structured Networks)



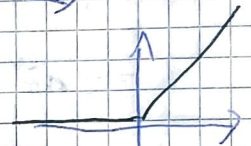
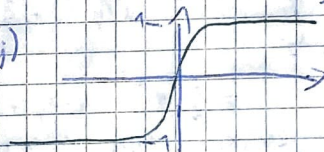
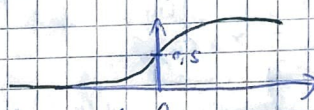
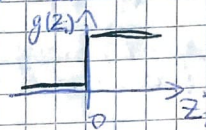
Nonlinearity through non-linear activation functions:

Threshold:  $g(z_j) = 1(z_j \geq 0)$

Sigmoid:  $g(z_j) = \frac{e^{z_j}}{1 + e^{z_j}}$

Tanh:  $g(z_j) = \tanh(z_j)$

ReLU:  $g(z_j) = \max(0, z_j)$



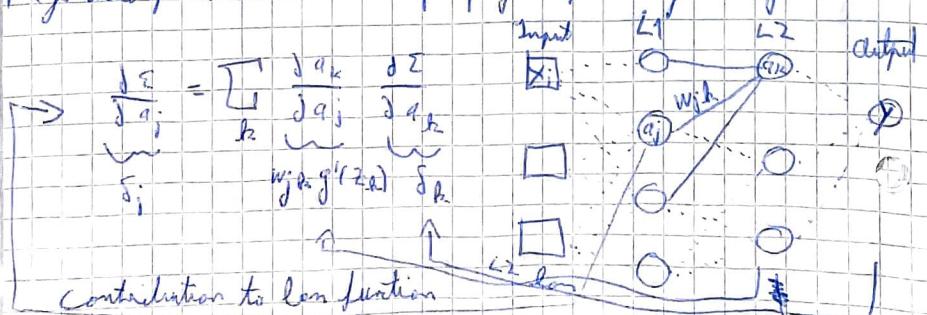
Learning:

Take partial derivative of objective  $\mathcal{E}$  with respect to  $w_{ij}$  to learn influence of  $w_{ij}$  on  $\mathcal{E}$

$$\frac{\partial \mathcal{E}}{\partial w_{ij}} = \frac{\partial z_j}{\partial w_{ij}} \cdot \frac{\partial a_j}{\partial z_j} \cdot \frac{\partial \mathcal{E}}{\partial a_j} = a_j \cdot g'(z_j) \cdot \frac{\partial \mathcal{E}}{\partial a_j}$$

Backpropagation Algorithm / Backward Pass:

The gradient from the error can be propagated from layer to layer with chain rule.



Implementation  
Store intermediate  
 $\delta_j, \delta_k, \dots$  results to  
save computation  
time!



Neural Networks with sufficiently many neurons  
can approximate any function  $f(x)$   $x \in x_1, x_2, \dots, x_d$

~~X~~

~~X~~

### Disadvantages of Neural Networks:

- non-convex
- many hyperparameters (initialisation, learning rate, ...)
- multiple layers can cause pathological curvature
- optimiser might get stuck on local optima

Universality: CNN model recognises most visual objects

Compactness: Representation finite dimensional at each layer

Convexity: Non-convex, but usually converges to good solutions

|                     | Universal | Compact | Convex |
|---------------------|-----------|---------|--------|
| Feature Engineering | X         | ✓       | ✓      |
| Kernels             | ✓         | X       | ✓      |
| Neural Networks     | ✓         | ✓       | X      |

### Convolutional Neural Networks

- consist of convolution and pooling layers
- backpropagation of errors
- CNNs are state-of-the-art (by wide margin) for image classification
- CNN is very resource expensive → pretrained model

#### Standard NN

- neurons and weights are scalars,  $a_i, w_{ij} \in \mathbb{R}$

- mapping between two layers is a sum of products:  $z_j = \sum_i w_{ij} a_i + b_j$

#### Convolutional NN

- neurons & weights are 2D arrays:  $a^{(i)} \in \mathbb{R}^{W \times H}$ ,  $w^{(ij)} \in \mathbb{R}^{W' \times H'}$
- mapping between two layers is sum of 2D cross-correlation  $\otimes$   
 $z_j = \sum_i w^{(ij)} \otimes a^{(i)} + b_j$
- pooling layers to reduce spatial resolution and increase # neurons



## Summary:

- NNs can learn compact representation of any task
- but more difficult optimization problem
- structured networks allow handling of real-world data (images, text, graphs)
- CNNs extract image representation by separating semantic content while compressing spatial content
- GNNs generalize NNs to any graph-structured input (trees, lattices, ...)

## Step 1: Cross-correlation:

$$z_{t+s} = [w * a]_{t+s} = \sum_{s'=0}^n \sum_{s''=0}^n w_{s'} \cdot a_{t+s, t+s'+s''}$$

## Step 2: The Convolution Layer:

$$\tilde{a}^{(j)} = g\left(\sum_i w^{(j,i)} * a^{(i)}\right)$$

Forward pass:  $z_t = [w * a]_t = \sum_{s=-\infty}^{\infty} w_s \cdot a_{t+s}$

Backward pass:  $\frac{\partial \mathcal{E}}{\partial a_n} = \sum_{t=-\infty}^{\infty} \frac{\partial \mathcal{E}}{\partial z_t} \cdot \frac{\partial z_t}{\partial a_n} = \sum_{t=-\infty}^{\infty} \sum_{s=-\infty}^{\infty} \frac{\partial \mathcal{E}}{\partial z_t} \cdot w_s \cdot 1_{t=t+s=n}$

$$= \sum_{t=-\infty}^{\infty} \frac{\partial \mathcal{E}}{\partial z_t} \cdot w_{n-t} = \left[ \frac{\partial \mathcal{E}}{\partial z} * w \right]_n$$

## Parameter Gradient:

$$\frac{\partial \mathcal{E}}{\partial w_n} = \sum_{t=-\infty}^{\infty} \frac{\partial \mathcal{E}}{\partial z_t} \frac{\partial z_t}{\partial w_n} = \sum_{t=-\infty}^{\infty} \sum_{s=-\infty}^{\infty} \frac{\partial \mathcal{E}}{\partial z_t} \cdot a_{t+s} \cdot 1_{t+s=n} = \sum_{t=-\infty}^{\infty} \frac{\partial \mathcal{E}}{\partial z_t} \cdot a_{n-t}$$

$$= \left[ \frac{\partial \mathcal{E}}{\partial z} * a \right]_n$$

## Text Input: Vector Embedding of a Sentence.

1. One-hot encoding: Represent word with indicator function

$$e[w] = \begin{pmatrix} 1_{w=w_1} \\ 1_{w=w_2} \\ \vdots \\ 1_{w=w_c} \end{pmatrix} \in \mathbb{R}^{\# \text{ words}}$$

2. PCA embedding:

structured kernel measures similarity between words (co-occuring)

$$e[w] = u_{w,d} \odot (1, \lambda^{0.5}, \dots, \lambda^{0.5}, 1)_d \in \mathbb{R}^d$$

3. Learned Embedding: Start with random embedding  $e[w] \in \mathbb{R}^d$

and learn embedding:  $e[w] \leftarrow e[w] - \eta \frac{\partial \mathcal{E}}{\partial e[w]}$

Recursive Neural Networks: make use of parsing tree of sentences

Graph Neural Networks: - specialised for classifying graphs

- graph as input, propagation step at each layer along edges
- CNN special case of GNN where graph is 2D lattice

