

PC3. ECOLE POLYTECHNIQUE. MAP 569. MACHINE LEARNING II.

EXERCISE 1 (ELASTIC-NET) Let $Y \in \mathbb{R}^n$ and $\mathbf{X} = [\mathbf{X}_1, \dots, \mathbf{X}_p] \in \mathbb{R}^{n \times p}$. The Elastic-Net estimator involves both a ℓ^2 and a ℓ^1 penalty. It is meant to improve the Lasso estimator when the columns of \mathbf{X} are "strongly" correlated. It is defined for $\lambda, \mu \geq 0$ by

$$\hat{\beta}_{\lambda, \mu} \in \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \mathcal{L}(\beta) \quad \text{with} \quad \mathcal{L}(\beta) = \|Y - \mathbf{X}\beta\|^2 + \lambda \|\beta\|^2 + \mu |\beta|_{\ell^1}.$$

In the following, we assume that the columns of \mathbf{X} have norm 1.

1. Check that the partial derivative of \mathcal{L} with respect to $\beta_j \neq 0$ is given by

$$\partial_j \mathcal{L}(\beta) = 2 \left((1 + \lambda) \beta_j - R_j + \frac{\mu}{2} \operatorname{sign}(\beta_j) \right) \quad \text{with} \quad R_j = \mathbf{X}_j^\top \left(Y - \sum_{k: k \neq j} \beta_k \mathbf{X}_k \right).$$

2. Prove that the minimum of $\beta_j \rightarrow \mathcal{L}(\beta_1, \dots, \beta_j, \dots, \beta_p)$ is reached at $\beta^* = (\beta_1^*, \dots, \beta_p^*)$ where

$$\beta_j^* = \frac{R_j^*}{1 + \lambda} \left(1 - \frac{\mu}{2|R_j^*|} \right)_+$$

$$\text{and } R_j^* = \mathbf{X}_j^\top \left(Y - \sum_{k: k \neq j} \beta_k^* \mathbf{X}_k \right)$$

3. Propose an algorithm to compute the Elastic-Net estimator.

The Elastic-Net procedure is implemented in the R package `glmnet` available at <http://cran.r-project.org/web/packages/glmnet/>.

EXERCISE 2 (SUPPORT VECTOR MACHINE (SVM)) Minimization of convex functions: Karush-Kuhn-Tucker sufficient conditions

Let $f, -g_1, \dots, -g_n$ be \mathcal{C}^1 convex functions and define the Lagrangian

$$\mathcal{L} : (x, \lambda) \mapsto f(x) - \sum_{i=1}^n \lambda_i g_i(x).$$

For any (x, λ) , the Karush-Kuhn-Tucker conditions read:

1. $\forall i \in \llbracket 1, n \rrbracket : g_i(x) \geq 0;$
2. $\forall i \in \llbracket 1, n \rrbracket : \lambda_i \geq 0;$
3. $\nabla_x \mathcal{L}(x, \lambda) = 0;$
4. $\min(\lambda_i, g_i(x)) = 0$ for $i = 1, \dots, n.$

We know that, under the previous assumptions, KKT conditions are sufficient: if a couple $(\hat{x}, \hat{\lambda})$ fulfills the KKT conditions, then

$$\hat{x} \in \underset{\forall i \in \llbracket 1, n \rrbracket : g_i(x) \geq 0}{\operatorname{argmin}} f(x) \quad \text{and} \quad \hat{\lambda} \in \underset{\lambda \geq 0}{\operatorname{argmax}} \inf_x \mathcal{L}(x, \lambda).$$

Also, still under the previous assumptions, weak duality holds:

$$\sup_{\lambda \geq 0} \inf_x \mathcal{L}(x, \lambda) \leq \inf_x \sup_{\lambda \geq 0} \mathcal{L}(x, \lambda) = \inf_{\forall i \in \llbracket 1, n \rrbracket : g_i(x) \geq 0} f(x).$$

Strong duality (i.e. equality holds) under additional assumptions.

Strong duality: If there exists a x such that $g_i(x) > 0$ for all $i \in \{1, \dots, n\}$, then the KKT conditions are also necessary (i.e. $\hat{\lambda}$ exists and KKT conditions are satisfied by $(\hat{x}, \hat{\lambda})$) and

$$\sup_{\lambda \geq 0} \inf_x \mathcal{L}(x, \lambda) = \inf_x \sup_{\lambda \geq 0} \mathcal{L}(x, \lambda).$$

Application to SVM

For any $w \in \mathbb{R}^p$, define the linear function $f_w(x) = \langle w, x \rangle$ from \mathbb{R}^p to \mathbb{R} . For a given $R > 0$, we consider the set of linear functions $\mathcal{F} = \{f_w : \|w\| \leq R\}$. The aim of this exercise is to investigate the classifier $\hat{h}_{\varphi, \mathcal{F}}(x) = \text{sign}(\hat{f}_{\varphi, \mathcal{F}}(x))$ where $\hat{f}_{\varphi, \mathcal{F}}$ is solution to the convex optimisation problem

$$\hat{f}_{\varphi, \mathcal{F}} \in \operatorname{argmin}_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \varphi(-y_i f(x_i)),$$

with $\varphi(x) = (1 + x)_+$ the *hinge loss*.

1. From the strong duality, prove that there exists $\lambda \geq 0$ such that

$$\hat{f}_{\varphi, \mathcal{F}} \in \operatorname{argmin}_{f_w} \left\{ \frac{1}{n} \sum_{i=1}^n (1 - y_i f_w(x_i))_+ + \lambda \|w\|^2 \right\}.$$

2. Prove that $\hat{f}_{\varphi, \mathcal{F}} = f_{\hat{w}}$ where \hat{w} belongs to $V = \text{Span}\{x_i : i = 1, \dots, n\}$.
3. Prove that $\hat{w} = \sum_{j=1}^n \hat{\beta}_j x_j$ where $\hat{\beta} = [\hat{\beta}_1, \dots, \hat{\beta}_n]^\top$ is solution to

$$\hat{\beta} = \operatorname{argmin}_{\beta \in \mathbb{R}^n} \left\{ \frac{1}{n} \sum_{i=1}^n (1 - y_i (K\beta)_i)_+ + \lambda \beta^\top K \beta \right\},$$

with K the Gram matrix $K = [\langle x_i, x_j \rangle]_{1 \leq i, j \leq n}$.

4. Check that this minimization problem is equivalent to

$$\hat{\beta} = \operatorname{argmin}_{\substack{\beta, \xi \in \mathbb{R}^n \text{ such that} \\ y_i (K\beta)_i \geq 1 - \xi_i \\ \xi_i \geq 0}} \left\{ \frac{1}{n} \sum_{i=1}^n \xi_i + \lambda \beta^\top K \beta \right\}.$$

5. Let us assume that K is not singular. From the KKT conditions, check that $\hat{\beta}_i = y_i \hat{\alpha}_i / (2\lambda)$, for $i = 1, \dots, n$ with $\hat{\alpha}_i$ fulfilling $\min(\hat{\alpha}_i, y_i (K\hat{\beta})_i - (1 - \hat{\xi}_i)) = 0$ et $\min(1/n - \hat{\alpha}_i, \hat{\xi}_i) = 0$.

6. Prove the following properties

- if $y_i \hat{f}_{\varphi, \mathcal{F}}(x_i) > 1$ then $\hat{\beta}_i = 0$;
- if $y_i \hat{f}_{\varphi, \mathcal{F}}(x_i) < 1$ then $\hat{\beta}_i = y_i / (2\lambda n)$;
- in any case (in particular if $y_i \hat{f}_{\varphi, \mathcal{F}}(x_i) = 1$), $0 \leq \hat{\beta}_i y_i \leq 1 / (2\lambda n)$.

7. From the strong duality, prove that $\hat{\alpha}_i$ is solution to the dual problem

$$\hat{\alpha} = \operatorname{argmax}_{0 \leq \alpha_i \leq 1/n} \left\{ \sum_{i=1}^n \alpha_i - \frac{1}{4\lambda} \sum_{i,j=1}^n K_{i,j} y_i y_j \alpha_i \alpha_j \right\}.$$