

# Apprentissage robuste de distance par géométrie riemannienne

Antoine COLLAS<sup>1</sup>, Arnaud BRELOY<sup>2</sup>, Guillaume GINOLHAC<sup>3</sup>, Chengfang REN<sup>1</sup>, Jean-Philippe OVARLEZ<sup>1,4</sup>,

<sup>1</sup>CentraleSupélec SONDRRA, Université Paris-Saclay, 3 rue Joliot Curie, 91190 Gif-sur-Yvette, France,

<sup>2</sup>LEME, Université Paris-Nanterre, 50 Rue de Sèvres, 92410 Ville-d'Avray, France,

<sup>3</sup>LISTIC, Université Savoie Mont-Blanc, 5 Chemin de Bellevue, Annecy-Le-Vieux, 74940 Annecy, France,

<sup>4</sup>DEMR, ONERA, Université Paris-Saclay, 8 Chemin de la Hunière, 91123 Palaiseau, France,

antoine.collas@centralesupelec.fr

**Résumé** – L'apprentissage de distance propose d'optimiser la distance de Mahalanobis pour des problèmes de classification. Plusieurs algorithmes associés à ce problème s'interprètent comme des estimateurs d'une matrice de covariance. Partant de ce constat, nous proposons un nouveau problème d'estimation appelé *Robust Geometric Metric Learning (RGML)* qui vise à estimer une matrice de covariance par classe et leur barycentre riemannien. Deux algorithmes associés à ce problème d'estimation sont développés : *RGML gaussien* et *RGML Tyler*. Ces algorithmes optimisent deux fonctions de coût sur la variété riemannienne des matrices symétriques définies positives et sa sous-variété à déterminant unitaire. Enfin, les algorithmes proposés sont appliqués sur des jeux de données réels. De bonnes performances ainsi qu'une robustesse aux données mal étiquetées sont obtenues.

**Abstract** – *Metric learning* proposes to optimize the *Mahalanobis* distance for classification problems. Several algorithms associated with this problem are interpreted as estimators of a covariance matrix. Based on this observation, we propose a new estimation problem called *Robust Geometric Metric Learning (RGML)* which aims to estimate a covariance matrix per class and their Riemannian barycenter. Two algorithms associated with this estimation problem are developed: *RGML Gaussian* and *RGML Tyler*. These algorithms optimize two cost functions on the Riemannian manifold of symmetric positive definite matrices and its submanifold of unitary determinant. Finally, the proposed algorithms are applied on real datasets. Good performance as well as robustness to mislabeled data are obtained.

## 1 Introduction

De nombreux algorithmes de classification et de partitionnement reposent sur une distance entre les données. Ces algorithmes incluent les classiques *K-means*, *Classifieur du centroïde le plus proche*, et *k-plus proches voisins*. La définition de cette distance a une importance cruciale puisqu'elle détermine la similarité entre points. En pratique, ces algorithmes reposent le plus souvent sur la distance euclidienne qui est  $d_{\mathbf{I}_p}(\mathbf{x}_i, \mathbf{x}_j) = \|\mathbf{x}_i - \mathbf{x}_j\|_2$  pour  $\mathbf{x}_i, \mathbf{x}_j \in \mathbb{R}^p$ . Cependant, cette distance est sujette à plusieurs problèmes. Un exemple est lorsque deux classes ont une variance élevée le long d'un axe commun : dans cette configuration, deux données de la même classe peuvent être éloignées l'une de l'autre, tandis que deux données de deux classes différentes peuvent être très proches.

Afin de trouver une distance plus pertinente pour la classification, le problème de l'*apprentissage de distance* a été proposé. Il vise à trouver une distance de Mahalanobis

$$d_{\mathbf{A}}(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{(\mathbf{x}_i - \mathbf{x}_j)^T \mathbf{A}^{-1} (\mathbf{x}_i - \mathbf{x}_j)} \quad (1)$$

qui rapproche les données de même classe et éloigne les données de classes différentes. L'*apprentissage de distance* consiste à formaliser un problème d'optimisation par rapport à  $\mathbf{A} \in \mathcal{S}_p^+$ , l'ensemble des matrices symétriques définies posi-

tives (SDP) d'ordre  $p$ , pour déterminer au mieux cette distance.

Soit un problème de classification à  $K$  classes :  $m$  données  $\{\mathbf{x}_1, \dots, \mathbf{x}_m\} \subset \mathbb{R}^p$  étiquetées dans  $[[1, K]]$ . Les données de la classe  $k$  sont notées  $\{\mathbf{x}_{kl}\}$  et  $n_k$  paires distinctes  $(\mathbf{x}_{kl}, \mathbf{x}_{kq})$  avec  $l \neq q$  sont formées. L'ensemble  $S_k$  contient tous ces couples et  $S$  contient les  $n_S = \sum_{k=1}^K n_k$  couples de toutes les classes. Le rapport  $\frac{n_k}{n_S}$  est noté  $\pi_k$ . Ensuite, pour chaque classe  $k$ ,  $n_k$  vecteurs  $\mathbf{s}_{ki}$  sont formés par soustraction des paires  $(\mathbf{x}_{kl}, \mathbf{x}_{kq}) \in S_k$ , i.e.  $\mathbf{s}_{ki} = \mathbf{x}_{kl} - \mathbf{x}_{kq}$ . L'ensemble  $\mathcal{D}$  contient  $n_{\mathcal{D}}$  paires de vecteurs qui n'appartiennent pas à la même classe. Chaque vecteur  $\mathbf{d}_i$  est défini comme la soustraction des éléments de chaque paire dans  $\mathcal{D}$ ,  $\mathbf{d}_i = \mathbf{x}_l - \mathbf{x}_q$  avec  $(\mathbf{x}_l, \mathbf{x}_q) \in \mathcal{D}$ .  $\mathcal{S}_p$  est l'ensemble des matrices symétriques d'ordre  $p$ , et  $\mathcal{SS}_p^+$  est le sous-ensemble de  $\mathcal{S}_p^+$  des matrices de déterminant unitaire. Enfin,  $|\cdot| : \mathcal{S}_p \rightarrow \mathbb{R}$  est l'opérateur déterminant.

### 1.1 État de l'art et estimation de la matrice de covariance

De nombreux travaux d'*apprentissage de distance* existent dans la littérature (voir [1] pour une étude complète). Dans ce qui suit, nous en présentons deux importants.

*ITML (Information-Theoretic Metric Learning)* [2] propose

de trouver une matrice  $\mathbf{A}$  qui reste proche d'une matrice prédéfinie  $\mathbf{A}_0$  tout en respectant des contraintes de similarités et de dissimilarités. La *divergence de Kullback-Leibler gaussienne*  $D_{KL}(\mathbf{A}_0, \mathbf{A}) = \text{Tr}(\mathbf{A}^{-1}\mathbf{A}_0) + \log |\mathbf{A}\mathbf{A}_0^{-1}|$  mesure la proximité entre  $\mathbf{A}$  et  $\mathbf{A}_0$ . Le problème d'optimisation *ITML* est

$$\begin{aligned} \underset{\mathbf{A} \in \mathcal{S}_p^+}{\text{minimiser}} \quad & \text{Tr}(\mathbf{A}^{-1}\mathbf{A}_0) + \log |\mathbf{A}| \\ & d_{\mathbf{A}}^2(\mathbf{x}_l, \mathbf{x}_q) \leq u, \quad (\mathbf{x}_l, \mathbf{x}_q) \in \mathcal{S}, \\ & d_{\mathbf{A}}^2(\mathbf{x}_l, \mathbf{x}_q) \geq l, \quad (\mathbf{x}_l, \mathbf{x}_q) \in \mathcal{D}, \end{aligned} \quad (2)$$

où  $u, l \in \mathbb{R}$  sont des seuils choisis pour rapprocher des points de même classe et éloigner des points de classes différentes. En pratique,  $\mathbf{A}_0$  est choisi comme la matrice identité ou bien la matrice de covariance empirique. Dans ce dernier cas, (2) revient à maximiser la vraisemblance gaussienne multivariée sous contraintes. Ainsi, *ITML* peut être vu comme un problème d'estimation de matrice de covariance.

Ensuite, *GMML* (*Geometric Mean Metric Learning*) [3] est un algorithme important en *apprentissage de distance*. Il atteint de bonnes performances sur plusieurs jeux de données avec un coût calculatoire faible. Le problème d'optimisation *GMML* est

$$\underset{\mathbf{A} \in \mathcal{S}_p^+}{\text{minimiser}} \frac{1}{n_{\mathcal{S}}} \sum_{(\mathbf{x}_l, \mathbf{x}_q) \in \mathcal{S}} d_{\mathbf{A}}^2(\mathbf{x}_l, \mathbf{x}_q) + \frac{1}{n_{\mathcal{D}}} \sum_{(\mathbf{x}_l, \mathbf{x}_q) \in \mathcal{D}} d_{\mathbf{A}^{-1}}^2(\mathbf{x}_l, \mathbf{x}_q). \quad (3)$$

L'intuition derrière ce problème est que  $d_{\mathbf{A}^{-1}}$  éloigne les points de classes différentes tandis que  $d_{\mathbf{A}}$  rapproche les points de même classe. La formulation (3) se réécrit

$$\underset{\mathbf{A} \in \mathcal{S}_p^+}{\text{minimiser}} \text{Tr}(\mathbf{A}^{-1}\mathbf{S}) + \text{Tr}(\mathbf{A}\mathbf{D}) \quad (4)$$

où  $\mathbf{S} = \frac{1}{n_{\mathcal{S}}} \sum_{k=1}^K \sum_{i=1}^{n_k} \mathbf{s}_{ki} \mathbf{s}_{ki}^T$  et  $\mathbf{D} = \frac{1}{n_{\mathcal{D}}} \sum_{i=1}^{n_{\mathcal{D}}} \mathbf{d}_i \mathbf{d}_i^T$ . Dans [3], la solution de (4) est calculée. Il s'agit du point milieu de la géodésique entre  $\mathbf{S}^{-1}$  et  $\mathbf{D}$ , i.e.  $\mathbf{A}^{-1} = \mathbf{S}^{-1} \#_{\frac{1}{2}} \mathbf{D} = \mathbf{S}^{-\frac{1}{2}} (\mathbf{S}^{\frac{1}{2}} \mathbf{D} \mathbf{S}^{\frac{1}{2}})^{\frac{1}{2}} \mathbf{S}^{-\frac{1}{2}}$ . Ensuite, [3] propose de généraliser cette solution par  $\mathbf{A}^{-1} = \mathbf{S}^{-1} \#_t \mathbf{D} = \mathbf{S}^{-\frac{1}{2}} (\mathbf{S}^{\frac{1}{2}} \mathbf{D} \mathbf{S}^{\frac{1}{2}})^t \mathbf{S}^{-\frac{1}{2}}$  avec  $t \in [0, 1]$ . Appliquée à des jeux de données réels, les auteurs obtiennent généralement de bonnes performances avec  $t$  petit ou nul (voir la Figure 3 de [3]) ce qui revient à  $\mathbf{A} \approx \mathbf{S}$ . Cette solution peut être interprétée statistiquement. Supposons que les données soient des réalisations de vecteurs aléatoires indépendants de moments un et deux dépendants de la classe,

$$\mathbf{x}_{kl} \stackrel{d}{=} \boldsymbol{\mu}_k + \boldsymbol{\Sigma}_k^{\frac{1}{2}} \mathbf{u}_{kl} \quad (5)$$

avec  $\boldsymbol{\mu}_k \in \mathbb{R}^p$ ,  $\boldsymbol{\Sigma}_k \in \mathcal{S}_p^+$ ,  $\mathbf{u}_{kl}$  le vecteur aléatoire centré réduit de la classe  $k$ , i.e.  $\mathbb{E}[\mathbf{u}_{kl}] = \mathbf{0}$  et  $\mathbb{E}[\mathbf{u}_{kl} \mathbf{u}_{kq}^T] = \mathbf{I}_p$  si  $kl = kq$ ,  $\mathbf{0}_p$  sinon. Il s'ensuit que  $\mathbf{s}_{ki} \stackrel{d}{=} \boldsymbol{\Sigma}_k^{\frac{1}{2}} (\mathbf{u}_l - \mathbf{u}_q)$ . Par conséquent, la matrice de covariance de  $\mathbf{s}_{ki}$  est égale à deux fois la matrice de covariance de la classe  $k$ ,  $\mathbb{E}[\mathbf{s}_{ki} \mathbf{s}_{ki}^T] \stackrel{d}{=} 2\boldsymbol{\Sigma}_k$ . Il en résulte que, en espérance,  $\mathbf{S}$  est deux fois la moyenne arithmétique des matrices de covariance des différentes classes,

$$\mathbb{E}[\mathbf{S}] = \frac{1}{n_{\mathcal{S}}} \sum_{k=1}^K \sum_{i=1}^{n_k} \mathbb{E}[\mathbf{s}_{ki} \mathbf{s}_{ki}^T] = 2 \sum_{k=1}^K \pi_k \boldsymbol{\Sigma}_k. \quad (6)$$

L'utilisation de  $\mathbf{S}$  dans la distance (1) rappelle en conséquence l'étape de préblanchiment des données de l'analyse discriminante linéaire.

## 1.2 Motivations et contributions

D'après la Section 1.1, *GMML* s'interprète comme une méthode à 2 étapes qui calcule, premièrement, la matrice de covariance empirique de chaque classe et, deuxièmement, leur moyenne arithmétique. Au lieu de faire ces 2 étapes, nous proposons dans ce papier, un nouvel algorithme, *Robust Geometric Metric Learning* (*RGML*) qui fera l'estimation jointe régularisée et robuste des matrices de covariance et de leur moyenne riemannienne.

Cette idée a initialement été proposée dans [4]. La contribution de ce papier est quadruple : **1)** *RGML* est formulé pour l'*apprentissage de distance*, **2)** *RGML* fait usage de la distance riemannienne sur  $\mathcal{S}_p^+$ , **3)** l'optimisation repose sur les géométries riemanniennes de  $\mathcal{S}_p^+$  et  $\mathcal{S}\mathcal{S}_p^+$  [5, 6], **4)** nous montrons que cette méthode est robuste et performante sur des données réelles partiellement mal étiquetées.

## 2 Formulation du problème

### 2.1 Formulation générale de *RGML*

La formulation du problème d'optimisation *RGML* est

$$\underset{\theta \in \mathcal{M}_{p,K}}{\text{minimiser}} \left\{ h(\theta) = \sum_{k=1}^K \pi_k [\mathcal{L}_k(\mathbf{A}_k) + \lambda d^2(\mathbf{A}, \mathbf{A}_k)] \right\} \quad (7)$$

où  $\theta = (\mathbf{A}, \{\mathbf{A}_k\})$ ,  $\mathcal{M}_{p,K} = (\mathcal{S}_p^+)^{K+1}$ ,  $\mathcal{L}_k$  est une fonction de coût d'estimation de la matrice de covariance dépendant de  $\{\mathbf{s}_{ki}\}$ ,  $\lambda > 0$  et  $d$  est une distance entre deux matrices. Dans les prochaines sous-sections, deux fonctions  $\mathcal{L}_k$  seront considérées : la log-vraisemblance négative gaussienne et la fonction de coût de Tyler. Une fois que (7) est résolu, la matrice centrale  $\mathbf{A}$  est utilisée dans la distance de Mahalanobis (1).

Nous continuons en expliquant plus en détails la fonction de coût (7). Pour cela, nous regardons les effets à variables fixées. Tout d'abord, pour une matrice centrale  $\mathbf{A}$  fixe, (7) se réduit à  $k$  problèmes séparables

$$\underset{\mathbf{A}_k}{\text{minimiser}} \mathcal{L}_k(\mathbf{A}_k) + \lambda d^2(\mathbf{A}, \mathbf{A}_k), \quad (8)$$

dont les solutions sont des estimations régularisées de  $\{\boldsymbol{\Sigma}_k\}$ .

Deuxièmement, pour  $\{\mathbf{A}_k\}$  fixes, résoudre (7) calcule le barycentre riemannien des matrices  $\{\mathbf{A}_k\}$ .

$$\underset{\mathbf{A} \in \mathcal{S}_p^+}{\text{minimiser}} \sum_{k=1}^K \pi_k d^2(\mathbf{A}, \mathbf{A}_k). \quad (9)$$

Par exemple, si  $d$  est la distance euclidienne  $d_E(\mathbf{A}, \mathbf{A}_k) = \|\mathbf{A} - \mathbf{A}_k\|_F$ , alors le minimum de (9) est la moyenne arithmétique  $\sum_{k=1}^K \pi_k \mathbf{A}_k$ . Dans la suite de l'article, nous considérons la distance riemannienne sur  $\mathcal{S}_p^+$  [5], c'est-à-dire

$$d_R(\mathbf{A}, \mathbf{A}_k) = \left\| \log_m \left( \mathbf{A}^{-\frac{1}{2}} \mathbf{A}_k \mathbf{A}^{-\frac{1}{2}} \right) \right\|_F. \quad (10)$$

Une propriété remarquable de  $d_R$  (10) est son invariance aux transformations affines. En effet, pour tout  $C$  inversible, on a  $d_R(CAC^T, CA_kC^T) = d_R(A, A_k)$ . Ainsi, si  $\{s_{ki}\}$  est transformé en  $\{Cs_{ki}\}$  alors le minimum  $(A, \{A_k\})$  de (7) devient  $(CAC^T, \{CA_kC^T\})$ . Enfin, nous soulignons que l'optimisation de (7) est effectuée par rapport à toutes les matrices  $A$  et  $\{A_k\}$  en même temps. Ainsi, les matrices de covariance régularisées  $\{A_k\}$  sont estimées tout en calculant leur barycentre inconnu  $A$ .

## 2.2 RGML gaussien

Pour obtenir une expression de la fonction de coût  $h$  (7), il reste à spécifier les fonctions  $\mathcal{L}_k$ . L'hypothèse la plus classique sur la distribution des données est l'hypothèse gaussienne (par exemple utilisée dans *ITML*). Ainsi, les premières fonctions  $\mathcal{L}_k$  considérées sont les log-vraisemblances négatives gaussiennes multivariées centrées

$$\mathcal{L}_{G,k}(A) = \frac{1}{n_k} \sum_{i=1}^{n_k} s_{ki}^T A^{-1} s_{ki} + \log |A|. \quad (11)$$

Avec cette log-vraisemblance négative, le problème *RGML* d'optimisation (7) devient

$$\underset{\theta \in \mathcal{M}_{p,K}}{\text{minimiser}} \left\{ h_G(\theta) = \sum_{k=1}^K \pi_k [\mathcal{L}_{G,k}(A_k) + \lambda d_R^2(A, A_k)] \right\} \quad (12)$$

et est appelé *RGML Gaussien*.

## 2.3 RGML Tyler

Lorsqu'une petite proportion des échantillons a un comportement aberrant, les méthodes d'estimation robuste de la matrice de covariance constituent un choix privilégié. Dans un contexte de classification, ce cas se produit lorsque les données sont mal étiquetées. Un estimateur classique est l'estimateur de Tyler [7] qui minimise la fonction de coût suivante

$$\mathcal{L}_{T,k}(A) = \frac{p}{n_k} \sum_{i=1}^{n_k} \log (s_{ki}^T A^{-1} s_{ki}) + \log |A|. \quad (13)$$

Une remarque importante est que (13) est invariant par facteur d'échelle de  $A$ . En effet, pour tout  $\alpha > 0$ , il est facile de vérifier que  $\mathcal{L}_{T,k}(\alpha A) = \mathcal{L}_{T,k}(A)$ . Ainsi, une contrainte de déterminant unitaire est ajoutée à (7) pour fixer les échelles de  $\{A_k\}$ . De plus, la distance riemannienne (10) est également celle de  $\mathcal{SS}_p^+$ . Ainsi, nous choisissons de contraindre également  $A$  pour qu'elle soit la moyenne riemannienne de  $\{A_k\}$  sur  $\mathcal{SS}_p^+$ . Nous désignons par  $\mathcal{SM}_{p,K}$  ce nouvel espace de paramètres

$$\mathcal{SM}_{p,K} = \{\theta \in \mathcal{M}_{p,K}, |A| = |A_k| = 1, \forall k \in \llbracket 1, K \rrbracket\}. \quad (14)$$

Ainsi, le problème d'optimisation (7) avec la fonction de coût de Tyler (13) devient

$$\underset{\theta \in \mathcal{SM}_{p,K}}{\text{minimiser}} \left\{ h_T(\theta) = \sum_{k=1}^K \pi_k [\mathcal{L}_{T,k}(A_k) + \lambda d_R^2(A, A_k)] \right\} \quad (15)$$

et est appelé *RGML Tyler*.

## 2.4 Optimisation

Les problèmes (12) et (15) sont résolus à l'aide de descentes de gradient riemanniennes [8]. Les variétés riemanniennes d'intérêt sont  $\mathcal{M}_{p,K}$  et  $\mathcal{SM}_{p,K}$  et leurs outils d'optimisation sont hérités directement de ceux de  $\mathcal{S}_p^+$  et  $\mathcal{SS}_p^+$  [5, 6]. Dans le reste de cette sous-section, seules quelques remarques sur l'optimisation sont données et les détails techniques sont omis pour des raisons d'espace<sup>1</sup>. Tout d'abord, les problèmes (12) et (15) sont géodésiquement convexes, *i.e.* convexes le long des géodésiques de  $\mathcal{M}_{p,K}$  et  $\mathcal{SM}_{p,K}$ . Par conséquent, les descentes de gradient sont peu sensibles à l'initialisation. Ensuite, la complexité d'une itération de la descente de gradient *RGML* est dominée par celle de l'évaluation du gradient, en pratique, calculé par différenciation automatique. Ainsi, nous donnons la complexité de l'évaluation de la fonction de coût :  $\mathcal{O}(n_S p^2 + K p^3)$ . La complexité de la forme fermée de *GMML* est  $\mathcal{O}(np^2 + p^3)$  où  $n = \max\{n_S, n_D\}$ .  $K$  étant le nombre de classes (et donc peu élevé), les deux algorithmes ont des complexités proches. Enfin, l'hyperparamètre  $\lambda$  dans (7) peut être choisi par validation croisée.

## 3 Expériences

Dans cette section, nous présentons une application de la méthode *RGML* développée dans la Section 2. En effet, nous l'appliquons sur des jeux de données du *UCI machine learning repository* [9]. Trois jeux de données de classification sont considérés : *Wine*, *Vehicle* et *Iris*. Leurs statistiques, nombre et taille des données ainsi que le nombre de classes, sont présentées dans la Table 1. Ces ensembles de données ont des classes quasiment équiprobables, *c'est-à-dire* que les classes ont des nombres de données proches. Les nombres de paires générées dans  $\mathcal{S}$  et  $\mathcal{D}$  sont  $n_S = n_D = 75 K(K-1)$  (voir [2] et [3]).

La validation se fait selon un protocole très classique en *apprentissage de distance*. **1)** Une matrice  $A$  est estimée via une méthode d'*apprentissage de distance*. **2)** Les données  $\{x_l\}$  sont multipliées par  $A^{-\frac{1}{2}}$  pour obtenir  $\{A^{-\frac{1}{2}}x_l\}$ . **3)** Les données  $\{A^{-\frac{1}{2}}x_l\}$  sont classées en utilisant un *k-plus proches voisins* avec 5 voisins. Ainsi, la classification est faite en utilisant la distance de Mahalanobis  $d_A$  (1) définie dans l'Introduction. Cette classification est répétée 200 fois via une validation croisée. La répartition des données dans les ensembles d'entraînement et de test est de 50/50%. L'erreur moyenne de mauvaise classification est reportée dans la Table 1. Afin de montrer la robustesse de la méthode proposée, des données mal étiquetées sont introduites. Pour ce faire, nous sélectionnons au hasard des données dans l'ensemble d'apprentissage dont les étiquettes sont ensuite changées au hasard.

Les implémentations de la validation croisée ainsi que du *k-plus proches voisins* proviennent de la bibliothèque *scikit-learn* [10]. La valeur choisie du paramètre  $\lambda$  est 0.05. Sa valeur a peu d'impact sur les performances tant qu'elle est ni trop petite ni trop grande. Ils sont comparés aux algorithmes

1. Un programme Python de ces méthodes est disponible ici : [https://github.com/antoinecollas/robust\\_metric\\_learning](https://github.com/antoinecollas/robust_metric_learning).

TABLE 1 – Taux d’erreurs de classification calculés sur 3 jeux de données : Wine, Vehicle, et Iris. Les meilleurs résultats sont en **gras**. Les taux d’étiquetage erroné indiquent les taux d’étiquettes modifiées de façon aléatoire dans l’ensemble d’apprentissage.

Méthodes	Wine $p = 13, n = 178, K = 3$				Vehicle $p = 18, n = 846, K = 4$				Iris $p = 4, n = 150, K = 3$			
	Taux d’étiquetage erroné				Taux d’étiquetage erroné				Taux d’étiquetage erroné			
	0%	5%	10%	15%	0%	5%	10%	15%	0%	5%	10%	15%
Identité	30.12	30.40	31.40	32.40	38.27	38.58	39.46	40.35	3.93	4.47	5.31	<b>6.70</b>
Cov. emp.	10.03	11.62	13.70	17.57	23.59	24.27	25.24	26.51	12.57	13.38	14.93	16.68
ITML - Identité	3.12	4.15	5.40	<b>7.74</b>	24.21	23.91	24.77	26.03	3.04	4.47	5.31	<b>6.70</b>
ITML - Cov. emp.	2.45	4.76	6.71	10.25	23.86	23.82	24.89	26.30	3.05	13.38	14.92	16.67
GMML	2.16	3.58	5.71	9.86	21.43	22.49	23.58	25.11	2.60	5.61	9.30	12.62
LMNN	4.27	6.47	7.83	9.86	20.96	24.23	26.28	28.89	3.53	9.59	11.19	12.22
RGML - gaussien	<b>2.07</b>	<b>2.93</b>	5.15	9.20	<b>19.76</b>	21.19	22.52	24.21	<b>2.47</b>	5.10	8.90	12.73
RGML - Tyler	<b>2.12</b>	<b>2.90</b>	<b>4.51</b>	8.31	19.90	<b>20.96</b>	<b>22.11</b>	<b>23.58</b>	<b>2.48</b>	<b>2.96</b>	<b>4.65</b>	7.83

d’apprentissage de distance classiques : la matrice identité, la matrice de covariance empirique calculée sur toutes les données, *ITML* [2], *GMML* [3], et *LMNN* [11]. Les implémentations de la bibliothèque *metric-learn* [12] sont utilisées pour les trois derniers algorithmes.

À partir de la Table 1, plusieurs observations sont faites. Tout d’abord, sur les données brutes (*c’est-à-dire* lorsque le taux de mauvais étiquetage est de 0%) le *RGML gaussien* est toujours l’algorithme le plus performant parmi ceux testés. De plus, le *RGML Tyler* est proche avec un écart maximum de 0.26% par rapport au *RGML gaussien*. Ensuite, le *RGML Tyler* est l’algorithme le plus performant lorsque le taux d’erreur d’étiquetage est de 5 % ou de 10 %. Lorsque le taux de mauvais étiquetage est de 15 %, *RGML Tyler* est l’algorithme le plus performant pour le jeu de données *Vehicle* et il n’est battu que par *ITML - identité* sur les deux autres jeux de données. Ceci montre l’intérêt de considérer la fonction de coût de Tyler (13) en présence d’un mauvais étiquetage. Enfin, les algorithmes *RGML* sont rapides. La Figure 1 montre que *RGML gaussien* et *RGML Tyler* convergent en moins de 20 itérations sur le jeu de données *Wine* (sans changement d’étiquettes).

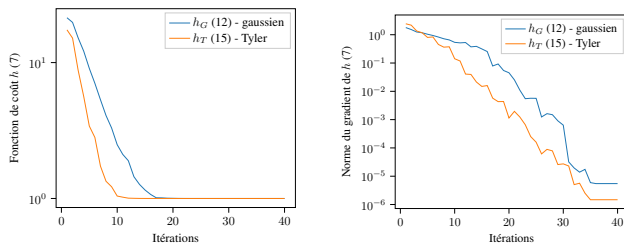


FIGURE 1 – Gauche : évolution des fonctions de coût *RGML gaussien* (12) et *RGML Tyler* (15) en fonction des itérations. Droite : évolution des normes des gradients des fonctions (12) et (15). L’optimisation est réalisée sur le jeu de données *Wine*.

## Références

[1] “A tutorial on distance metric learning : Mathematical foundations, algorithms, experimental analysis, prospects and challenges,” *Neurocomputing*, vol. 425.

[2] J. V. Davis, B. Kulis, P. Jain, S. Sra, and I. S. Dhillon, “Information-theoretic metric learning,” in *Proceedings of the 24th International Conference on Machine Learning*. New York, NY, USA : Association for Computing Machinery, 2007.

[3] P. H. Zadeh, R. Hosseini, and S. Sra, “Geometric mean metric learning,” in *Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48*. JMLR.org, 2016.

[4] E. Ollila, I. Soloveychik, D. E. Tyler, and A. Wiesel, “Simultaneous penalized m-estimation of covariance matrices using geodesically convex optimization,” 2016.

[5] L. T. Skovgaard, *Scandinavian Journal of Statistics*, no. 4.

[6] X. Pennec, P. Fillard, and N. Ayache, “A riemannian framework for tensor computing,” *International Journal of Computer Vision*, vol. 66, 2005.

[7] D. E. Tyler, “A Distribution-Free  $M$ -Estimator of Multivariate Scatter,” *The Annals of Statistics*, vol. 15, no. 1.

[8] P.-A. Absil, R. Mahony, and R. Sepulchre, *Optimization algorithms on matrix manifolds*. Princeton University Press, 2008.

[9] D. Dua and C. Graff, “UCI machine learning repository,” 2017.

[10] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, “Scikit-learn : Machine learning in Python,” *Journal of Machine Learning Research*, vol. 12.

[11] K. Q. Weinberger and L. K. Saul, “Distance metric learning for large margin nearest neighbor classification,” *The Journal of Machine Learning Research*, vol. 10.

[12] W. de Vazelhes, C. Carey, Y. Tang, N. Vauquier, and A. Bellet, “metric-learn : Metric Learning Algorithms in Python,” *Journal of Machine Learning Research*, vol. 21, no. 138, 2020.