

Optim^o - SGD → Backprop
 Theory MLP

Not^o: Data $(x_i, y_i)_{i=1}^n$ #pts
 $\begin{matrix} \nearrow & \nwarrow \\ \text{Features} & \text{Label} \\ \in \mathbb{R}^d & \end{matrix}$

Parameters: Θ

Supervised Learning (Regress^o)

$$y_i \approx f_{\Theta}(x_i)$$

↑
!?!?

Classification : $y_i \in \{-1, +1\}$.

$y_i \approx \text{sign}(f_\theta(x_i))$

Score

-1

+1

ERM : $\min_{\theta} \frac{1}{n} \sum_i l(f_\theta(x_i), y_i)$

Unregularized

$\mathcal{E}(\theta)$

$l(y, y') = (y - y')^2$ least square

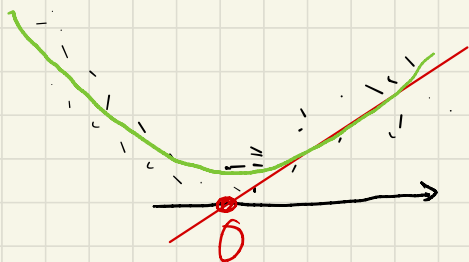
$l(y, y') = \log(1 + e^{-yy'})$ (logistic loss)

$y, y' \in \{-1, +1\}$

kernel method : $x_i \in \mathbb{R}^d \rightarrow \varphi(x_i) \in \mathbb{R}^D$

Linear method: $y_i \approx f_{\theta}(x_i) = \langle \theta, \varphi(x_i) \rangle$

$$\langle x, \theta \rangle = \sum_k x[k] \theta[k]$$



$$\varphi(x) = [1, x, x^2]$$

$$d=1$$

$$\begin{aligned} \langle \theta, \varphi(x) \rangle &= \theta[1] \cdot 1 \\ &+ \theta[2] \cdot x \\ &+ \theta[3] \cdot x^2 \end{aligned}$$

kernel Trick

$$\mathcal{E}(\theta) = \frac{1}{n} \sum_{i=1}^n (\langle \varphi(x_i), \theta \rangle - y_i)^2 + \lambda \|\theta\|^2$$

$$\Phi = \begin{bmatrix} \varphi(x_1) \\ \varphi(x_2) \\ \vdots \\ \varphi(x_n) \end{bmatrix} \in \mathbb{R}^{n \times D}$$

$$\mathcal{E}(\theta) = \frac{1}{n} \left\| \begin{array}{c} \Phi \\ \uparrow \end{array} \times \theta - y \right\|^2 + \lambda \|\theta\|^2$$

$$\min_{\Theta} \mathcal{E}(\Theta) \quad \Leftrightarrow \quad \nabla \mathcal{E}(\Theta) = 0$$

$$\nabla \mathcal{E}(\Theta) = \frac{2}{n} \Phi^T (\Phi \Theta - \underline{y}) + 2\lambda \Theta = 0$$

$$\left(\frac{\Phi^T \Phi}{n} + \lambda \text{Id} \right) \Theta = \frac{\Phi^T \underline{y}}{n}$$

$$\Theta_{\text{LS}} = \left(\frac{\Phi^T \Phi}{n} + \lambda \text{Id}_{D \times D} \right)^{-1} \Phi^T \underline{y}$$

kernel-trick / Woodbury formula

Thm $\Theta_{\text{LS}} = \Phi^T \left(\frac{\Phi \Phi^T}{n} + \lambda \text{Id}_{m \times m} \right)^{-1} \underline{y}$
kernel method

$$\hat{K} = \left(\underbrace{\langle \varphi(x_i), \varphi(x_j) \rangle}_{k(x_i, x_j)} \right)_{i, j=1}^m$$

Works if k is "valid" kernel

Def: k is an SDP kernel iff

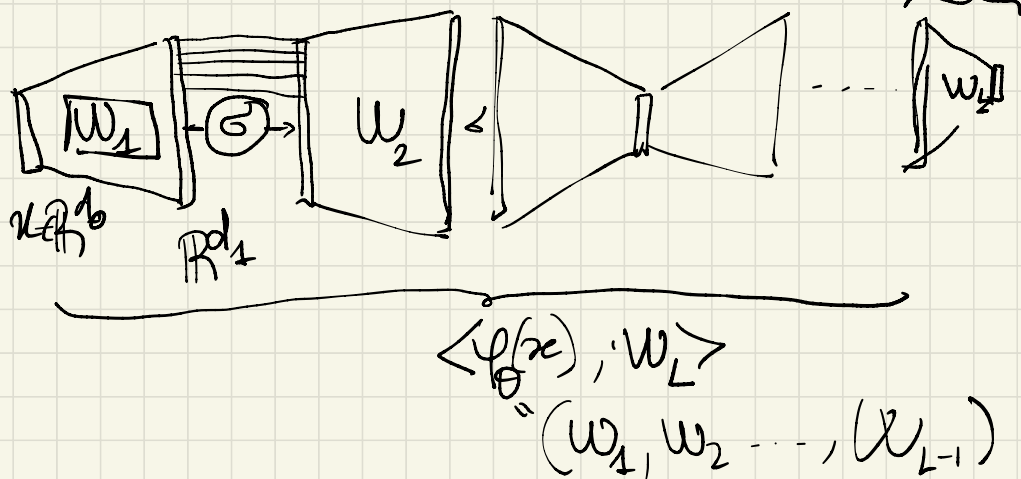
$(k(x_i, x_j))_{i,j=1}^n$ has positive eigen

$$k(x, x') = \exp\left(-\frac{\|x-x'\|^p}{\sigma}\right) \quad (p=1)$$

$$p \geq 1$$

x is graph $k(x, x')$

Neural Net: MLP

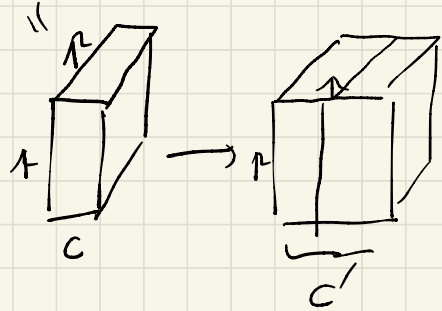
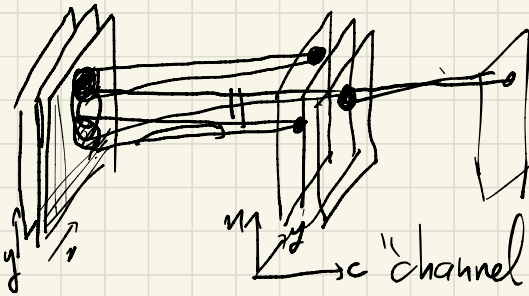


$$x_0 = x$$

$$x_{t+1} = \delta(\omega_k x_t)$$

$$f_{\theta}(x_0) = x_L$$

CNN: "Weight sharing"



Convolve: ψ "kernel" $a * \psi = \psi * a$

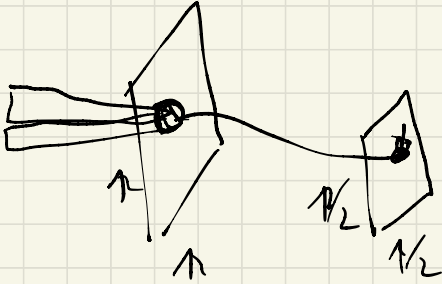
$$(a * \psi)[r] = \sum_{\delta \in \{-1, 0, 1\} \times \{-1, 0, 1\}} a[r - \delta] \psi[\delta]$$

Each CNN layer: $c \rightarrow c'$

$$\{\psi_{st}\}_{s=1 \dots c, t=1 \dots c'}$$

$$x_{t+1}[i, j, t] = \delta \left[\sum_{s=1}^c x_t[i, j, s] * \psi_{st} \right]$$

$$W_k \rightarrow \{ \varphi_{st}^k \}_{s,t}$$



CNN = conv. + sub

ResNet Block

