

# TD2: SVM, optimization and Neural Nets

## Ex1: Positive definite (p.d) kernels

### Notes:

- A kernel  $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  is said to be p.d kernel if for any  $\{x_1, \dots, x_m\} \subseteq \mathcal{X}$  the matrix  $K = [k(x_i, x_j)]_{i,j} \in \mathbb{R}^{m \times m}$  is p.d
- $K$  is a p.d. kernel if one of the following 2 equivalent conditions holds:
  - (i) the eigenvalues of  $K$  are non-negative;
  - (ii)  $\forall a \in \mathbb{R}^m, a^T K a = \sum_{i,j} a_i a_j k(x_i, x_j) \geq 0$

prop: if  $k_i$  p.d. kernel and  $d_i \geq 0 \forall i \in \mathbb{N}$

(i)  $d_1 k_1 + d_2 k_2$  and  $k_1 k_2$  p.d kernels

(ii)  $\lim_n k_n$  is p.d. kernel if the limit exists

(iii)  $k_1(x, y) = k_2(\phi(x), \phi(y))$  p.d. kernel

(power series)  
 $\Rightarrow \sum_{n \geq 0} d_n k^n$  p.d. kernel

### proof:

(i) Trivial for  $d_1 k_1 + d_2 k_2$ . For the other, it suffices to show that 2 Hermitian matrices  $A$  and  $B$  are p.d, so is their component-wise product

(ii) Just notice that the non-negativity in the definition holds also for the limit

(iii) Trivial

prop: (i)  $k = \exp(k_1)$  p.d. kernel

(ii)  $k(x, y) = e^{-\|x-y\|^2 / 2\sigma^2}$  p.d kernel

### proof:

(i)  $\exp(k_1) = \sum_{n \geq 0} \frac{1}{n!} k_1^n$  p.d. k

(ii) see course

Thm: (Aronszajn, 1950)  $k$  p.d kernel iff  $\exists \mathcal{H}$  Hilbert space and mapping  $\phi: \mathcal{X} \rightarrow \mathcal{H}$  s.t.

$$\forall x, x' \in \mathcal{X} \quad k(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathcal{H}}$$

$$\underline{1)} \quad \forall x, y \in \mathbb{R}, \begin{cases} K_1(x, y) = 10^{xy} \\ K_2(x, y) = 10^{x+y} \end{cases}$$

$K_1(x, y) = 10^{xy} = e^{xy \ln(10)}$  with  $\mathcal{X} = \mathbb{R}$  is a exponential of a scaled (by factor  $\ln(10) > 0$ ) linear kernel

$\Rightarrow$  p.d. kernel

•  $x = (x_1, \dots, x_n) \in \mathbb{R}^n$  and  $a = (a_1, \dots, a_n) \in \mathbb{R}^n$ ,

$$\begin{aligned} a^T K_2 a &= \sum_{i=1}^n \sum_{j=1}^n a_i a_j K_2(x_i, x_j) = \sum_{i=1}^n \sum_{j=1}^n a_i a_j 10^{x_i + x_j} \\ &= \sum_{i=1}^n a_i 10^{x_i} \cdot \sum_{j=1}^n a_j 10^{x_j} = \left( \sum_{i=1}^n a_i 10^{x_i} \right)^2 \geq 0 \end{aligned}$$

$\Rightarrow$  p.d. kernel

$$\underline{2)} \quad x, y \in [0, 1[ , K_3(x, y) = -\log(1 - xy) = \sum_{n \geq 1} \frac{(xy)^n}{n} \Rightarrow K_3 \text{ p.d. kernel}$$

$\underline{3)}$  an example before :  $\tilde{K}(x, y) = \min(x, y)$  with  $\mathcal{X} = \mathbb{R}_+$

Let us define  $\phi: \mathbb{R} \rightarrow L^2(\mathbb{R}_+)$

$$x \mapsto \mathbb{1}_{[0, x]}$$

$$\forall a, b \in \mathbb{R}^+$$

$$\int_0^\infty \phi(a) \phi(b) = \int_0^a \mathbb{1} = a$$

$$[0, a] \cap [0, b] = [0, \min(a, b)]$$

$$\phi(\min(a, b)) = \phi(a) \phi(b)$$

$$\left. \begin{array}{l} \int_0^\infty \phi(a) \phi(b) = a \\ [0, a] \cap [0, b] = [0, \min(a, b)] \\ \phi(\min(a, b)) = \phi(a) \phi(b) \end{array} \right\} \Rightarrow \min(a, b) = \int_0^{+\infty} \mathbb{1}_{[0, a]} \cdot \mathbb{1}_{[0, b]} = \langle \phi(a), \phi(b) \rangle_{L^2(\mathbb{R}_+)}$$

Thanks to Aronszajn (1950) Thm,  $\tilde{K}$  p.d kernel

exercise:  $\max(x, y)$  with  $\mathcal{X} = \mathbb{R}_+$  p.d kernel ?

Back to the question:

$$K_4(x, y) = \min \{ f(x)g(y), f(y)g(x) \} = \min \left\{ \frac{f(x)}{g(x)}, \frac{f(y)}{g(y)} \right\} g(x)g(y)$$

By using  $\phi_4(x) = \mathbb{1}_{[0, \frac{f(x)}{g(x)}]}$ ,

$$\min \left\{ \frac{f(x)}{g(x)}, \frac{f(y)}{g(y)} \right\} = \int \mathbb{1}_{[0, \frac{f(x)}{g(x)}]} \mathbb{1}_{[0, \frac{f(y)}{g(y)}]} = \langle \phi_4(x), \phi_4(y) \rangle_{L^2(\mathbb{R}_+)}$$

$\min \left\{ \frac{f(x)}{g(x)}, \frac{f(y)}{g(y)} \right\}$  is a p.d kernel }  $\Rightarrow K_4$  p.d kernel  
Moreover,  $g(x)g(y)$  is a p.d kernel (trivial)

4) a) Let  $S$  be a sequence of words on a finite alphabet  $A$

For example  $A = \{A, B, C, \dots, Z\}$

$x = \underline{CGGSLI} \underline{AMMWF} \underline{GV}$

Subchains of length 5 in  $x = (\bullet, \bullet, \dots, \bullet)$

Let  $\phi_u(x)$  denote the number of subchains  $u$  in  $x$

Let us consider  $k(x, y) = \#$  subchains of length 5 that  $x$  and  $y \in S$  share

$\Rightarrow k(x, y) = \sum_{u \in A^5} \phi_u(x) \phi_u(y)$  which is clearly a p.d kernel:

$$\begin{aligned} \sum_{i, j=1}^m a_i a_j k(x_i, x_j) &= \sum_{i, j=1}^m a_i a_j \sum_{u \in A^5} \phi_u(x_i) \phi_u(x_j) \\ &= \sum_u \left( \sum_i a_i \phi_u(x_i) \right)^2 \geq 0 \end{aligned}$$

b) Left as an exercise

## Ex 2: Computations of RKHS

Definition of reproducing kernel (r.k) and RKHS (Reproducing Kernel Hilbert Space)

$K: X \times X \rightarrow \mathbb{R}$  is a r.k of a Hilbert space  $\mathcal{H}_0 \subset \mathbb{R}^X$  if

(i)  $\mathcal{H}_0$  contains all functions of the form

$$\forall x \in X, K_x : t \mapsto K(x, t)$$

(ii)  $\forall x \in X, h \in \mathcal{H}_0$  the reproducing property holds:

$$h(x) = \langle h, K_x \rangle_{\mathcal{H}_0}$$

If a r.k exists, then  $\mathcal{H}_0$  is called RKHS

In particular,  $\forall x, y \in X, K(x, y) = \langle K_x, K_y \rangle_{\mathcal{H}_0}$

Thm (see course):

For any set  $X$ , a function  $K: X \times X$  is a p.d iff it is a r.k

$$\Downarrow \begin{cases} K_5(x, y) = (xy+1)^2 & \text{which are clearly p.d. kernels.} \\ K_6(x, y) = (xy-1)^2 \\ K_7(x, y) = K_5(x, y) + K_6(x, y) \\ \quad = (xy+1)^2 + (xy-1)^2 \\ \quad = 2((xy)^2 + 1) \end{cases}$$

These are special cases of the **polynomial kernel** of degree  $p$ :

$$K_{\text{poly}}(x, y) = (\langle x, y \rangle_{\mathbb{R}^d} + c)^p$$

Let us find its RKHS  $\mathcal{H}$  for  $p=2$  and  $c=0$

(Left as an exercise = - what is the RKHS of the general  $K_{\text{poly}}$ ?  
- what about  $K_5, K_6$  and  $K_7$ )

• **step 1**: Look for an inner-product

$$\begin{aligned} K_{\text{poly}}(x, y) &= x^T y x^T y \\ &= \text{Trace}(x^T y x^T y) \\ &= \text{Trace}(x x^T y y^T) \\ &= \langle x x^T, y y^T \rangle_F \end{aligned}$$

$$\left( \begin{array}{l} \text{rks:} \\ \text{Tr}(AB) = \text{Tr}(BA) \\ \text{Tr}(A^t) = \text{Tr}(A) \end{array} \right)$$

where  $F$  is the Frobenius norm for matrices  $\mathbb{R}^{d \times d}$

• **step 2**: propose a candidate RKHS

We know that  $\mathcal{H}$  contains all the functions

$$f(x) = \sum_i a_i K(x_i, x) = \sum_i a_i \langle x_i x_i^T, x x^T \rangle_F = \langle \sum_i a_i x_i x_i^T, x x^T \rangle_F$$

Any symmetric matrix in  $\mathbb{R}^{d \times d}$  may be decomposed as  $\sum_i a_i x_i x_i^T$

$\Rightarrow$  Candidate RKHS  $\mathcal{H}$ : the set of quadratic functions

$$f_S(x) = \langle S, xx^T \rangle_F = x^T S x \quad \text{for } S \in S^{d \times d}$$

where  $S^{d \times d}$  is the set of symmetric matrices in  $\mathbb{R}^{d \times d}$ , endowed with the inner product  $\langle f_{S_1}, f_{S_2} \rangle_{\mathcal{H}} = \langle S_1, S_2 \rangle_F$

- **step 3:** check that the candidate is a Hilbert Space  
Here it is trivial because  $\mathcal{H}$  is an Euclidean space
- **step 4:** check that  $\mathcal{H}$  is the RKHS

$$(i) \quad \{ k_x : t \mapsto k(x, t) = \langle xx^T, tt^T \rangle_F \} \subset \mathcal{H}$$

(ii) Let us consider  $f_S \in \mathcal{H}$  and  $x \in \mathcal{X}$ . Then,

$$f_S(x) = \langle S, xx^T \rangle_F = \langle f_S, f_{xx^T} \rangle_{\mathcal{H}} = \langle f_S, k_x \rangle_{\mathcal{H}}$$

2] see ex 1 for the p.d kernels property.

## Ex4: Neural Nets

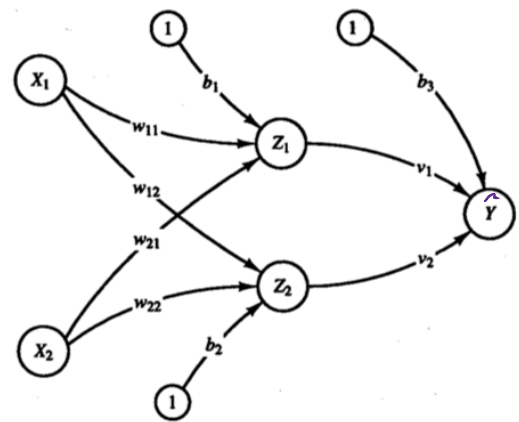
Algorithm for this architecture:

step 0: initialize weights

step 1: while "stopping criterion is false"

do step 2-6

② For each training step:



feed forward / forward propagation

③ compute output of each hidden unit: *activation function*

$$z_{-in_1} = b_1 + x_1 w_{11} + x_2 w_{21} \Rightarrow z_1 = \sigma(z_{-in_1})$$

$$z_{-in_2} = b_2 + x_1 w_{12} + x_2 w_{22} \Rightarrow z_2 = \sigma(z_{-in_2})$$

④ compute output of net

$$y_{-in} = b_3 + z_1 v_1 + z_2 v_2 \Rightarrow \hat{y} = \sigma(y_{-in})$$

⑤ compute error and weights / biases correction terms:

$$\delta = (y - \hat{y}) \sigma'(y_{-in})$$

(error on net output)

$$\Delta v_j = \alpha \cdot \delta \cdot z_j \quad \text{and} \quad \Delta b_3 = \alpha \cdot \delta$$

(weights correction term)

$$\delta_j = \delta_{-in_j} \cdot \sigma'(z_{-in_j}) \quad \text{with} \quad \delta_{-in_j} = \delta \cdot v_j$$

(error on each hidden unit)

$$\Delta w_{ij} = \alpha \cdot \delta_j \cdot x_i \quad \text{and} \quad \Delta b_j = \alpha \cdot \delta_j$$

(weights correction term)

⑥ updates weights and biases:

$$v_j^{new} \leftarrow v_j + \Delta v_j \quad \text{and} \quad b_3^{new} \leftarrow b_3 + \Delta b_3$$

$$w_{ij}^{new} \leftarrow w_{ij} + \Delta w_{ij} \quad \text{and} \quad b_j^{new} \leftarrow b_j + \Delta b_j$$

(gradient descent)

tips:

Let us define a loss function  $l = \frac{1}{2} (\hat{y} - y)^2$ ,

$$\frac{\partial l}{\partial \hat{y}} = (\hat{y} - y) \quad ; \quad \frac{\partial l}{\partial v_j} = \frac{\partial l}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial y_{-in}} \cdot \frac{\partial y_{-in}}{\partial v_j} = \underbrace{(\hat{y} - y) \sigma'(y_{-in})}_{-\delta} z_j$$

1)

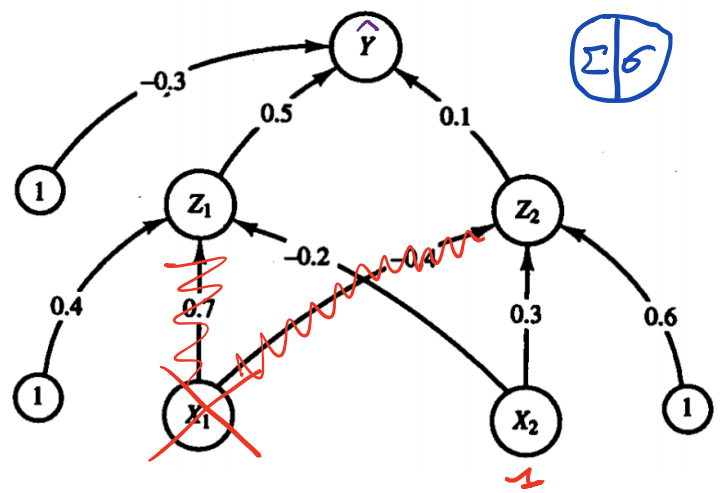
$(x_1, x_2) = (0, 1)$

$Y = 1$

stepsize :  $\alpha = 0.25$

activation function :  $\sigma(x) = \frac{1}{1 + e^{-x}}$

new weights ?



$z_{-in_1} = 0.4 + (-0.2) = 0.2 \Rightarrow z_1 = \frac{1}{1 + e^{-0.2}} \approx 0.55$

$z_{-in_2} = 0.6 + 0.3 = 0.9 \Rightarrow z_2 = \frac{1}{1 + e^{-0.9}} \approx 0.71$

$y_{-in} = -0.3 + 0.5 \cdot z_1 + 0.1 \cdot z_2 \approx 0.05 \Rightarrow \hat{y} \approx \frac{1}{1 + e^{-0.05}} \approx 0.51$

$\delta \approx (1 - 0.51) \sigma'(0.05) \approx 0.12$  }  $\left( \begin{aligned} \sigma'(x) &= \sigma(x)(1 - \sigma(x)) \\ &= \frac{1}{1 + e^{-x}} \left( 1 - \frac{1}{1 + e^{-x}} \right) \end{aligned} \right)$

$\Delta w_1 \approx 0.25 \times 0.12 \times 0.55 \approx 0.018$

$\Delta w_2 \approx 0.25 \times 0.12 \times 0.71 \approx 0.02$

$\Delta b_3 \approx 0.25 \times 0.12 \approx 0.03$

$\delta_1 = \delta \cdot w_1 \cdot \sigma'(z_{-in_1}) \approx 0.12 \times 0.5 \times \sigma'(0.2) \approx 0.02$

$\delta_2 = \delta \cdot w_2 \cdot \sigma'(z_{-in_2}) \approx 0.12 \times 0.1 \times \sigma'(0.9) \approx 0.003$

$\Delta w_{11} = \Delta w_{12} = 0 \quad \Delta w_{21} \approx 0.05$

$\Delta w_{22} \approx 0.00075$

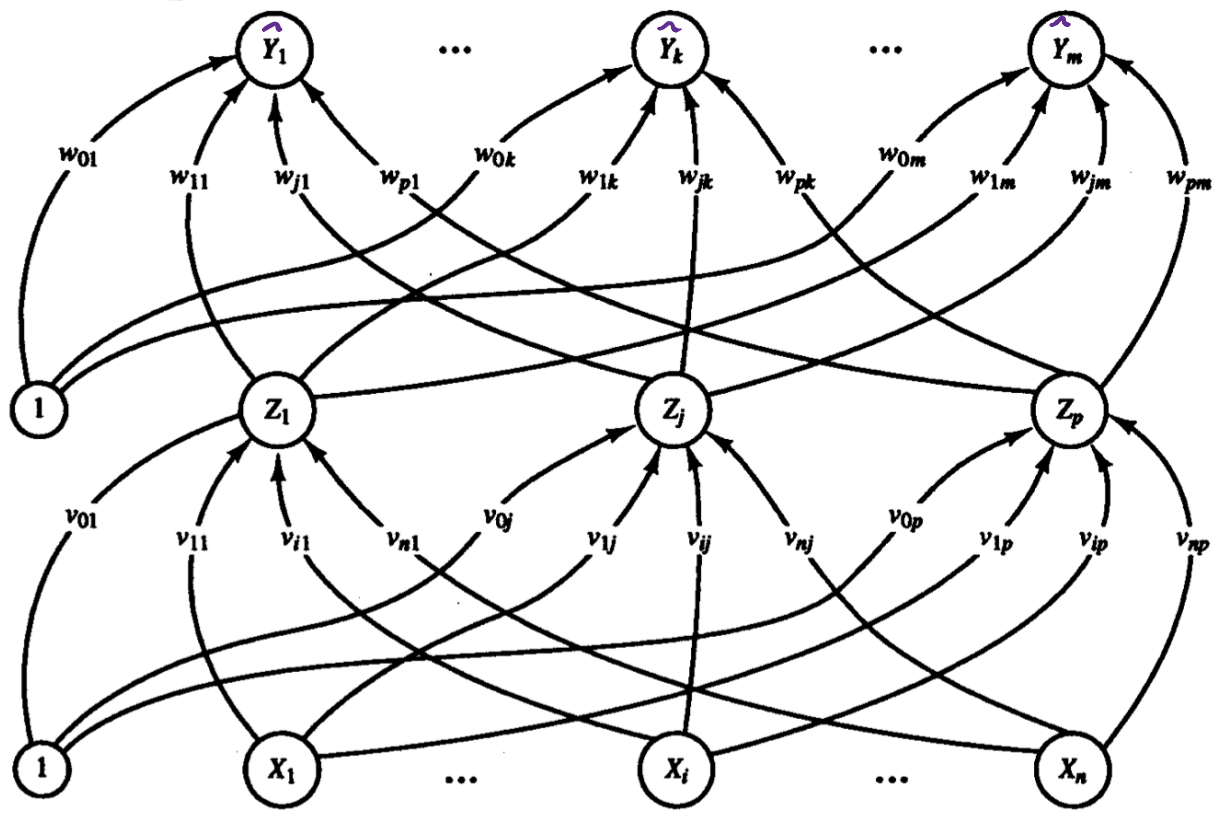
$w_1 \approx 0.5 + 0.018 \approx 0.517$      $w_2 \approx 0.1 + 0.02 \approx 0.12$

$b_3 \approx -0.3 + 0.03 \approx -0.27$     ... until stopping criterion

Exercise : implement with python feedforward and backprop  
answering question 1, 2 and 3

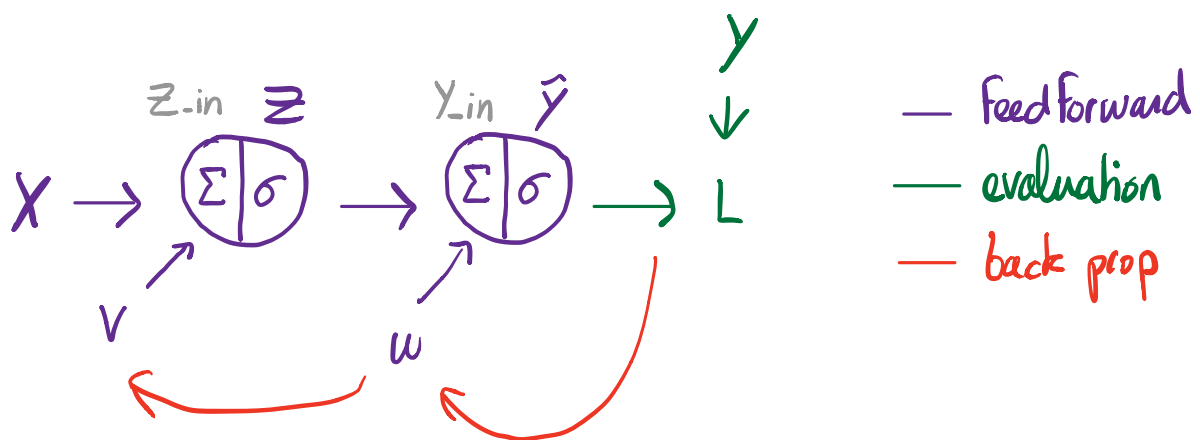


# An overview



Graphs taken from "Fundamentals of Neural Networks" by L. Faussett

Let us simplify the graph :



- update its weights by gradient descent : 
$$\begin{cases} W_{ij} \leftarrow W_{ij} - \alpha \frac{\partial L}{\partial W_{ij}} \\ v_{ij} \leftarrow v_{ij} - \alpha \frac{\partial L}{\partial v_{ij}} \end{cases}$$

- compute  $\begin{cases} \frac{\partial L}{\partial W_{ij}} \\ \frac{\partial L}{\partial v_{ij}} \end{cases}$  by chain-rule