

**FACULTÉ DES SCIENCES  
DÉPARTEMENT D'INFORMATIQUE**

**IFT712 - Techniques d'apprentissage  
TP3**

*préparé par*

**Agathe Le Boulter (leba3207)**

**Philippe Spino (spip2401)**

**Gabriel Gibeau Sanchez (gibg2501)**

*présenté à*

**Pierre-Marc Jodoin**

13 mars 2020

# 1 Question 1

Dans me cas d'une régression, avec l'utilisation d'une optimisation de type Maximisation A Priori, on dispose de la fonction de coût  $E_D$  à minimiser telle que :

$$E_D(\vec{w}) = \frac{\beta}{2} \sum_{n=1}^N (y(\vec{x}_n, w) - t_n)^2 + \frac{\alpha}{2} \vec{w}^T \vec{w}$$

Si l'on écrit :  $\lambda = \frac{\alpha}{\beta}$  et  $y(\vec{x}_n, w) = \vec{w}^T \vec{\phi}(\vec{x}_n)$

$$\iff \frac{1}{2} \sum_{n=1}^N \left( \vec{w}^T \vec{\phi}(\vec{x}_n) - t_n \right)^2 + \frac{\lambda}{2} \vec{w}^T \vec{w}$$

Cette expression est la représentation primale que l'on note  $J(\vec{w})$ .

$$\begin{aligned} J(\vec{w}) &= \frac{1}{2} \sum_{n=1}^N \left( \vec{w}^T \vec{\phi}(\vec{x}_n) - t_n \right)^2 + \frac{\lambda}{2} \vec{w}^T \vec{w} \\ &= \frac{1}{2} \sum_{n=1}^N \left( \vec{w}^T \vec{\phi}(\vec{x}_n)^2 \vec{w} - 2 \vec{w}^T \vec{\phi}(\vec{x}_n) t_n + t_n^2 \right) + \frac{\lambda}{2} \vec{w}^T \vec{w} \end{aligned} \quad (1)$$

Afin de déterminer sa représentation duale dans un second temps, on va déterminer  $\vec{w}$  par le calcul du gradient de  $E_D$  par rapport à  $\vec{w}$  et le forcer à 0.

$$\begin{aligned} \nabla_{\vec{w}} E_D(\vec{w}) &= \frac{\partial}{\partial \vec{w}} \left( \frac{1}{2} \sum_{n=1}^N \left( \vec{w}^T \vec{\phi}(\vec{x}_n)^2 \vec{w} - 2 \vec{w}^T \vec{\phi}(\vec{x}_n) t_n + t_n^2 \right) + \frac{\lambda}{2} \vec{w}^T \vec{w} \right) = 0 \\ &\iff \frac{1}{2} \sum_{n=1}^N \left( 2 \vec{w}^T \vec{\phi}(\vec{x}_n)^2 - 2 \vec{\phi}(\vec{x}_n) t_n \right) + \lambda \vec{w} = 0 \\ &\iff \sum_{n=1}^N \left( \vec{w}^T \vec{\phi}(\vec{x}_n)^2 - \vec{\phi}(\vec{x}_n) t_n \right) + \lambda \vec{w} = 0 \\ &\iff \sum_{n=1}^N \left( \vec{w}^T \vec{\phi}(\vec{x}_n)^2 - \vec{\phi}(\vec{x}_n) t_n \right) = -\lambda \vec{w} \\ &\iff \vec{w} = -\frac{1}{\lambda} \sum_{n=1}^N \left( \vec{w}^T \vec{\phi}(\vec{x}_n) - t_n \right) \vec{\phi}(\vec{x}_n) \\ &\iff \vec{w} = \sum_{n=1}^N \left( -\frac{\vec{w}^T \vec{\phi}(\vec{x}_n) - t_n}{\lambda} \vec{\phi}(\vec{x}_n) \right) \end{aligned} \quad (2)$$

On peut simplifier cette expression par :

$$\Longleftrightarrow \vec{w} = \sum_{n=1}^N a_n \vec{\phi}(\vec{x}_n) \quad (3)$$

en considérant  $a_n$  tel que :  $a_n = -\frac{\vec{w}^T \vec{\phi}(\vec{x}_n) - t_n}{\lambda}$

Aussi, on peut exprimer  $\sum_{n=1}^N \vec{\phi}(\vec{x}_n) = \Phi^T$ , matrice de design où la  $n^{ieme}$  ligne est donnée par  $\vec{\phi}(\vec{x}_n)^T$ . Et  $\vec{a}$  tel que  $\vec{a} = (a_1, \dots, a_N)^T$ . On écrit alors :  $\vec{w} = \Phi^T \vec{a}$ .

Maintenant, revenons à notre représentation primale (sous forme développée) de  $J(\vec{w})$  et exprimons  $\vec{w}$  tel que  $\vec{w} = \Phi^T \vec{a}$ .  $J$  se devient être exprimé en fonction de  $\vec{a}$ . On inscrit :  $\Phi^T = \sum_{n=1}^N \vec{\phi}(\vec{x}_n)$  et  $\vec{t} = \sum_{n=1}^N t_n$ .

$$J(\vec{a}) = \frac{1}{2} [(\Phi \vec{a}^T)(\Phi^T \Phi)(\Phi^T \vec{a}) - 2(\Phi \vec{a}^T)\Phi^T \vec{t} + \vec{t}^T \vec{t}] + \frac{\lambda}{2} (\Phi \vec{a}^T)(\Phi^T \vec{a}) \quad (4)$$

que l'on peut réécrire tel que :

$$\Longleftrightarrow \frac{1}{2} \vec{a}^T \Phi \Phi^T \Phi \vec{a} - \vec{a}^T \Phi \Phi^T \vec{t} + \frac{1}{2} \vec{t}^T \vec{t} + \lambda \vec{a}^T \Phi \Phi^T \vec{a} \quad (5)$$

qui est la représentation duale de  $J(\vec{w})$ .

On définit la matrice de Gram  $K$  telle qu'une matrice symétrique  $N \times N$  où chaque élément est une fonction noyau telle que :  $K_{nm} = \vec{\phi}(\vec{x}_n)^T \vec{\phi}(\vec{x}_m) = k(\vec{x}_n, \vec{x}_m)$ .

On exprime donc  $K$  tel que :  $K = \Phi \Phi^T$

On peut alors réécrire la représentation duale de  $J(\vec{w})$  telle que :

$$J(\vec{a}) = \frac{1}{2} \vec{a}^T K K \vec{a} - \vec{a}^T K \vec{t} + \frac{1}{2} \vec{t}^T \vec{t} + \frac{\lambda}{2} \vec{a}^T K \vec{a} \quad (6)$$

On détermine le gradient de  $J(\vec{a})$  par rapport à  $\vec{a}$  et on le force à 0 pour déterminer le meilleur vecteur tel que les données sont les mieux classifiées.

$$\begin{aligned}
& \nabla_{\vec{a}} J(\vec{a}) = 0 \\
\iff & \frac{\partial}{\partial \vec{a}} \left( \frac{1}{2} \vec{a}^T K K \vec{a} - \vec{a}^T K \vec{t} + \frac{1}{2} \vec{t}^T \vec{t} + \frac{\lambda}{2} \vec{a}^T K \vec{a} \right) = 0 \\
& \iff \vec{a} K K - K \vec{t} + \lambda \vec{a} K = 0 \\
& \iff \vec{a} (K K + \lambda K) = K \vec{t} \\
& \iff \vec{a} = \frac{K \vec{t}}{K K + \lambda K} \\
& \iff \vec{a} = \frac{K \vec{t}}{K(K + \lambda I_N)} \\
& \iff \vec{a} = (K + \lambda I_N)^{-1} \vec{t} \quad (7)
\end{aligned}$$

Reprenons maintenant notre modèle de régression linéaire tel que :  $y(\vec{x}) = \vec{w}^T \vec{\phi}(\vec{x})$ , où nous pouvons exprimer  $\vec{w}^T$  tel que :  $\vec{w}^T = \Phi \vec{a}^T$ .

$$y(\vec{x}) = \vec{w}^T \vec{\phi}(\vec{x}) = \Phi \vec{a}^T \vec{\phi}(\vec{x}) \quad (8)$$

Nous pouvons remplacer  $\vec{a}^T$  par l'expression trouvée via le calcul du gradient.

$$\iff y(\vec{x}) = \Phi \left( (K + \lambda I_N)^{-1} \vec{t} \right)^T \vec{\phi}(\vec{x}) \quad (9)$$

On sait que :  $\Phi = \sum_{n=1}^N \vec{\phi}(\vec{x}_n)$

$$\iff y(\vec{x}) = \left( (K + \lambda I_N)^{-1} \vec{t} \right)^T \sum_{n=1}^N \vec{\phi}(\vec{x}_n)^T \vec{\phi}(\vec{x}) \quad (10)$$

Or on sait que la fonction noyau est définie telle que :  $k(\vec{x}_n, \vec{x}_m) = \vec{\phi}(\vec{x}_n)^T \vec{\phi}(\vec{x}_m)$

$$\iff y(\vec{x}) = \left( (K + \lambda I_N)^{-1} \vec{t} \right)^T k(\vec{x}_n, \vec{x}) \quad (11)$$

On définit maintenant  $k(\vec{x})$  comme la forme vectorisée des éléments de fonctions noyaux  $k_n(\vec{x}) = k(\vec{x}_n, \vec{x})$

$$\begin{aligned}
& \iff y(\vec{x}) = \left( (K + \lambda I_N)^{-1} \vec{t} \right)^T k(\vec{x}) \\
& \iff y(\vec{x}) = k(\vec{x})^T (K + \lambda I_N)^{-1} \vec{t} \quad (12)
\end{aligned}$$

## 2 Question 2

Un vecteur de support est un exemple de l'ensemble des données d'entraînement qui est situé sur la marge d'un modèle Machine à Vecteur de Support. On définit la marge comme la plus petite distance entre la surface de séparation du modèle entre les classes et l'ensemble des données d'entraînement.

Pour un modèle, on exprime ses prédictions telles que :  $y_{\vec{w}}(\vec{\phi}(\vec{x}_n)) = \vec{w}^T \vec{\phi}(\vec{x}_n)$ , où  $\vec{\phi}(\vec{x}_n)$  peut se révéler être un vecteur de support.

On se place dans un cas où les données sont non linéairement séparables et on envisage une utilisation d'une représentation primal. Aussi, on accepte qu'un nombre restreint de données soient en-deçà de la marge afin d'éviter au maximum le sur-apprentissage. On dispose alors des contraintes telles que :

$$\begin{cases} t_n \left( \vec{w}^T \vec{\phi}(\vec{x}_n) + w_0 \right) = 1 \\ C \sum_{n=1}^N \xi_n = 1 \end{cases} \quad (13)$$

$\xi_n$  représente la variable de ressort permettant d'accepter des données mal classées pour notre modèle et C est un hyper-paramètre.

On a alors :

$$\operatorname{argmin}_{\vec{w}, w_0} \left( \frac{1}{2} \|\vec{w}\|^2 + C \sum_{n=1}^N \xi_n \right) \quad (14)$$

tel que :  $t_n \left( y_{\vec{w}}(\vec{\phi}(\vec{x}_n)) \right) \geq 1 - \xi_n$  et  $\forall n, \xi_n \geq 0$

La contrainte  $t_n \left( y_{\vec{w}}(\vec{\phi}(\vec{x}_n)) \right) \geq 1 - \xi_n$  peut être réécrite telle que :

$$\xi_n \geq 1 - t_n \left( y_{\vec{w}}(\vec{\phi}(\vec{x}_n)) \right)$$

Si on remplace  $\xi_n$  dans l'expression de la représentation primale (équation 7.21 dans Bishop)

$$C \sum_{n=1}^N \xi_n + \frac{1}{2} \|\vec{w}\|^2 \quad (15)$$

on a alors :

$$\iff \operatorname{argmin}_{\vec{w}, w_0} \left( \frac{1}{2} \|\vec{w}\|^2 \right) + C \sum_{n=1}^N 1 - t_n \left( y_{\vec{w}}(\vec{\phi}(\vec{x}_n)) \right) \quad (16)$$

considérant la contrainte  $1 - t_n \left( y_{\vec{w}}(\vec{\phi}(\vec{x}_n)) \right) \geq 0$

$$\begin{aligned}
&\Longleftrightarrow \operatorname{argmin}_{\vec{w}, w_0} \left( \frac{1}{2} \|\vec{w}\|^2 \right) + C \sum_{n=1}^N \max \left( 0, 1 - t_n \left( y_{\vec{w}}(\vec{\phi}(\vec{x}_n)) \right) \right) \\
&\Longleftrightarrow \operatorname{argmin}_{\vec{w}, w_0} \sum_{n=1}^N \max \left( 0, 1 - t_n \left( y_{\vec{w}}(\vec{\phi}(\vec{x}_n)) \right) \right) + \lambda \|\vec{w}\|^2
\end{aligned} \tag{17}$$

avec  $\lambda = \frac{C}{2}$

On reconnaît ici l'expression de la fonction d'erreur de Hinge.