

Master Masef - Mido 2018-2019
Exam : Machine Learning in Finance¹ : Duration 1h30

Exercise 1 (QCM) : [10pts]

1. If $(X_i, Y_i)_{i \in \llbracket 1, n \rrbracket}$ is a learning sample, how is the calibration error generally defined for a classification problem ?
 - a) $\frac{1}{n} \sum_{i=1}^{i=n} 1_{f(X_i) \neq Y_i}$
 - b) $\frac{1}{n} \sum_{i=1}^{i=n} |f(X_i) - Y_i|$
 - c) $E[|f(X_{n+1}) - Y_{n+1}|]$
 - d) none of the answers above

2. If $(X_i, Y_i)_{i \in \llbracket 1, n \rrbracket}$ is a learning sample, which of the hypothesis always made on the variables ?
 - a) the X_i are all Gaussian
 - b) the (X_i, Y_i) have all the same laws
 - c) the X_i have all the same laws but not necessarily the Y_i
 - d) the (X_i, Y_i) are all independent
 - e) none of the answers above

3. If $R_n(f_n)$ is the calibration error for the optimal classifier f_n obtained for the learning sample $(X_i, Y_i)_{i \in \llbracket 1, n \rrbracket}$ and if we note $R(f_n) = E[1_{f(X_{n+1}) \neq Y_{n+1}}]$ the prediction error, which properties are always verified ?
 - a) $R_n(f_n) = R(f_n)$
 - b) $R_n(f_n) \leq R(f_n)$
 - c) $E(R_n(f_n)) \leq R(f_n)$
 - d) none of the answers above

4. If there are k particular points in \mathbb{R}^d which can be classified in all possible ways by the family of classifiers \mathcal{F}_1 and if it is not possible to classify them in all possible ways by the family of classifiers \mathcal{F}_2 then which assertions is/are necessarily true :
 - a) $VC(\mathcal{F}_1) = k$
 - b) $VC(\mathcal{F}_2) < k$
 - c) $VC(\mathcal{F}_1) > VC(\mathcal{F}_2)$
 - d) none of the answers above

5. If for two families of classifiers/machines \mathcal{F}_1 and \mathcal{F}_2 the error of calibration is the same on a learning sample, which machine does the Vapnik Chernovenkis inequality and SRM principle encourage to use to predict :

1. Pierre Brugière University Paris 9 Dauphine

- a) \mathcal{F}_1 if $VC(\mathcal{F}_1) > VC(\mathcal{F}_2)$
 b) not \mathcal{F}_2 if $VC(\mathcal{F}_1) > VC(\mathcal{F}_2)$
 c) none of the answers above
6. If $x_1, x_2 \cdots x_n$ form a family of independent vectors, and if 0 is the null vector, in how many different ways is it possible to classify : $0, x_1, x_2 \cdots x_n$?
 a) 2^n
 b) $n - 1$
 c) 2^{n+1}
 d) none of the answers above
7. The inequality of Vapnik Chervonenkis enables, knowing the complexity of the family of classifiers/machine used and the error on calibration to :
 a) define a confidence interval for $R(f_n)$
 b) calculate the exact value of $R(f_n)$
 c) define a boundary $\epsilon < 1$ for $R(f_n)$
 d) none of the answers above
8. Which assertions is/are true ?
 a) a high VC for a family of classifiers implies necessarily a bad quality of prediction
 b) an infinite VC for a family of classifiers implies always a perfect calibration
 c) the VC for a parametric family of classifiers is always close to the number of parameters of the family of classifiers
 d) none of the answers above
9. Which assertions are true
 a) It is unlikely that a machine which calibrates badly will predict accurately
 b) VC gives no guarantee that a complex machine which calibrates well will predict accurately
 c) SRM means Structural Risk Maximization
 d) a machine which is very complex in a VC sense predicts always very accurately
10. In \mathbb{R}^d for $w \neq 0$ let $H_{w,b} = \{x \in \mathbb{R}^d, \langle w, x \rangle + b = 0\}$. Which assertions are true :
 a) $\forall x \in \mathbb{R}^d d(x, H_{w,b}) = \frac{|\langle w, x \rangle + b|}{\|w\|^2}$
 b) $\forall x \in \mathbb{R}^d d(x, H_{w,b}) = \frac{|\langle w, x \rangle + b|}{\|w\|}$
 c) $\forall x \in \mathbb{R}^d d(x, H_{w,b}) = \frac{|\langle w, x \rangle - b|}{\|w\|}$
 d) none of the answers above

11. Let $w \in \mathbb{R}^d \setminus \{0\}$. Let $H_{w,b} = \{(x, y) \in \mathbb{R}^{d+1}, \langle w, x \rangle + b - y = 0\}$. Which assertions are true :
- $d(H_{w,b_1}, H_{w,b_2}) = \frac{|b_2+b_1|}{\|w\|_d}$
 - $d(H_{w,b_1}, H_{w,b_2}) = \frac{|b_2-b_1|}{\sqrt{1+\|w\|_d^2}}$
 - $d(H_{w,b_1}, H_{w,b_2}) = \frac{|b_2-b_1|}{\|w\|_d}$
 - none of the answers above
12. If you can classify perfectly a learning sample with the hyperplanes $H_{w_1,b_1} = \{x \in \mathbb{R}^d, \langle w_1, x \rangle + b_1 = 0\}$ of margin Δ_1 and $H_{w_2,b_2} = \{x \in \mathbb{R}^d, \langle w_2, x \rangle + b_2 = 0\}$ of margin Δ_2 you have a good reason to choose H_{w_1,b_1} if :
- $\Delta_1 < \Delta_2$
 - $\Delta_1 > \Delta_2$
 - $b_1 > b_2$
 - $b_2 > b_1$
13. If you can classify perfectly a learning sample with two classes with an hyperplane H of margin Δ and if we call \mathcal{C}_1 and \mathcal{C}_2 the convex envelopes of the two classes of points then :
- for sure $d(\mathcal{C}_1, \mathcal{C}_2) \geq \Delta$
 - for sure $d(\mathcal{C}_1, \mathcal{C}_2) > \Delta$
 - H of maximum margin $\iff (H \cap \mathcal{C}_1 \neq \emptyset \text{ and } H \cap \mathcal{C}_2 \neq \emptyset)$
 - none of the answers above
14. If \mathcal{F} is a family of classifiers of \mathbb{R}^{10^6} of diameters 1 of hyperplanes of margin 0.1 then :
- $VC(\mathcal{F}) = 10^6$
 - $VC(\mathcal{F}) = 10^6 + 1$
 - $VC(\mathcal{F}) \approx 100$
 - $VC(\mathcal{F}) \approx 10$
15. According to the minimax theorem for any function $g : \mathcal{Y} \times \mathcal{Z} \rightarrow \mathbb{R}$
- $\inf_{y \in \mathcal{Y}} \left[\sup_{z \in \mathcal{Z}} g(y, z) \right] \leq \sup_{z \in \mathcal{Z}} \left[\inf_{y \in \mathcal{Y}} g(y, z) \right]$
 - $\sup_{z \in \mathcal{Z}} \left[\inf_{y \in \mathcal{Y}} g(y, z) \right] \leq \inf_{y \in \mathcal{Y}} \left[\sup_{z \in \mathcal{Z}} g(y, z) \right]$
 - $\sup_{z \in \mathcal{Z}} \left[\inf_{y \in \mathcal{Y}} g(y, z) \right] = \inf_{y \in \mathcal{Y}} \left[\sup_{z \in \mathcal{Z}} g(y, z) \right]$
16. When solving a SVM for a classification $\{-1, 1\}$ which assertions are true :
- The Primal and Dual problems have the same solution if and only if the KKT conditions are added to the constraints of the dual problem
 - The Primal and Dual problems have the same solution because of the

particular nature of the problem

c) The KKT conditions are automatically satisfied because of the particular nature of the problem

d) none of the answers above

17. When solving a C-SVM for a classification $\{-1, 1\}$ which assertions are true :

a) all the support vectors are necessarily classified correctly

b) some support vectors may be classified incorrectly

c) the support vectors are necessarily on the border of the maximum margin hyperplane

d) if $0 < \alpha_i < C$, x_i is a support vector on the margin of the maximum margin hyperplane

18. Among the following functions from $\mathbb{R}^d \times \mathbb{R}^d$ to \mathbb{R} which ones are Kernel

a) $\exp(-\frac{\|x-y\|^2}{2}) + \exp(-\frac{\|x-y\|^2}{4})$ (because sum of Kernels implies strictly positive as well)

b) $\exp(-\|x\| - \|y\|)$ (criteria of strict positiveness easy to verify)

c) $\langle x, y \rangle$ (because of the form $\langle \phi(x), \phi(y) \rangle$)

19. If $\phi(\cdot)$ is the transformation associated to the kernel $K(x, y) = \exp(-\frac{\|x-y\|^2}{2})$ by the relationship $\langle \phi(x), \phi(y) \rangle = K(x, y)$ among these properties which ones are true for the image points $\phi(x_i)$ from a learning sample.

a) the $\phi(x_i)$ are necessarily on a sphere of radius 1

b) the $\phi(x_i)$ for sure can be separated from 0 by an hyperplane

c) the $\phi(x_i)$ for sure are independent if the x_i are distinct

20. If $K(\cdot, \cdot)$ is a kernel and $(H, \langle \cdot, \cdot \rangle_{RK})$ is a Reproducing Hilbert Space for K then which of the following assertions are true :

a) for all $f \in H$, $\langle K(x, \cdot), f \rangle_{RK} = f(x)$

b) $K(x_1, x_2)^2 \leq K(x_1, x_1)K(x_2, x_2)$

c) for all $f \in H$, $f(x)^2 \leq K(x, x)\langle f, f \rangle_{RK}$

Exercise : [4pts]

We consider the following Ridge Regression problem :

$$(P) \arg \min_{\beta \in \mathbb{R}^d} \|Y - X\beta\|_n^2 + \lambda \|\beta\|_d^2$$

where Y is a vector of \mathbb{R}^n , X is a $n \times d$ matrix and $\|\cdot\|$ is the euclidean norm.

We remark that (P) can be written in the form :

$$(Q) \begin{cases} \arg \min_{z \in \mathbb{R}^n, \beta \in \mathbb{R}^d} \|z\|_n^2 + \lambda \|\beta\|_d^2 \\ z = Y - X\beta \end{cases}$$

1. [0.5pt] Determine the expression of the Lagrangian $L(z, \beta, \gamma)$ of (Q) with $\gamma \in \mathbb{R}^n$.

Correction : $L(z, \beta, \gamma) = \|z\|^2 + \lambda \|\beta\|^2 + \gamma'(Y - X\beta - z)$

2. [0.5pt] Calculate $\frac{\partial L}{\partial z}$ and $\frac{\partial L}{\partial \beta}$
Correction : $\frac{\partial L}{\partial z} = 2z' - \gamma'$ and $\frac{\partial L}{\partial \beta} = 2\lambda\beta' - \gamma'X$
3. [1pt] Write the dual formulation (D) of (Q)
Correction : $\frac{\partial L}{\partial z} = 0$ and $\frac{\partial L}{\partial \beta} = 0$ implies $z(\gamma) = \frac{1}{2}\gamma$ and $\beta(\gamma) = \frac{1}{2\lambda}X'\gamma$
so, $\min_{z,r} L(z, r, \gamma) = L(z(\gamma), r(\gamma), \gamma)$
 $= \frac{1}{4}\|\gamma\|^2 + \frac{1}{4\lambda}\|X'\gamma\|^2 + \gamma'Y - \frac{1}{2\lambda}\gamma'XX'\gamma - \gamma'\frac{1}{2}\gamma$
 $= -\frac{1}{4}\|\gamma\|^2 - \frac{1}{4\lambda}\|X'\gamma\|^2 + \gamma'Y$ that we note $L^*(\gamma)$.
and the dual problem to solve is : $\max_{\gamma} L^*(\gamma)$
4. [1pt] Solve (D) by determining its solution γ^* and calculate as well the corresponding β^* it produces.
Correction : to solve (D) we solve $\frac{\partial L^*}{\partial \gamma} = 0$.
 $\frac{\partial L^*}{\partial \gamma} = -\frac{1}{2}\gamma' - \frac{1}{2\lambda}\gamma'XX' + Y'$
so $\frac{\partial L^*}{\partial \gamma} = 0 \implies \frac{1}{2\lambda}XX'\gamma + \frac{1}{2}\gamma = Y \implies XX'\gamma + \lambda\gamma = 2\lambda Y$
 $\implies (XX' + \lambda Id_n)\gamma = 2\lambda Y \implies \gamma = 2\lambda(XX' + \lambda Id_n)^{-1}Y$.
Now using the fact that $\beta(\gamma) = \frac{1}{2\lambda}X'\gamma$ we get $\beta^* = X'(XX' + \lambda Id_n)^{-1}Y$.
5. [1pt] Is β^* the expression you expected knowing that the solution of (P) is known to be $(X'X + \lambda Id_d)^{-1}X'Y$? Comment.
Correction : For this type of problem the primal and dual problem have the same solution. So, even if the expressions are not the same they are in fact equal as $(X'X + \lambda Id_d)^{-1}X' = X'(XX' + \lambda Id_n)^{-1}$ which can be proved easily by multiplying both matrices by $XX' + \lambda Id_n$.

Exercise : [8.5pts]

Let $\{(x_i, y_i)\}_{i \in [1, n]}$ be a learning sample where $x_i \in \mathbb{R}^d$ and $y_i \in \mathbb{R}$. Let $\epsilon \in \mathbb{R}^{+*}$. We consider the problem :

$$(P_\epsilon) \left\{ \begin{array}{l} \arg \min_{w \in \mathbb{R}^d, b \in \mathbb{R}, \xi \in \mathbb{R}^n, \tilde{\xi} \in \mathbb{R}^n} \frac{1}{2}\|w\|^2 + C \sum_{i=1}^n (\xi_i + \tilde{\xi}_i) \\ y_i - \langle w, x_i \rangle - b \leq \epsilon + \xi_i \\ \langle w, x_i \rangle + b - y_i \leq \epsilon + \tilde{\xi}_i \\ \xi_i \geq 0 \\ \tilde{\xi}_i \geq 0 \end{array} \right.$$

1. [1pt] Write the Lagrangian $L(w, b, \xi, \tilde{\xi}, \alpha, \tilde{\alpha}, \mu, \tilde{\mu})$ of the problem (P_ϵ)

Correction :

$$L(w, b, \xi, \tilde{\xi}, \alpha, \tilde{\alpha}, \mu, \tilde{\mu}) = \frac{1}{2}\|w\|^2 + C \sum_{i=1}^n (\xi_i + \tilde{\xi}_i) \\ + \sum_{i=1}^{i=n} \alpha_i [y_i - \langle w, x_i \rangle - b - \epsilon - \xi_i] + \sum_{i=1}^{i=n} \tilde{\alpha}_i [\langle w, x_i \rangle + b - y_i - \epsilon - \tilde{\xi}_i] \\ - \sum_{i=1}^{i=n} \xi_i \mu_i - \sum_{i=1}^{i=n} \tilde{\xi}_i \tilde{\mu}_i.$$

$$\begin{aligned}
&= \frac{1}{2}\|w\|^2 + \langle w, \sum_{i=1}^{i=n} (\tilde{\alpha}_i x_i - \alpha_i x_i) \rangle \\
&+ \sum_{i=1}^{i=n} \left[b(-\alpha_i + \tilde{\alpha}_i) + \xi_i(C - \alpha_i - \mu_i) + \tilde{\xi}_i(C - \tilde{\alpha}_i - \tilde{\mu}_i) \right] \\
&+ \sum_{i=1}^{i=n} [\alpha_i(y_i - \epsilon) - \tilde{\alpha}_i(y_i + \epsilon)].
\end{aligned}$$

2. [1pt] Calculate :

- (a) $\frac{\partial L}{\partial w}$
- (b) $\frac{\partial L}{\partial b}$
- (c) $\frac{\partial L}{\partial \xi_i}$
- (d) $\frac{\partial L}{\partial \tilde{\xi}_i}$

Correction :

$$\frac{\partial L}{\partial w} = w' + \sum_{i=1}^{i=n} (\tilde{\alpha}_i - \alpha_i) x_i'$$

$$\frac{\partial L}{\partial b} = \sum_{i=1}^{i=n} (\tilde{\alpha}_i - \alpha_i)$$

$$\frac{\partial L}{\partial \xi_i} = C - \alpha_i - \mu_i$$

$$\frac{\partial L}{\partial \tilde{\xi}_i} = C - \tilde{\alpha}_i - \tilde{\mu}_i$$

3. [1pt] Show that the dual (D_ϵ) of (P_ϵ) is of the form

$$(D_\epsilon) \left\{ \begin{array}{l} \arg \max_{\alpha \in \mathbb{R}^n, \tilde{\alpha} \in \mathbb{R}^n} -\frac{1}{2}(\alpha - \tilde{\alpha})' K(\alpha - \tilde{\alpha}) + \sum_{i=1}^{i=n} [\alpha_i(y_i - \epsilon) - \tilde{\alpha}_i(y_i + \epsilon)] \\ 0 \leq \alpha_i \leq C \\ 0 \leq \tilde{\alpha}_i \leq C \end{array} \right.$$

- (a) [0.25pt] what is the expression of K ?
- (b) [0.25pt] say, without demonstration, under what conditions on the $\{(x_i)\}_{i \in [1, n]}$ K is invertible.
- (c) [0.50pt] express the solution w^* of (P_ϵ) as a function of the solutions α^* and $\tilde{\alpha}^*$ of (D_ϵ) and of the x_i .

Correction :

Taking into account the nullity of the partial derivatives, the Lagrangian expression is reduced to :

$$\frac{1}{2}\|w\|^2 + \langle w, \sum_{i=1}^{i=n} (\tilde{\alpha}_i x_i - \alpha_i x_i) \rangle + \sum_{i=1}^{i=n} [\alpha_i(y_i - \epsilon) - \tilde{\alpha}_i(y_i + \epsilon)].$$

replacing w by $\sum_{i=1}^{i=n} (\alpha_i - \tilde{\alpha}_i) x_i$ the expression becomes :

$$-\frac{1}{2}(\alpha - \tilde{\alpha})' K(\alpha - \tilde{\alpha}) + \sum_{i=1}^{i=n} [\alpha_i(y_i - \epsilon) - \tilde{\alpha}_i(y_i + \epsilon)]$$

where K is the symmetric matrix of general term $\langle x_i, x_j \rangle$.

In the constraints :

$(C - \alpha_i - \mu_i = 0$ with $\alpha_i \geq 0$ and $\mu_i \leq 0) \iff 0 \leq \alpha_i \leq C$ and
 $(C - \tilde{\alpha}_i - \tilde{\mu}_i = 0$ with $\tilde{\alpha}_i \geq 0$ and $\tilde{\mu}_i \leq 0) \iff 0 \leq \tilde{\alpha}_i \leq C$ So,

$$(D_\epsilon) \begin{cases} \arg \max_{\alpha \in \mathbb{R}^n, \tilde{\alpha} \in \mathbb{R}^n} -\frac{1}{2}(\alpha - \tilde{\alpha})'K(\alpha - \tilde{\alpha}) + \sum_{i=1}^{i=n} [\alpha_i(y_i - \epsilon) - \tilde{\alpha}_i(y_i + \epsilon)] \\ 0 \leq \alpha_i \leq C \\ 0 \leq \tilde{\alpha}_i \leq C \end{cases}$$

- (a) K is the symmetric matrix of general term $\langle x_i, x_j \rangle$
 - (b) K is invertible if and only if the $\{(x_i)\}_{i \in [1, n]}$ form an independent family of \mathbb{R}^d
 - (c) the condition $\frac{\partial L}{\partial w} = 0$ implies $w^* = \sum_{i=1}^{i=n} (\alpha_i^* - \tilde{\alpha}_i^*)x_i'$
4. Due to the form of the problem, we assume in all what follows that the KKT conditions are satisfied.
- (a) **[0.50pt]** write the KKT conditions
 - (b) **[0.50pt]** show that $\alpha_i \tilde{\alpha}_i = 0$
 - (c) **[0.50pt]** show that $\xi_i \tilde{\xi}_i = 0$

Correction :

- (a) the KKT conditions are

$$(D_\epsilon) \begin{cases} \alpha_i (y_i - \langle w, x_i \rangle - b - \epsilon - \xi_i) = 0 \\ \tilde{\alpha}_i (y_i - \langle w, x_i \rangle - b + \epsilon + \tilde{\xi}_i) = 0 \\ (C - \alpha_i) \xi_i = 0 \\ (C - \tilde{\alpha}_i) \tilde{\xi}_i = 0 \end{cases}$$

- (b) by the absurd : $\alpha_i \tilde{\alpha}_i \neq 0 \implies$

$$\begin{cases} y_i - \langle w, x_i \rangle - b - \epsilon - \xi_i = 0 \\ y_i - \langle w, x_i \rangle - b + \epsilon + \tilde{\xi}_i = 0 \end{cases}$$

$\implies -\epsilon - \xi_i = \epsilon + \tilde{\xi}_i$ which is impossible

- (c) by the absurd : $\xi_i \tilde{\xi}_i \neq 0 \implies \alpha_i = C$ and $\tilde{\alpha}_i = C$
 $\implies \alpha_i \tilde{\alpha}_i \neq 0$ which is impossible
- (d) according to the KKT conditions

5. We call "support vectors" the x_i which appears with a non zero coefficient in the expression of w^* .
- (a) **[0.5pt]** show that, x_i is a support vector if and only if :
 $\alpha_i \neq 0$ or $\tilde{\alpha}_i \neq 0$
 - (b) **[0.5pt]** show that, $0 < \alpha_i < C$ or $0 < \tilde{\alpha}_i < C$ implies :
 $|\langle w, x_i \rangle + b - y_i| = \epsilon$

- (c) **[0.5pt]** show that, $\alpha_i = C$ or $\tilde{\alpha}_i = C$ implies $|\langle w, x_i \rangle + b - y_i| \geq \epsilon$
- (d) **[0.5pt]** conclude on the location of the points (x_i, y_i) for the support vectors x_i .

Correction :

- (a) as $\alpha_i \tilde{\alpha}_i = 0$ then $(\alpha_i - \tilde{\alpha}_i) \neq 0 \iff \alpha_i \neq 0$ or $\tilde{\alpha}_i \neq 0$
- (b) $0 < \alpha_i < C \implies C \xi_i = 0 \implies \xi_i = 0$.
 now, $\xi_i = 0$ and $\alpha_i \neq 0$ imply that :
 $y_i - \langle w, x_i \rangle - b - \epsilon = 0$ and $|\langle w, x_i \rangle + b - y_i| = \epsilon$.
 Same reasoning for $0 < \tilde{\alpha}_i < C$.
- (c) Let assume $\alpha_i = C$ then $y_i - \langle w, x_i \rangle - b - \epsilon - \xi_i = 0$ implies
 $|\langle w, x_i \rangle + b - y_i| = \epsilon + \xi_i \geq \epsilon$.
 Same reasoning for the case $\tilde{\alpha}_i = C$.
- (d) according to the previous questions for the support vectors the points (x_i, y_i) are on the border or outside the domain :
 $\{(x, y), |\langle w, x \rangle + b - y| = \epsilon\}$
6. **[1pt]** What is the aim of the problem (P_ϵ) ? is there a geometric interpretation in \mathbb{R}^{d+1} ? and what is the interest to have $\|w^*\|$ small?

Correction : It is a regression problem where we search for a margin hyperplane $|\langle w, x \rangle - y + b| \leq \epsilon$ in \mathbb{R}^{d+1} which contains the (x_i, y_i) . The thickness of the hyperplane is $\frac{2\epsilon}{\sqrt{1+\|w\|^2}} \leq 2\epsilon$. The minimizing of $\|w\|$ has a negative impact on the slimness but has for objective to reduce the complexity of the model and to give a solution to the problem which is unique.