

# 確率的ダイクストラ法による複数マハラノビス行列の学習 Stochastic Dykstra Algorithms for Learning Multiple Mahalanobis Matrices

大沼 由弥<sup>†</sup>  
Yuya Onuma

佐藤 貴亮<sup>†</sup>  
Takaaki Sato

松澤 知己<sup>†</sup>  
Tomoki Matsuzawa

加藤 毅<sup>† § ¶</sup>  
Tsuayoshi Kato

## 1. はじめに

パターン認識問題において、特徴ベクトル間の距離計量を判別的に学習することによって識別性能が向上することは、数多くの研究報告によって裏付けられている。今日においては、対象がベクトル形式ではない場合における計量学習算法に研究の関心が拡大している [3]。本論文では、解析対象  $\mathbf{x} \in \mathcal{X}$  が  $M$  個の行列の集合  $\{\Phi_m(\mathbf{x})\}_{m=1}^M$  で表現される場合を考える。ただし、 $m$  個目の行列のサイズは  $n_m \times n'_m$  であり、 $\Phi_m: \mathcal{X} \rightarrow \mathbb{R}^{n_m \times n'_m}$  は対象  $\mathbf{x}$  から  $m$  個目の行列を取り出す特徴抽出器である。本論文では、次の形式で表される距離関数  $D_{\Phi}(\cdot, \cdot; \mathcal{W}): \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  を議論する：

$$D_{\Phi}(\mathbf{x}_1, \mathbf{x}_2; \mathcal{W}) := \frac{1}{M} \sum_{m=1}^M \left\langle \mathbf{W}_m, (\Phi_m(\mathbf{x}_1) - \Phi_m(\mathbf{x}_2)) (\Phi_m(\mathbf{x}_1) - \Phi_m(\mathbf{x}_2))^{\top} \right\rangle. \quad (1)$$

ただし、 $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}$  であり、 $\mathcal{W} := (\mathbf{W}_1, \dots, \mathbf{W}_M)$  は距離関数のパラメータで、 $M$  個の正定値行列  $\mathbf{W}_1 \in \mathbb{S}_{++}^{n_1}, \dots, \mathbf{W}_M \in \mathbb{S}_{++}^{n_M}$  からなる。 $n = \max_m n_m$  とおく。たとえば、 $M=1$  で  $\Phi_m(\mathbf{x}) \in \mathbb{R}^{n \times 1}$  の場合、標準的なマハラノビス距離に等しい。 $M=1$  で  $\Phi_m(\mathbf{x}) \in \mathbb{S}_+^n$  が共分散記述子もしくはその *spectral variants* の場合に関しては、すでに幾つかの論文で議論されている (e.g. [3])。  $M$  を 2 以上にすれば、複数の共分散記述子間距離を統合した距離となる。入力空間  $\mathcal{X}$  が  $M$  モードのテンソルの集合のとき、 $\Phi_m(\mathbf{x})$  を mode- $m$  flattening とすることで、テンソル間の距離関数を表すこともできる。

本論文では、上述の距離関数 (1) のパラメータ  $\mathcal{W}$  の値を判別的に学習する算法を提案する。本研究では、ベクトル間の計量学習算法として有名な ITML [2] に倣って、ブレグマン距離で目的関数を構成した。ITML は、ダイクストラ法 [1] を使って目的関数の最小化を行うものである。ダイクストラ法は、複数の半空間の共通集合内で最も近い点を探すための反復解法であり、各反復で半空間を巡回的に選択して解を射影することで解を最適解に収束させるものである。ITML では、 $\Phi_m(\mathbf{x})$  のランクが 1 であることから半空間への射影は閉形式で求められることを利用していた。一方、距離関数を (1) のように一般化してしまうと、半空間の射影は閉形式では求まらない。本研究では、(a) 計算量  $O(Mn^3)$  で半空間への射影を求めることができる効率的な算法を発見した。この理論的発見のみならず、本研究では次の 2 点の実験的発見を得た：(b) ダイクストラ法に

おける半空間の選択を巡回的にではなく無作為に行うと、Objective Gap (最適値との差) が劇的に早く 0 に近づく；(c) 共分散行列を複数化して計量学習することでパターン認識における汎化性能が向上する。

## 2. ブレグマン距離

本研究では、2 種類のブレグマン距離を用いる。ブレグマン距離は、シード関数  $\varphi: \text{dom}(\varphi) \rightarrow \mathbb{R}$  を通じて一般に次のように定義されている：

$$\text{BD}_{\varphi}(\Theta, \Theta_0) = \varphi(\Theta) - \varphi(\Theta_0) - \langle \nabla \varphi(\Theta_0), \Theta - \Theta_0 \rangle,$$

ただし、シード関数は連続微分可能で狭義凸でなければならない。例えば、 $\varphi((\Phi_m(\mathbf{x}))_{m=1}^M) := \sum_{m=1}^M \langle \mathbf{W}_m, \Phi_m(\mathbf{x}) \Phi_m(\mathbf{x})^{\top} \rangle / M$  とおくと距離関数 (1) を得る。ブレグマン距離は、非対称なため、一般に距離の公理を満たさないが、写像  $\mathbf{x} \mapsto \{\Phi_m(\mathbf{x})\}_{m=1}^M$  が単射のときは  $D_{\Phi}$  は距離の公理を満たす。 $D_{\Phi}$  は、本研究で用いる 2 種類のブレグマン距離のうち、1 種類目となる。

2 種類目は次のようなブレグマン距離である。後述する  $\mathcal{W}$  の学習算法では、スラック変数  $\xi = [\xi_k]_{k=1}^K \in \mathbb{R}^K$  を導入し、 $K$  個の半空間

$$\mathcal{C}_k := \left\{ (\mathcal{W}, \xi) \mid \frac{y_k}{M} \sum_{m=1}^M \langle \mathbf{A}_{m,k}, \mathbf{W}_m \rangle \leq y_k b_k \xi_k \right\} \quad (2)$$

の共通集合から解を探す。ただし、 $(\mathbf{A}_{m,k}, b_k, y_k) \in \mathbb{S}^{n_m} \times \mathbb{R}_{++} \times \{\pm 1\}$  とする。後述する提案学習算法では、シード関数として、

$$\varphi(\mathcal{W}, \xi) := -\frac{1}{M} \sum_{m=1}^M \log \det(\mathbf{W}_m) + c \sum_{k=1}^K \varphi_1(\xi_k)$$

を用いた。 $c$  は正の定数である。第 2 項の  $\varphi_1: \mathbb{R}_{++} \rightarrow \mathbb{R}$  には、次の 3 種類の関数を試した： $\varphi_{\text{is}}(\xi_k) := -\log(\xi_k)$ 、 $\varphi_{\text{e}}(\xi_k) := \xi_k^2/2$ 、 $\varphi_{\text{e}}(\xi_k) := (\log \xi_k - 1)\xi_k$ 。1 番目のシード関数  $\varphi_{\text{is}}$  から生成されるブレグマン距離は、 $M=1$ 、 $n'_m=1$  のとき ITML [2] の目的関数と等価になる。

## 3. 計量学習問題の定式化

多クラス分類問題を考える。 $\ell$  個の訓練用例題の集合  $\{(\mathbf{x}_i, \omega_i)\}_{i=1}^{\ell}$  が所与とする。ただし、 $\mathbf{x}_i \in \mathcal{X}$  は入力データであり、 $\omega_i \in \mathbb{N}$  はクラス番号である。これら  $\ell$  個の訓練用例題に対し、 $K$  個のペアを  $(i_1, j_1), \dots, (i_K, j_K) \in \mathbb{N}_{\ell} \times \mathbb{N}_{\ell}$  を考え、各ペアに対して次の制約を与える：

$$D_{\Phi}(\mathbf{x}_{i_k}, \mathbf{x}_{j_k}; \mathcal{W}) \begin{cases} \leq b_{\text{ub}} \xi_k & \text{if } \omega_{i_k} = \omega_{j_k}, \\ \geq b_{\text{lb}} \xi_k & \text{if } \omega_{i_k} \neq \omega_{j_k}. \end{cases} \quad (3)$$

<sup>†</sup>群馬大学大学院理工学府

<sup>§</sup>群馬大学次世代モビリティ社会実装センター (CRANTS)

<sup>¶</sup>早稲田大学規範科学総合研究所 (IIRS)

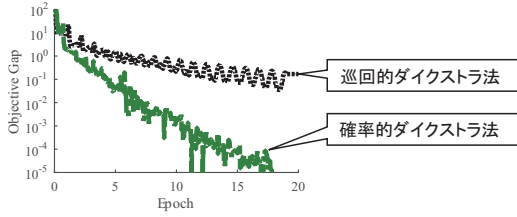


図1: 収束の様子. 半平面の選択を確率的にすると早く収束する.

ただし,  $\xi_k = 1$  のとき,  $b_{ub}$  および  $b_{lb}$  は, それぞれ, 同じクラス間の距離の上限, 異なるクラス間の距離の下限を定める定数となる. ここで,  $\omega_{i_k} = \omega_{j_k}$  なる  $k \in \mathbb{N}_K$  には,  $(y_k, b_k) = (+1, b_{ub})$ ,  $\omega_{i_k} \neq \omega_{j_k}$  なる  $k \in \mathbb{N}_K$  には,  $(y_k, b_k) = (-1, b_{lb})$  とおいて,  $k = 1, \dots, K$ ,  $m = 1, \dots, M$  に対して,

$$\mathbf{A}_{m,k} := (\Phi_m(\mathbf{x}_{i_k}) - \Phi_m(\mathbf{x}_{j_k})) (\Phi_m(\mathbf{x}_{i_k}) - \Phi_m(\mathbf{x}_{j_k}))^\top$$

とおくと, (3) で表される第  $k$  制約を満たす  $(\mathcal{W}, \xi)$  の集合は半空間  $\mathcal{C}_k$  になることがわかる.

$(\mathcal{W}, \xi)$  を点  $(\mathcal{W}_0, \mathbf{1}_K)$  から共通集合  $\mathcal{C}_1 \cap \dots \cap \mathcal{C}_K$  へのブレグマン射影, すなわち, 最適化問題

$$\begin{aligned} \min \quad & \text{BD}_\varphi((\mathcal{W}, \xi), (\mathcal{W}_0, \mathbf{1})), \\ \text{wrt} \quad & (\mathcal{W}, \xi) \in \mathcal{C}_1 \cap \dots \cap \mathcal{C}_K \end{aligned} \quad (4)$$

を解くことで決定することとする. ただし,  $\mathcal{W}_0 := (\mathbf{I}_{n_1}, \dots, \mathbf{I}_{n_M})$  とする. すると, マハラノビス行列は単なる単位行列に, また, 距離制約は単なる  $\xi_k = 1$  に近づくように  $\mathcal{W}$  の値が決めることになる.

#### 4. 確率的ダイクストラ法

射影問題 (4) をダイクストラ法で解くとすると, 各反復で,  $K$  個の半空間から一つの半空間  $\mathcal{C}_k$  を選んで, 現在の解を  $\mathcal{C}_k$  にブレグマン射影することになる. この射影を繰り返すことで最適解に収束することが理論的に保証されている [1]. 第  $t$  反復における解を  $(\mathcal{W}_t, \xi_t)$ , 第  $m$  マハラノビス行列を  $\mathbf{W}_{m,t}$  とおき,  $\mathcal{C}_k$  に現在の解を射影するとする. さらに,  $J_t: \mathbb{R} \rightarrow \mathbb{R}$  を

$$\begin{aligned} J_t(\delta) := & \sum_{m=1}^M \langle \mathbf{A}_{m,k}, (\mathbf{W}_{m,t-1}^{-1} + y_k \delta \mathbf{A}_{m,k})^{-1} \rangle \\ & - M b_k \nabla \varphi_1^* (\nabla \varphi_1(\xi_{t-1}) + y_k \delta b_k / c), \end{aligned}$$

とおくと,  $\mathcal{C}_k$  の境界への射影は各  $m$  に対し,  $\mathbf{W}_{m,t-1}^{-1} + \delta y_k \mathbf{A}_{m,k}$  が正定値である範囲で, 1 変数非線形方程式  $J_t(\delta) = 0$  の根を求める問題に帰着できる. ただし,  $\varphi_1^*$  は  $\varphi_1$  の凸共役である.  $J_t(\delta)$  の第1項には  $M$  個の逆行列を含んでいるため, ナイーブにニュートン法などで根を求めると, 1 回の射影に  $O(LMn^3)$  の計算量がかかってしまう. ただし,  $L$  はニュートン法内部で関数  $J_t(\delta)$  の値を評価する回数である. 本研究では, 次の定理を発見した:

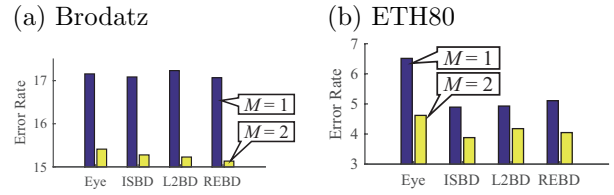


図2: 多クラス分類の誤分類率. Eye:  $\mathbf{W}_m = \mathbf{I}$  に固定, ISBD, L2BD, REBD: 計量学習あり. それぞれ  $\varphi_1 = \varphi_{is}, \varphi_{l2}, \varphi_e$ .

**Theorem 1.**  $L \in O(n^2)$  のとき, 1 回のブレグマン射影は  $O(Mn^3)$  の計算量で済む.

**Proof Sketch:** 各  $m = 1, \dots, M$  に対し,

$$\langle \mathbf{A}_{m,k}, (\mathbf{W}_{m,t-1}^{-1} + \delta y_k \mathbf{A}_{m,k})^{-1} \rangle = \sum_{i=1}^{n_m} \frac{1}{d_{i,m,k} + y_k \delta} \quad (5)$$

を満たす  $n_m$  個の実数  $\{d_{i,m,k}\}_{i=1}^{n_m}$  を  $O(n_m^3)$  で見つけることができる. いったん  $\{d_{i,m,k}\}_{i=1}^{n_m}$  を見つければ,  $J(\delta) = 0$  の根は  $O(LMn)$  で見つけられる. 1 回のブレグマン射影の計算量は  $O(Mn^3)$  となる.  $\square$

**半空間乱択の効果** 本研究のもう一つの貢献は, 半空間の選び方を巡回的にするのではなく, 無作為に選ぶと劇的に早く最適解に収束することができることを経験的に発見したことである. 図1は,  $\varphi_1 = \varphi_{is}$ ,  $M = 3$ ,  $n_1 = 10$ ,  $n_2 = 15$ ,  $n_3 = 20$  を用いたときの例である. この例では, 巡回的に半空間を選んだ場合は, 20 エポック経っても  $10^{-1}$  近くの Objective gap が残っているが, 半空間の選択法を確率的にすると Objective gap に  $10^{-5}$  未満にまで到達している.

**パターン認識性能** 2 個の画像データセット Brodatz および ETH80 を用いて 3-近傍識別による多クラス分類のベンチマークをとった. クラス数はそれぞれ 112 および 8 である. 各画像に対し,  $(M = 2)$  個の共分散記述子を得た. Brodatz では, 1 個目の共分散記述子は各画素から得られる  $(n_1 = 5)$  次元ベクトルから構成, 2 個目は画像解像度を半分にしたときの各画素から得られる  $(n_2 = 5)$  次元ベクトルから抽出した. ETH80 では  $n_1 = n_2 = 9$  の記述子を抽出した. 図2に,  $M = 1$  の場合との誤分類率の比較を示す. 計量学習により多クラス分類の汎化性能が向上することが確認できる.

謝辞: 本研究は JSPS 科研費 26249075, 40401236 の助成を受けたものである.

#### 参考文献

- [1] Y. Censor and S. Reich. The dykstra algorithm with bregman projections. *Comm. in Applied Anal.*, 2:407–419, 1998.
- [2] J. V. Davis et al. Information-theoretic metric learning. In *ICML*, pages 209–216, New York, NY, USA, 2007. ACM.
- [3] T. Matsuzawa, R. Relator, J. Sese, and T. Kato. Stochastic dykstra algorithms for metric learning with positive definite covariance descriptors. In *ECCV2016*, pages 786–799, 2016.