# Neural networks' architectures

$\vec{x}$ - input , $y$ - output

$$f(\vec{x}) \approx y$$

Ex. 1. Regression $\quad y \in \mathbb{R} \quad f(\vec{x}) \in \mathbb{R}$

$$(f(\vec{x}) - y)^2 \to \min$$

2. Classification $\quad y \in \{1, \ldots, c\}, \quad f(\vec{x}) \in \mathbb{R}^c$

$$\vec{p} = f(\vec{x}), \quad \vec{p} = \{p_i\}_{i=1}^{c}, \quad p_i \geq 0 \quad \sum p_i = 1$$

$$y = c \to \begin{pmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{pmatrix} = \vec{y}$$

$$- \sum y_i \log p_i \to \min \quad - \text{cross-entropy}$$

$$f(\vec{x}) \in \mathcal{F}$$

$$\mathcal{F} = \{ f(\vec{x}; \vec{\Theta}) = f_{\vec{\Theta}}(\vec{x}) = f(\vec{x} \mid \vec{\Theta}), \vec{\Theta} \in (H) \subseteq \mathbb{R}^P \}$$

$$\underset{\Theta}{f(\vec{x})} = \vec{\Theta}^T \cdot \vec{x}$$

<span style="color:orange">I layer, $\{\vec{\Theta}_i\}_{i=1}^{c}$, $c \cdot d$ parameters</span>

<span style="color:green">logit</span>

$$f_{\vec{\Theta}}(\vec{x}) = \text{softmax}\left( \{ \vec{\Theta}_i^T \vec{x} \} \right), \quad i = 1 \ldots c$$

$$\underset{\vec{f}}{\vec{p}} = \text{softmax}(\vec{\ell}) = \left\{ \frac{\exp(\ell_i)}{\sum\limits_{j=1}^{c} \exp(\ell_j)} \right\}$$

<span style="color:orange">II layer, nonlinear layer, no parameters</span>

<span style="color:green">$$\vec{\Theta}_1^T \vec{x} \quad < \quad \vec{\Theta}_2^T \vec{x} \qquad 0 < (\vec{\Theta}_2 - \vec{\Theta}_1)^T \vec{x}$$</span>

<span style="color:green">hyperplane</span>

Multilayer (fully-connected) neural network

$$\sigma(t) = \frac{1}{1 + \exp(-t)}$$

Ex. 2 layer neural network

I layer $\quad \{ (\vec{\Theta}_i)^T \vec{x} \} \to (H)_1 \vec{x} \to \vec{\ell}_1 \to \sigma(\vec{\ell}_1) \to \vec{z}_1$

II layer $\quad \{ (\vec{\Theta}_i)^T \vec{z}_1 \} \to (H)_2 \vec{z}_1 \to \vec{\ell}_2 \to \sigma(\vec{\ell}_2) \to \vec{z}_2$

$\vdots$

k layer $\quad \{ (\vec{\Theta}_i)^T \vec{z}_{k-1} \} = (H)_k \vec{z}_{k-1} \to \vec{\ell}_k \to \sigma(\vec{\ell}_k) \to \vec{z}_k$

<span style="color:orange">softmax $(\vec{\ell}_k) \to \vec{p}$</span>

$$\left[ \begin{array}{l} \Theta \sim \text{theta} \qquad\qquad t \quad T \\ (H) - \text{capital theta} \quad T \quad \text{JTI} \end{array} \right.$$

$$L(\mathbb{H}) = \frac{1}{m} \sum_{i=1}^{m} CE(\vec{p_i}, y_i) \to \min_{\mathbb{H}}, \quad \mathbb{H} = \{\mathbb{H}_1, \dots, \mathbb{H}_k\}$$

$$\hookrightarrow \frac{\partial L(\mathbb{H})}{\partial \mathbb{H}} \to \text{apply a variant of gradient ascent}$$

Backpropagation to calculate derivatives

$\mathcal{F}$ – big enough? - optimization -good or not?

### Universal approximation theorem

$$k \geq 2 \quad - \quad |g(\vec{x}) - f_{NN}(\vec{x})| < \varepsilon \quad \forall \vec{x} \in X \subset \mathbb{R}^d$$

$$d_1 = \dim(\vec{z_1}) \qquad \mathbb{H}_1 \in \mathbb{R}^{d_1 \times d} \qquad \vec{z_k} - \text{a representation}$$
$$\text{of an object at layer } k$$

$$\text{soft max}\left(\mathbb{H}_2 \, g(\mathbb{H}_1 \vec{x})\right)$$
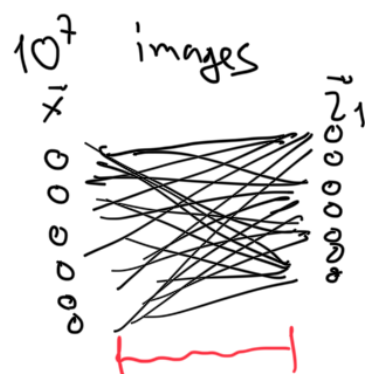$$\vec{z_0} = \vec{x}$$
$$d_0 = d$$

$$P = \sum_{i=1}^{k} d_i \, d_{i-1} = \underline{k \cdot d^2}, \qquad m - \text{a sample size}$$

$\vec{x}$ – an image $\quad 32 \times 32 = 1024 = d_0$

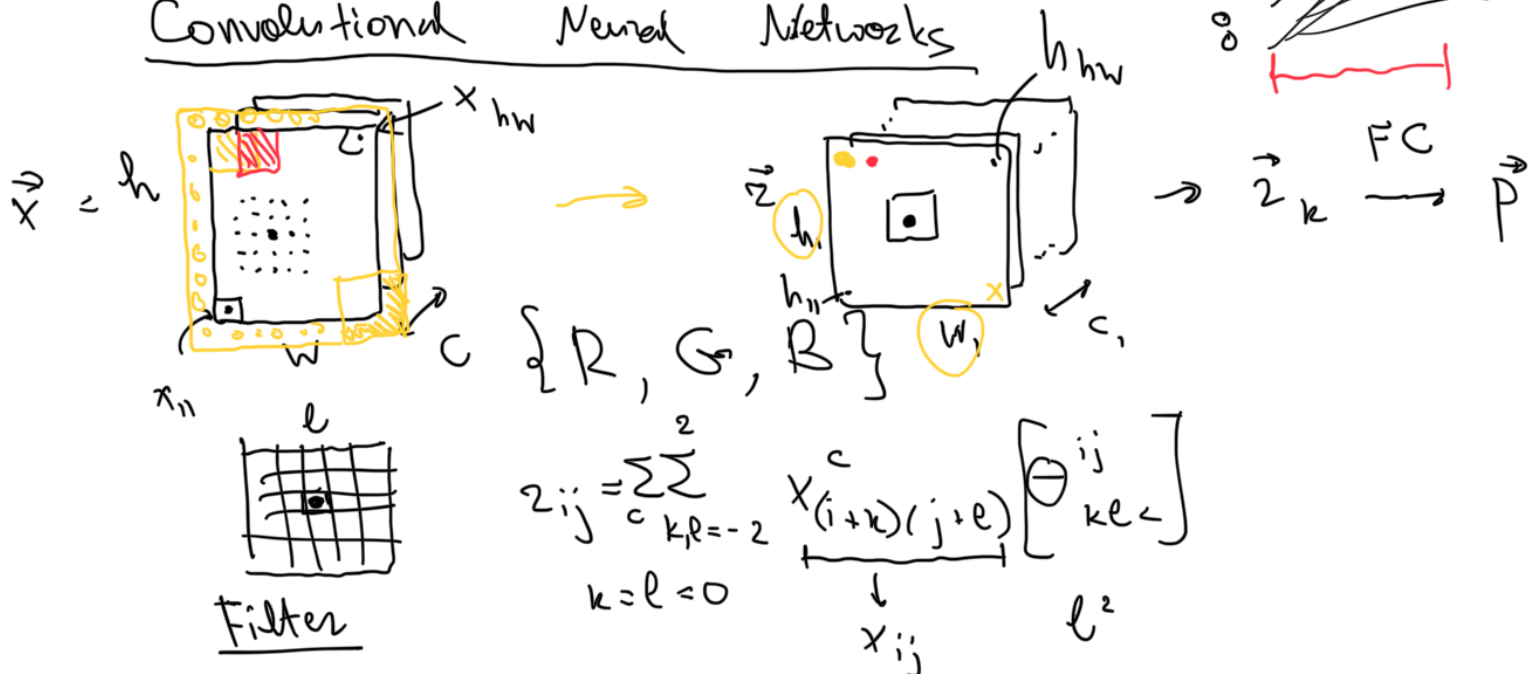$$d_1 \approx d_0 \qquad d_1 = d_0 \quad \dots \quad d_k = d_{u-1}$$

$$\boxed{P \leq m}$$

$$10\, p < m$$

ImageNet $\quad 10^7$ images

### Convolutional Neural Networks



$$x^0 = h$$

$$x_{hw}$$
$$h_{hw}$$

$$C \{R, G, B\}$$

$$\text{FC}$$
$$\to \vec{z_k} \longrightarrow \vec{p}$$

Filter

$$z_{ij} = \sum_{c}\sum_{\substack{k,\ell=-2 \\ k=\ell=0}}^{2} x_{(i+k)(j+\ell)} \left[\Theta^{ij}_{k\ell c}\right]$$
$$x_{ij} \qquad \ell^2$$

$$\underline{\ell^2 \cdot hw} \qquad \underline{\ell^2 \cdot d} < d^2$$

$$\overbrace{\ell^2}$$

$$z_{ij} = \sum x_{(i+k)(j+\ell)}^{c} \cdot \underline{\Theta_{k\ell c}}$$
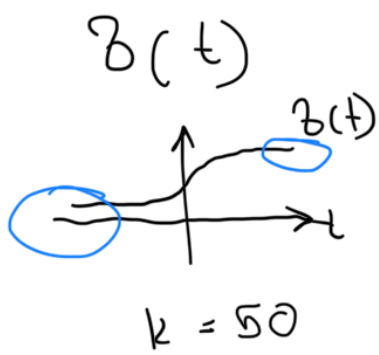
$$- k \cdot \ell^2$$

Ex. 
$$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} \quad z_{ij} = x_{ij} \qquad \begin{bmatrix} 1/9 & 1/9 & \cdots \\ \cdots & 1/9 & \\ \cdots & 1/9 & \end{bmatrix} \quad z_{ij} = \sum x_{(i+k)(j+\ell)} \cdot \frac{1}{9}$$

$z(t)$



$z(t)$

$k = 50$

$$\vec{z}_1 = \text{ReLU}(F(\vec{x}))$$

$$\frac{2014}{\text{Res Net}} \quad \frac{8-10}{50}$$

2) Pooling layers

dropout



$\frac{\partial t}{\partial \theta}$

---

$\begin{bmatrix} \text{ReLU}(t) - \text{rectified linear unit} \\ \text{ReLU}(t) = \begin{cases} t, & t \geq 0 \\ 0, & \text{otherwise} \end{cases} \end{bmatrix}$



ReLU(t)

smoothly

$$\vec{z}_1 = \text{ReLU}(F(\vec{x}) + \vec{x})$$

1) skip connection



$\vec{x} \longrightarrow F \longrightarrow \vec{z} = F(\vec{x}) + \vec{x}$

skip connection

$$\frac{\partial}{\partial \theta} g_3(g_2(g_1(\vec{x}, \theta))) =$$

$$= \frac{\partial g_3}{\partial g_2} \cdot \frac{\partial g_2}{\partial g_1} \cdot \frac{\partial g_1}{\partial \theta}$$

~ Regression
regularization



$S$

$S$ | $25$

$625$ $\longrightarrow$ FC $\longrightarrow$ FC $\longrightarrow \vec{p}$