

距離しきい値を自動調整できる ITML 型計量学習法

大沼 由弥*

加藤 毅^{*,†,‡}

* 群馬大学大学院理工学府 電子情報・数理教育プログラム

† 群馬大学次世代モビリティ社会実装研究センター

‡ 早稲田大学規範科学総合研究所

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。本論文では、2つの入力 $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}$ に対するマハラノビス距離を一般化させた距離計量

$$D_{\Phi}(\mathbf{x}_1, \mathbf{x}_2; \mathbf{W}) :=$$

$$\text{tr}(\mathbf{W}(\Phi(\mathbf{x}_1) - \Phi(\mathbf{x}_2))(\Phi(\mathbf{x}_1) - \Phi(\mathbf{x}_2))^{\top}) \quad (1)$$

を考え、マハラノビス行列 $\mathbf{W} \in \mathbb{S}_{++}^n$ の教師あり学習を行う方法を議論する。ただし、 \mathbb{S}_{++}^n は n 次狭義正定値行列の集合である。例えば、 $\Phi: \mathcal{X} \rightarrow \mathbb{R}^n$ とすると (1) は、簡単な式変形により、標準的なベクトル間のマハラノビス距離になることがわかる。筆者らの研究グループでは $\Phi(\mathbf{x})$ を \mathbf{x} から得られる共分散記述子に一般化した場合の計量学習法を開発していた [2]。本研究では、複数のマハラノビス行列を含むように拡張した距離における計量 (1) の重みパラメータ \mathbf{W} を学習するアルゴリズムを開発した。しかし簡単のため、本論文では単一のマハラノビス行列を用いる場合を中心に議論する。複数にした場合の理論の詳細は、文献 [3] に記した。

ベクトル空間における計量学習の方法としては ITML [1] が広く用いられており、多くの派生法が生まれている。ITML はブレグマン射影問題で学習法を構成していることから、最適化にダイクストラ法 [2] を適用できる。ITML やその派生法は距離しきい値 b_0 というハイパーパラメータがあり、 b_0 の値を利用者があらかじめ決めなくてはならないという短所を持っている。その距離しきい値 b_0 が汎化性能に敏感であるため慎重に定める必要があることが、筆者らの実験で判明した。

本研究では、(i) ブレグマン射影問題の枠組みから逸脱することなく、距離しきい値を自動設定できる新しい計量学習法を開発した。 (ii) 複数の共分散記述子を扱えるように距離関数を拡張した。 さらに、(iii) ダイクストラ法の各反復に必要な半空間への射影を $O(Mn^3)$ で求められることを発見した。 ただし、 M はマハラノビス行列の個数である。また、(iv) 距離しきい値を自動設定しても、手動で距離しきい値を設定する従来の計量学習の性能と competitive であること、および、共分散記述子を複数化して計量学習を行うことでパターン認識にお

ける汎化性能が向上することを実データを使った実験により確認した。なお、理論の詳細と実験結果は、文献 [3] を参照されたい。

2 ブレグマン距離

ブレグマン距離はシード関数 $\varphi: \text{dom}(\varphi) \rightarrow \mathbb{R}$ を通じて

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

のように定義される。シード関数は連続微分可能かつ狭義凸であることが仮定されている。たとえば、 $\varphi(\Phi(\mathbf{x})) := \text{tr}(\mathbf{W}\Phi(\mathbf{x})\Phi(\mathbf{x})^{\top})$ とおくと距離関数 (1) を得る。ITML では、計量学習問題をブレグマン射影問題で定式化した。ある凸集合 \mathcal{S} へのブレグマン射影とは、 \mathcal{S} の元のうち、ある点 Θ_0 からのブレグマン距離 $\text{BD}(\cdot, \Theta_0; \varphi)$ が最も近くなるもの、すなわち、 $\arg \min_{\Theta \in \mathcal{S}} \text{BD}(\Theta, \Theta_0)$ である。ブレグマン距離は狭義凸関数なので、射影点は一意に定まる。

3 従来法：距離しきい値 b_0 を固定した計量学習

距離関数のパラメータ \mathbf{W} の値を教師あり学習によって決定するために、 ℓ 個のラベルありデータ $\mathbf{x}_1, \dots, \mathbf{x}_{\ell} \in \mathcal{X}$ を使う。ITML 型の距離計量学習は、 ℓ 個の例題から K 個の例題ペアを選択する。最初の K_+ ($< K$) 個の例題ペア $(i_k, j_k) \in \mathbb{N}^2$ ($k = 1, \dots, K_+$) は、同じクラスに所属する例題 \mathbf{x}_{i_k} および \mathbf{x}_{j_k} を選ぶ。残りの K_- ($:= K - K_+$) 個の例題ペア $(i_k, j_k) \in \mathbb{N}^2$ ($k = K_+ + 1, \dots, K$) は、異なるクラスに所属する例題 \mathbf{x}_{i_k} および \mathbf{x}_{j_k} を選ぶ。そのうえで、 \mathbf{W} を K 個の制約

$$\begin{cases} D_{\Phi}(\mathbf{x}_{i_k}, \mathbf{x}_{j_k}; \mathbf{W}) \leq \xi_k, & \text{for } k = 1, \dots, K_+, \\ D_{\Phi}(\mathbf{x}_{i_k}, \mathbf{x}_{j_k}; \mathbf{W}) \geq \xi_k, & \text{for } k = K_+ + 1, \dots, K. \end{cases}$$

を満たす中から探す。ただし、 $\boldsymbol{\xi} := [\xi_1, \dots, \xi_K]^{\top}$ はスラック変数である。ITML 型の計量学習法では、 \mathbf{W} の値を決定する前に、利用者が距離しきい値 b_0 の値を決定しておく。スラック変数 $\boldsymbol{\xi}$ については、 $k = 1, \dots, K_+$ に対しては ξ_k が $b_k(b_0) := 2b_0$ から離れるほど、 $k = K_+ + 1, \dots, K$ に対しては ξ_k が $b_k(b_0) := b_0/2$ から離れるほどペナルティがかかるように目的関数が設計されている。具体的には、 $\mathbf{b}(b_0) := [b_1(b_0), \dots, b_K(b_0)]^{\top}$ とおいて、次のように目的関数を定義している：

$$P_o(\mathbf{W}, \boldsymbol{\xi}; b_0) := \text{BD}(\mathbf{W}, \mathbf{I}_n; \varphi_r) + c \text{BD}(\boldsymbol{\xi}, \mathbf{b}(b_0); \varphi_{\ell}). \quad (2)$$

ただし、 c は正則化パラメータである。 $P_o(\cdot, \cdot; b_0)$ の第2項はスラック変数 $\boldsymbol{\xi}$ が距離しきい値 $\mathbf{b}(b_0)$ から離れるほどかかるペナルティで、ITML ではそのシード関数に $\varphi_{\ell}(\boldsymbol{\xi}) := -\sum_{k=1}^K \log(\xi_k)$ が用いられている。第1項は \mathbf{W} の正定値性を保ち過学習を防ぐための正則化項であり、そのシード関数には $\varphi_r(\mathbf{W}) := -\log \det(\mathbf{W})$ が用いられている。

ITML-type metric learning algorithm without manual tuning of distance threshold

Yuya Onuma*, Tsuyoshi Kato^{*,†,‡}

* Education Program of Electronics and Informatics, Mathematics and Physics, Graduate School of Science and Technology, Gunma University

† Center for Research on Adoption of NextGen Transportation Systems, Gunma University

‡ Integrated Institute for Regulatory Science, Waseda University

半空間 C_k を

$$C_k := \{(\mathbf{W}, \boldsymbol{\xi}) \mid y_k D_{\Phi}(\mathbf{x}_{i_k}, \mathbf{x}_{j_k}; \mathbf{W}) \leq y_k \xi_k\},$$

と定義する。ただし、 $k \leq K_+$ には $y_k = +1$, $k > K_+$ には $y_k = -1$ とする。ITML 型学習法では、 K 個の半空間 C_k の共通集合の中から $P_o(\cdot, \cdot; b_0)$ を最小にする $(\mathbf{W}, \boldsymbol{\xi})$ を探す。 $P_o(\cdot, \cdot; b_0)$ はブレッグマン距離の和になっていることから1つのブレッグマン距離になる。よって、この最適化問題はブレッグマン射影問題になり、ダイクストラ法による数値的に安定した最適化が可能になる。しかし、距離しきい値 b_0 の値は、利用者が事前に定めるか、交差検証法などで決めるしかなかった。

4 提案法：距離しきい値 b_0 を自動調整する計量学習

本節では、距離しきい値 b_0 を定数ではなく、 \mathbf{W} や $\boldsymbol{\xi}$ と同時に最適化する変数として扱うことを考える。 b_0 に対する正則化項を加えて、目的関数を

$$P_{b_0}(\mathbf{W}, \boldsymbol{\xi}, b_0) := P_o(\mathbf{W}, \boldsymbol{\xi}; b_0) + c_0 \text{BD}(b_0, \mu_0; \varphi_{\ell_0}).$$

と定義する。ただし $c_0 > 0$ および $\mu_0 > 0$ は定数である。この目的関数で \mathbf{W} , $\boldsymbol{\xi}$ および b_0 を同時最適化する問題は、定数は増えたが、データへの依存性の高い定数は存在していない（実データを用いた実験により確認済み [3]）。しかし、 $P_{b_0}(\mathbf{W}, \boldsymbol{\xi}, b_0)$ は、もはや一般に1つのブレッグマン距離には変換できないためブレッグマン射影問題としては扱えず、さらに、凸関数でもないために最適化が困難になってしまう。

本研究で最も大きな理論的成果は、 \mathbf{W} , $\boldsymbol{\xi}$ および b_0 の同時最適化問題について、ある設定においては再びブレッグマン射影問題になることを発見したことである。実行可能領域 $\bigcap_{k=1}^K C_k$ 内で目的関数 $P_{b_0}(\mathbf{W}, \boldsymbol{\xi}, b_0)$ を最小化する問題は、 $P_b(\mathbf{W}, \boldsymbol{\xi}) := \min_{b_0 \in \text{dom}(\varphi_{\ell_0})} P_{b_0}(\mathbf{W}, \boldsymbol{\xi}, b_0)$ を使って、

$$\min P_b(\mathbf{W}, \boldsymbol{\xi}) \text{ wrt } (\mathbf{W}, \boldsymbol{\xi}) \in \bigcap_{k=1}^K C_k \quad (3)$$

と等価的に書き換えることができる。本研究では、2つのシード関数を

$$\varphi_{\ell_0}(b_0) := \frac{1}{2} b_0^2, \quad \varphi_{\ell}(\boldsymbol{\xi}) := \frac{1}{2} \|\boldsymbol{\xi}\|^2 \quad (4)$$

とおいたとき、目的関数 $P_b(\cdot, \cdot)$ は、定数を除いて1つのブレッグマン距離になることを見つけた。

Theorem 1. φ_{ℓ_0} および φ_{ℓ} を (4) のように定義すると、任意の $\boldsymbol{\xi} \in \mathbb{R}^K$ について

$$P_b(\mathbf{W}, \boldsymbol{\xi}) + C = \text{BD}((\mathbf{W}, \boldsymbol{\xi}), (\mathbf{I}, \boldsymbol{\xi}_0); \varphi_{\text{tot}}), \quad (5)$$

を満たす $\mathbf{G} \in \mathbb{S}_{++}^K$, $\boldsymbol{\xi}_0 \in \mathbb{R}^K$, および $C \in \mathbb{R}$ が存在する。ただし、右辺のブレッグマン距離のシード関数は

$$\varphi_{\text{tot}}(\mathbf{W}, \boldsymbol{\xi}) := \varphi_{\text{r}}(\mathbf{W}) + \frac{1}{2} \langle \boldsymbol{\xi}, \mathbf{G}\boldsymbol{\xi} \rangle. \quad (6)$$

とする。（証明は [3] 参照。）

複数共分散記述子への拡張 ここまでは簡単のため、距離関数を従来から用いられてきた (1) の形式のまま議論してきたが、実験により、複数の共分散記述子を抽出して次の距離を用いることでパターン認識性能が向上することが分かった：

$$D_{\Phi}(\mathbf{x}_1, \mathbf{x}_2; \mathbf{W}) := \frac{1}{M} \sum_{m=1}^M \left\langle \mathbf{W}_m, (\Phi_m(\mathbf{x}_1) - \Phi_m(\mathbf{x}_2))(\Phi_m(\mathbf{x}_1) - \Phi_m(\mathbf{x}_2))^{\top} \right\rangle. \quad (7)$$

ただし、 Φ_m は第 m 特徴抽出器であり、距離関数のパラメータは M 個のマハラノビス行列 $\mathbf{W} := (\mathbf{W}_1, \dots, \mathbf{W}_M)$ で構成している。Theorem 1 は (7) を使用した場合にも容易に拡張できる [3]。

最適化アルゴリズム Theorem 1 に示したとおり、最適化問題 (3) はブレッグマン射影になるので、ダイクストラ法を使用することができる。確率的ダイクストラ法は反復法のひとつで、反復 $t-1$ において、解が $(\mathbf{W}_{t-1}, \boldsymbol{\xi}_{t-1})$ にいるとすると、反復 t では無作為に選んだ第 k 半空間の境界に解 $(\mathbf{W}_{t-1}, \boldsymbol{\xi}_{t-1})$ を射影する。その射影点は次の $\bar{\delta}$ の非線形一元方程式

$$\langle \mathbf{e}_k, \boldsymbol{\xi}_{t-1} \rangle + \bar{\delta} h_{k,k} = \frac{1}{M} \sum_{m=1}^M \left\langle \mathbf{A}_{m,k}, \left(\mathbf{W}_{t-1,m}^{-1} + \bar{\delta} \mathbf{A}_{m,k} \right)^{-1} \right\rangle \quad (8)$$

を解くと得られる [3]。ただし、

$\mathbf{A}_{m,k} := (\Phi_m(\mathbf{x}_{i_k}) - \Phi_m(\mathbf{x}_{j_k})) (\Phi_m(\mathbf{x}_{i_k}) - \Phi_m(\mathbf{x}_{j_k}))^{\top}$, $h_{k,k}$ は \mathbf{G}^{-1} の第 k 対角成分である。非線形方程式 (8) の解は閉形式では求まらないので、二分探索やニュートン法などを使って数値的に探さなくてはならない。ナイーブに探そうとすると、いくつかの $\bar{\delta}$ の値における (8) の両辺を評価することになる。数値解法内部で L 回両辺を評価すると、(8) の右辺は計算に $O(n^3)$ にかかる逆行列を M 個含んでいるため、この非線形方程式の解を見つけるのに $O(LMn^3)$ の計算量がかかる。本研究では、文献 [2] で用いたトリックを再利用することで、 $O(Mn^3)$ の計算量で非線形方程式を解くことができることを見出した [3]。

実験結果 実験結果は文献 [3] にて報告する。

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

参考文献

- [1] Davis J. V. et al.: Information-theoretic metric learning, *ICML*, ACM, pp. 209–216 (2007).
- [2] Matsuzawa, T., Relator, R., Sese, J. and Kato, T.: Stochastic Dykstra Algorithms for Metric Learning with Positive Definite Covariance Descriptors, *ECCV*, pp. 786–799 (2016).
- [3] Onuma, Y., Rivero, R. and Kato, T.: Threshold Auto-Tuning Metric Learning, <https://arxiv.org/abs/1801.02125>.