

確率密度微分の直接推定と機械学習への応用

奈良先端科学技術大学院大学・情報科学研究科・佐々木 博昭

Hiroaki Sasaki

Graduate School of Information Science
Nara Institute of Science and Technology
hsasaki@is.naist.jp

東京大学・大学院新領域創成科学研究科・杉山 将

Masashi Sugiyama

Graduate School of Frontier Sciences
The University of Tokyo
sugi@k.u-tokyo.ac.jp

概要

確率密度関数の微分を推定することは、機械学習の研究分野において、一般性のある研究課題である。確率密度関数の微分を推定する上で、単純なアプローチは、最初に、確率密度関数を推定し、次に、その微分を計算することであろう。しかし、この2段階推定は適切なアプローチではない。なぜならば、良い確率密度関数の推定結果が、必ずしも良い確率密度関数の微分推定結果をもたらすとは限らないからである。より適切なアプローチは、確率密度関数の推定を実行することなく、直接、確率密度関数の微分を推定することであろう。本稿の目的は、この直接推定のアプローチに沿ったいくつかの確率密度関数の微分推定法を解説することである。そして、それら微分推定法の機械学習問題への応用例についても紹介する。

1 はじめに

入力データの確率密度関数を推定することは、機械学習における問題を解く上で最も一般性のあるアプローチの1つである。その一方、確率密度関数ではなく、本質的には、そ

の微分を必要とする問題が数多く存在する。例えば、平均シフトクラスタリング [1, 2, 3] では、まず、推定された確率密度関数の勾配を用いて、データ点を確率密度関数のモード点（極大点）へ向けて更新する。そして、同じモード点に収束したデータ点に対して、同じクラスラベルを割り当てることでクラスタリングを行う。他の例として、最近傍法を用いたカルバック・ライブラー情報量推定法の推定バイアスは、確率密度関数のヘシアン行列に依存することが知られている [4, 5]。その他にも、カーネル密度推定における最適バンド幅パラメータ決定 [6] もしくは最適バンド幅行列決定 [7, 8]、非ガウス成分分析 [9, 10]、十分次元削減 [11, 12] などの問題が確率密度関数の微分推定を用いて解くことができる。さらに、文献 [13] では、確率密度関数の微分に関連した統計データ解析問題が紹介されている。したがって、確率密度関数の微分を推定することは、機械学習研究において、一般性のある研究課題という言うことができるであろう。

確率密度関数を推定する上で、最も単純なアプローチは、(1) 入力データの確率密度関数を推定し、(2) その微分を計算するという2段階推定法である。しかしながら、この2段階推定は、確率密度関数の微分を推定する上で、適切なアプローチではない。なぜならば、良い確率密度関数の推定結果が、必ずしも良い確率密度関数の微分推定結果をもたらすとは限らないからである。加えて、このアプローチでは、高階微分を計算する際に、低次微分推定結果に基づくため、より大きな推定誤差が生じると考えられる。

この問題に対処する上で、より適切なアプローチは、確率密度関数の推定を行うことなく、直接、その微分を推定することであろう。本稿の目的は、この直接推定のアプローチに基づいた確率密度関数の微分推定法とその機械学習問題への応用結果を解説・紹介することである [14, 5, 12, 10]。直接推定法の基本的な考え方は、2乗損失関数の下で確率密度関数の微分モデルを真の確率密度関数の微分に直接適合することである。この単純な定式化によって、微分が解析的に推定でき、モデル選択法も組み込まれているといった利点も合わせもつ。実際に、このアプローチによって、高精度な微分推定が可能となり、既存手法を上回る性能をもつ機械学習手法がいくつか構築されている。

本稿では、第2節で、確率密度関数の対数勾配の直接推定法とモード探索クラスタリングへの応用結果について示す。第3節では、条件付き確率密度関数の対数勾配の推定法と十分次元削減への応用について紹介する。その他の微分推定法や応用例は、第4節で簡潔に紹介・解説する。第5節は、本稿のまとめである。

2 対数密度勾配推定とモード探索クラスタリングへの応用

本節では、最初に平均シフトクラスタリングについて解説し、次に、確率密度関数の対数勾配の直接推定法 [15, 14] を紹介する。そして、その直接推定法をモード探索クラスタリングへ応用し、その有用性を数値実験によって示す。

2.1 平均シフトクラスタリング

平均シフトクラスタリング (mean shift clustering) [1, 2, 3] は、モード探索型のクラスタリング法であり、推定された確率密度関数のモード点に応じて、クラスタリングが実行される。平均シフトクラスタリングでは、最初に、推定された確率密度関数の勾配を用いた最急上昇法によって、データ点を近傍のモード点へと更新する。次に、同じモード点に収束したデータ点に同じクラスタラベルを割り当てる。推定された確率密度関数のモード点の数がクラスタ数に対応するため、平均シフトクラスタリングは、クラスタ数の事前入力が不要といった利点をもつ。平均シフトクラスタリングにおいて、重要なステップは、確率密度関数の勾配推定である。

確率密度関数の勾配を推定するために、平均シフトクラスタリングでは、最初に、カーネル関数 k を用いて、カーネル密度推定

$$\hat{p}_h(\mathbf{x}) = \frac{1}{nh^{d_x}} \sum_{i=1}^n k\left(\frac{\|\mathbf{x} - \mathbf{x}_i\|^2}{h^2}\right)$$

を実行することにより確率密度関数を推定する。次に、その勾配 $\nabla_{\mathbf{x}} \hat{p}_h(\mathbf{x})$ を計算する。この計算された確率密度関数の勾配から、データ点の更新に関して、次のように、不動点アルゴリズムを導出する。

$$\begin{aligned} \nabla_{\mathbf{x}} \hat{p}_h(\mathbf{x}) &= \frac{2}{nh^{d_x+2}} \sum_{i=1}^n [\mathbf{x}_i - \mathbf{x}] g\left(\frac{\|\mathbf{x} - \mathbf{x}_i\|^2}{h^2}\right) \\ &= \frac{2}{nh^{d_x+2}} \sum_{i=1}^n \mathbf{x}_i g\left(\frac{\|\mathbf{x} - \mathbf{x}_i\|^2}{h^2}\right) - \mathbf{x} \frac{2}{nh^{d_x+2}} \sum_{i=1}^n g\left(\frac{\|\mathbf{x} - \mathbf{x}_i\|^2}{h^2}\right) \\ &= \frac{2}{nh^{d_x+2}} \sum_{i=1}^n g\left(\frac{\|\mathbf{x} - \mathbf{x}_i\|^2}{h^2}\right) \left[\frac{\sum_{i=1}^n \mathbf{x}_i g\left(\frac{\|\mathbf{x} - \mathbf{x}_i\|^2}{h^2}\right)}{\sum_{i=1}^n g\left(\frac{\|\mathbf{x} - \mathbf{x}_i\|^2}{h^2}\right)} - \mathbf{x} \right] \end{aligned}$$

上式において $g(t) = -\frac{d}{dt} k(t)$ である。右辺の括弧内がゼロであるとき、 $\nabla_{\mathbf{x}} \hat{p}_h(\mathbf{x}) = \mathbf{0}$ で

あるため、次のような更新式が導出される。

$$\boldsymbol{x} \leftarrow \frac{\sum_{i=1}^n \boldsymbol{x}_i g\left(\frac{\|\boldsymbol{x}-\boldsymbol{x}_i\|^2}{h^2}\right)}{\sum_{i=1}^n g\left(\frac{\|\boldsymbol{x}-\boldsymbol{x}_i\|^2}{h^2}\right)} \quad (1)$$

単純な計算により、更新式 (1) は、 \boldsymbol{x} に依存したステップ幅をもつ最急上昇法と等価であることを確認できる。

平均シフトクラスタリングはクラスタ数の事前設定が不要という利点をもつが、その一方、入力データの次元が高いときに、クラスタリング性能精度が低くなることが知られている [3]。1つの理由として挙げられるのは、確率密度関数の勾配推定における2段階推定である。なぜならば、良い確率密度関数の推定結果が必ずしも良い確率密度関数の勾配推定結果をもたらすとは限らないからである。より適切なアプローチは、確率密度関数の推定を行うことなく、直接、確率密度関数の勾配を推定することであろう。以下では、この直接推定のアプローチに沿う確率密度関数の対数勾配の推定法 [15, 14] を紹介する。

2.2 最小2乗対数密度勾配

最初に、確率密度関数の対数勾配（以下、対数密度勾配と呼ぶ）推定の問題設定について述べる。確率密度関数 $p(\boldsymbol{x})$ より生成された n 個のデータ点 $\{\boldsymbol{x}_i = (x_i^{(1)}, x_i^{(2)}, \dots, x_i^{(d_x)})^\top\}_{i=1}^n$ が与えられているとしよう。ただし、 \top は転置を表す。ここでの目的は、 $\{\boldsymbol{x}_i\}_{i=1}^n$ から

$$\partial_j \log p(\boldsymbol{x}) = \frac{\partial_j p(\boldsymbol{x})}{p(\boldsymbol{x})}$$

を推定することである。上式において、 $\partial_j = \frac{\partial}{\partial x^{(j)}}$ である。

基本的なアプローチは、勾配モデル $\boldsymbol{r}(\boldsymbol{x}) = (r^{(1)}(\boldsymbol{x}), r^{(2)}(\boldsymbol{x}), \dots, r^{(d_x)}(\boldsymbol{x}))^\top$ を真の対数密度勾配に対して、2乗損失関数下で、適合することである。

$$J(\boldsymbol{r}^{(j)}) = \int \left\{ r^{(j)}(\boldsymbol{x}) - \frac{\partial_j p(\boldsymbol{x})}{p(\boldsymbol{x})} \right\}^2 p(\boldsymbol{x}) d\boldsymbol{x} \quad (2)$$

式 (2) を展開することで3つの項を得る。

$$J(\boldsymbol{r}^{(j)}) = \int \left\{ r^{(j)}(\boldsymbol{x}) \right\}^2 p(\boldsymbol{x}) d\boldsymbol{x} - 2 \int r^{(j)}(\boldsymbol{x}) \partial_j p(\boldsymbol{x}) d\boldsymbol{x} + \int \left\{ \frac{\partial_j p(\boldsymbol{x})}{p(\boldsymbol{x})} \right\}^2 p(\boldsymbol{x}) d\boldsymbol{x} \quad (3)$$

式 (3) における第1項はデータ点より推定可能、第3項はモデル $\boldsymbol{r}^{(j)}$ に依存しないため省略可能である。しかし、第2項は真の確率密度関数の微分 $\partial_j p(\boldsymbol{x})$ を含んでいるため、

推定は一見容易ではない。しかしながら、次のように部分積分を実行することで容易に推定可能な形式へと変換することができる。

$$\begin{aligned} \int r^{(j)}(\mathbf{x}) \partial_j p(\mathbf{x}) d\mathbf{x} &= \int \left[r^{(j)}(\mathbf{x}) p(\mathbf{x}) \right]_{x^{(j)}=-\infty}^{x^{(j)}=\infty} d\mathbf{x}_{\setminus x^{(j)}} - \int \left\{ \partial_j r^{(j)}(\mathbf{x}) \right\} p(\mathbf{x}) d\mathbf{x} \\ &= - \int \left\{ \partial_j r^{(j)}(\mathbf{x}) \right\} p(\mathbf{x}) d\mathbf{x} \end{aligned} \quad (4)$$

式 (4) における $d\mathbf{x}_{\setminus x^{(j)}}$ は、 $x^{(j)}$ を除いた変数についての積分を意味し、さらに、 $\lim_{|x^{(j)}| \rightarrow \infty} r^{(j)}(\mathbf{x}) p(\mathbf{x}) = 0$ を仮定している。同様の操作は、これまでもスタインの補題 [16] やスコアマッチング法 [17] を導くためにも用いられている。式 (3) と (4) より、次のように経験損失関数を導出することができる。

$$\hat{J}(r^{(j)}) = \frac{1}{n} \sum_{i=1}^n \left\{ r^{(j)}(\mathbf{x}_i) \right\}^2 + 2\partial_j r^{(j)}(\mathbf{x}_i) \quad (5)$$

ただし、式 (5) において、式 (3) の第 3 項は省略されている。

$r^{(j)}$ を推定するために、次のような線形モデルを考える。

$$r^{(j)}(\mathbf{x}) = \sum_{i=1}^b \theta_{ij} \psi_{ij}(\mathbf{x}) = \boldsymbol{\theta}_j^\top \boldsymbol{\psi}_j(\mathbf{x}) \quad (6)$$

式 (6) 内の θ_{ij} は推定されるパラメータ、 $\psi_{ij}(\mathbf{x})$ は基底関数、 b は基底関数の数を示し、本稿では $b = \min(100, n)$ に固定することにする。式 (6) を経験損失関数 (5) に代入し、 ℓ_2 正則化項を加えた後で、 $\boldsymbol{\theta}_j$ に関して、次のように解析解を計算できる。

$$\hat{\boldsymbol{\theta}}_j = \underset{\boldsymbol{\theta}_j}{\operatorname{argmin}} \left[\boldsymbol{\theta}_j^\top \mathbf{G}_j \boldsymbol{\theta}_j + 2\boldsymbol{\theta}_j^\top \mathbf{h}_j + \lambda_j \boldsymbol{\theta}_j^\top \boldsymbol{\theta}_j \right] = -(\mathbf{G}_j + \lambda_j \mathbf{I}_{d_x})^{-1} \mathbf{h}_j$$

上式における \mathbf{I}_{d_x} は d_x 行 d_x 列の単位行列、 λ_j は非負の正則化パラメータ、

$$\mathbf{G}_j = \frac{1}{n} \sum_{i=1}^n \boldsymbol{\psi}_j(\mathbf{x}_i) \boldsymbol{\psi}_j(\mathbf{x}_i)^\top, \quad \mathbf{h}_j = \frac{1}{n} \sum_{i=1}^n \partial_j \boldsymbol{\psi}_j(\mathbf{x}_i)$$

最終的に、対数密度勾配の推定結果は次のような形式で与えられる。

$$\hat{\mathbf{r}}(\mathbf{x}) = (\hat{r}^{(1)}(\mathbf{x}), \hat{r}^{(2)}(\mathbf{x}), \dots, \hat{r}^{(d_x)}(\mathbf{x}))^\top = (\hat{\boldsymbol{\theta}}_1^\top \boldsymbol{\psi}_1(\mathbf{x}), \hat{\boldsymbol{\theta}}_2^\top \boldsymbol{\psi}_2(\mathbf{x}), \dots, \hat{\boldsymbol{\theta}}_{d_x}^\top \boldsymbol{\psi}_{d_x}(\mathbf{x}))^\top$$

この手法を最小 2 乗対数密度勾配 (least-squares log-density gradient: LSLDG) と呼ぶ。

2.3 モード探索クラスタリングへの応用

ここでは、LSLDG をモード探索クラスタリングへ応用する。平均シフトクラスタリングと同様の不動点法に基づき、クラスタリングアルゴリズムを導出する。そのアルゴリズムを導出する上で、次のような基底関数を用いる。

$$\psi_{ij}(\mathbf{x}) = \partial_j \exp\left(-\frac{\|\mathbf{x} - \mathbf{c}_i\|^2}{2\sigma^2}\right) = \frac{c_i^{(j)} - x^{(j)}}{\sigma^2} \exp\left(-\frac{\|\mathbf{x} - \mathbf{c}_i\|^2}{2\sigma^2}\right) = \frac{c_i^{(j)} - x^{(j)}}{\sigma^2} \phi_i(\mathbf{x})$$

上式において、 $\sigma (> 0)$ はバンド幅パラメータ、 $\mathbf{c}_i = (c_i^{(1)}, c_i^{(2)}, \dots, c_i^{(d_x)})^\top$ はガウス関数の中心点であり、データ点 $\{\mathbf{x}_i\}_{i=1}^b$ からランダムに b 個選択することにする。この基底関数を用いて、

$$\begin{aligned} \hat{r}^{(j)}(\mathbf{x}) &= \sum_{i=1}^b \hat{\theta}_{ij} \psi_i(\mathbf{x}) = \frac{1}{\sigma^2} \left[\sum_{i=1}^b \hat{\theta}_{ij} c_i^{(j)} \phi_i(\mathbf{x}) - x^{(j)} \sum_{i=1}^b \hat{\theta}_{ij} \phi_i(\mathbf{x}) \right] \\ &= \frac{\sum_{i=1}^b \hat{\theta}_{ij} \phi_i(\mathbf{x})}{\sigma^2} \left[\frac{\sum_{i=1}^b \hat{\theta}_{ij} c_i^{(j)} \phi_i(\mathbf{x})}{\sum_{i=1}^b \hat{\theta}_{ij} \phi_i(\mathbf{x})} - x^{(j)} \right] \end{aligned} \quad (7)$$

式 (7) より、右辺の括弧内がゼロのとき、 $\hat{r}^{(j)}(\mathbf{x}) = 0$ となる。これにより、次のような更新式を得る。

$$x^{(j)} \leftarrow \frac{\sum_{i=1}^b \hat{\theta}_{ij} c_i^{(j)} \phi_i(\mathbf{x})}{\sum_{i=1}^b \hat{\theta}_{ij} \phi_i(\mathbf{x})} \quad (8)$$

本稿では、この手法を最小二乗対数密度勾配クラスタリング (LSLDG clustering: LSLDGC) と呼ぶ。平均シフトクラスタリングの更新式 (1) と比較すると、最大の違いは、重み $\hat{\theta}_{ij}$ である。したがって、LSLDGC は重み付き平均シフトクラスタリングと解釈できる。

2.4 数値実験

ここで、LSLDGC の性能を簡単な数値実験で示す。本数値実験において、入力データは、次のような混合ガウス分布より生成する。

$$p(\mathbf{x}) = \frac{1}{2} \mathcal{N}(\boldsymbol{\mu}_1, \mathbf{I}_{d_x}) + \frac{1}{2} \mathcal{N}(\boldsymbol{\mu}_2, \mathbf{I}_{d_x})$$

$\mathcal{N}(\boldsymbol{\mu}, \mathbf{C})$ は期待値 $\boldsymbol{\mu}$ と共分散行列 \mathbf{C} をもつガウス分布、 $\boldsymbol{\mu}_1 = (3, 3, 0, \dots, 0)^\top$ 、 $\boldsymbol{\mu}_2 = (-3, -3, 0, \dots, 0)^\top$ とした。LSLDG のパラメータ σ と λ_j は損失関数 J を評価基準と

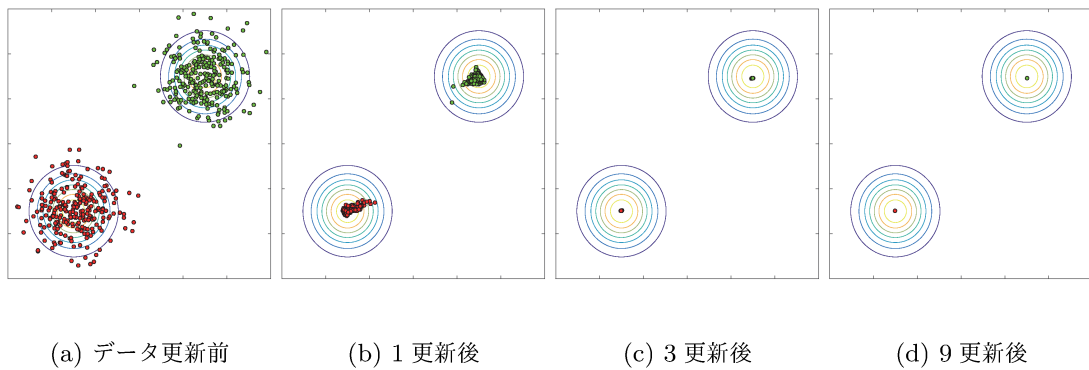


図1 LSLDGCによるモード探索過程におけるデータ点の分布。

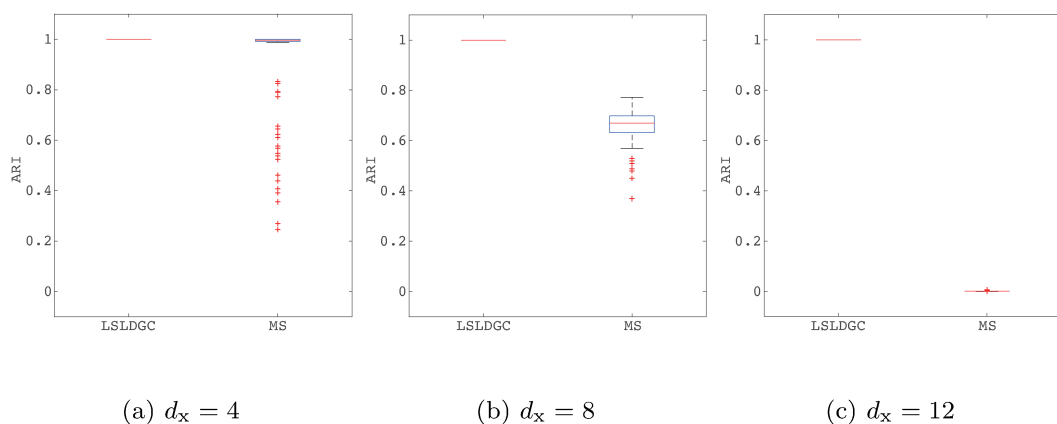


図2 LSLDGCと平均シフトクラスタリング (MS) の比較。箱ひげ図は、100 試行の数値実験結果である。ただし、 $n = 500$ 。

して、交差検定によって決定した。また、カーネル関数を $k(t) = \exp(-t)$ とし、平均シフトクラスタリングを同じデータに適用し、LSLDGC に対する比較を行った。平均シフトクラスタリングにおけるバンド幅パラメータ h もまた交差検定により決定した。クラスタリング性能を評価するために、調整ランド指数 (adjusted Rand index: ARI) [18] を用いた。ARI は 1 以下の値をとり、大きな ARI 値は良いクラスタリング結果を意味する。

更新式 (8) によるデータ点の更新過程を図 1 に示す。最初の数更新で、急速にデータ点が近傍のモードに向かい、その後、各モード点にデータ点が収束していく様子が分かる。このデータ点が更新されていく様子は、平均シフトクラスタリングと定性的に同じである。しかしながら、入力データの次元が高くなると、平均シフトクラスタリングの ARI

値は急速に減少していくことが分かる (図 2)。これに対して, LSLDGC はデータの次元が上がっても, 高い ARI 値を維持し続ける。これらの結果は, LSLDGC が, より高次元のデータに対して, 高精度なクラスタリングを実現することを意味している。

3 条件付き対数密度勾配推定と十分次元削減への応用

本節では, 最初に, 十分次元削減における条件付き確率密度関数の勾配に基づいたアプローチについて解説し, 次に, [12] で提案された条件付き確率密度関数の対数密度勾配推定法とその十分次元削減への応用例について紹介する。

3.1 十分次元削減と条件付き確率密度関数の勾配に基づいたアプローチ

十分次元削減 (sufficient dimension reduction) [11, 19, 20] とは, 教師付き次元削減の枠組みであり, 出力と入力に関するデータの組を (y, \boldsymbol{x}) で表すと, 次のような関係を満たす変換行列 $\boldsymbol{B} \in \mathbb{R}^{d_x \times d_z}$ を推定することを目的とする。

$$p(y|\boldsymbol{x}) = p(y|\boldsymbol{B}^\top \boldsymbol{x}) \quad (9)$$

$p(y|\boldsymbol{x})$ は入力データ \boldsymbol{x} が与えられたときの出力データ y の条件付き確率であり, $1 \leq d_z < d_x$, $\boldsymbol{B}^\top \boldsymbol{B} = \boldsymbol{I}_{d_z}$ とする。条件式 (9) は, 次元削減後の入力データ $\boldsymbol{B}^\top \boldsymbol{x}$ は, 出力データ y について, 次元削減前の入力データ \boldsymbol{x} と同じ情報をもつことを意味する。したがって, なるべく条件式 (9) を満たすような変換行列 \boldsymbol{B} を推定することで, 出力データ y の情報を失うことなく, 入力データ \boldsymbol{x} の次元を削減できることが期待できる。

これまでいくつかの十分次元削減法が提案されている。初期の手法は, スライス逆回帰 (sliced inverse regression) と呼ばれる手法である [11]。しかしながら, スライス逆回帰は, \boldsymbol{x} の周辺密度関数 $p(\boldsymbol{x})$ が楕円型であることを仮定しているため, 確率密度関数に関して, 強い制約をもつ。これに対して, 近年, 確率密度関数に関して, 制約の少ない手法が提案されている。例えば, カーネル次元削減 (kernel dimension reduction) [21] や最小 2 乗次元削減 (least-squares dimensionality reduction) [22] などがある。しかし, これら手法は, \boldsymbol{B} を推定する上で, 非凸最適化問題を解く必要があり, 一般的には局所最適解しか得られない。

この局所最適解を避けるアプローチとして, 条件付き確率密度関数の勾配に基づいたアプローチがある [23]。条件式 (9) の両辺の \boldsymbol{x} に関する勾配 $\nabla_{\boldsymbol{x}}$ を計算することで,

$$\nabla_{\boldsymbol{x}} p(y|\boldsymbol{x}) = \boldsymbol{B} \nabla_{\boldsymbol{B}^\top \boldsymbol{x}} p(y|\boldsymbol{B}^\top \boldsymbol{x}) \quad (10)$$

が得られる。式 (10) より、条件付き確率密度関数の勾配 $\nabla_{\mathbf{x}}p(y|\mathbf{x})$ は、 B の列空間に存在することが分かる。したがって、主成分分析と同様に、 B の列ベクトルは、

$$E\{(\nabla_{\mathbf{x}}p(y|\mathbf{x}))(\nabla_{\mathbf{x}}p(y|\mathbf{x}))^{\top}\}$$

の上位 d_x 個の固有値に対応した固有ベクトルとして得られる。このアプローチは、固有値分解に基づくため、大域最適解が保障される。ただし、このアプローチにおける困難な点は、 $\nabla_{\mathbf{x}}p(y|\mathbf{x})$ の推定である。

先行研究 [23] では、 $\nabla_{\mathbf{x}}p(y|\mathbf{x})$ を局所線形平滑化 (local linear smoother) [24] と呼ばれる手法で推定している。しかしながら、局所線形平滑化は、1 次のテイラー近似に基づいており、データ点が疎になる高次元データに対して、推定精度が落ちることが考えられる。また、パラメータ調整が困難なことや、データ点数が大きくなると計算効率が悪いといった問題点もある。そこで、[12] では、条件付き確率密度関数の対数勾配について、テイラー近似を使用することなく、モデル選択方法も含む推定手法を提案し、新たな十分次元削減法を構築している。以下では、その手法について解説する。

3.2 最小 2 乗条件付き対数密度勾配と十分次元削減への応用

最初に、条件付き確率密度関数の対数勾配 (以下、条件付き対数密度勾配と呼ぶ) 推定の問題設定について述べる。同時確率密度関数 $p(y, \mathbf{x})$ より生成された n 個の出入力データ点の組 $\{(y_i, \mathbf{x}_i)\}_{i=1}^n$ が与えられているとしよう。ここでの目的は、 $\{(y_i, \mathbf{x}_i)\}_{i=1}^n$ から

$$\partial_j \log p(y|\mathbf{x}) = \frac{\partial_j p(y|\mathbf{x})}{p(y|\mathbf{x})} = \frac{\partial_j p(y, \mathbf{x})p(\mathbf{x}) - p(y, \mathbf{x})\partial_j p(\mathbf{x})}{p(y, \mathbf{x})p(\mathbf{x})} = \frac{\partial_j p(y, \mathbf{x})}{p(y, \mathbf{x})} - \frac{\partial_j p(\mathbf{x})}{p(\mathbf{x})}$$

を推定することである。

LSLDG と同様に、基本的なアプローチは、勾配モデル $\mathbf{g}(y, \mathbf{x}) = (g^{(1)}(y, \mathbf{x}), \dots, g^{(d_x)}(y, \mathbf{x}))^{\top}$ を 2 乗損失関数の下で、真の条件付き対数密度勾配に直

接適合することである。

$$\begin{aligned}
J(g^{(j)}) &= \iint \left\{ g^{(j)}(y, \mathbf{x}) - \partial_j \log p(y|\mathbf{x}) \right\}^2 p(y, \mathbf{x}) dy d\mathbf{x} \\
&= \iint \left\{ g^{(j)}(y, \mathbf{x}) \right\}^2 p(y, \mathbf{x}) dy d\mathbf{x} - 2 \iint g^{(j)}(y, \mathbf{x}) \partial_j p(y, \mathbf{x}) dy d\mathbf{x} \\
&\quad + 2 \iint g^{(j)}(y, \mathbf{x}) \left\{ \frac{\partial_j p(\mathbf{x})}{p(\mathbf{x})} \right\} p(y, \mathbf{x}) dy d\mathbf{x} + \iint \left\{ \partial_j \log p(y|\mathbf{x}) \right\}^2 p(y, \mathbf{x}) dy d\mathbf{x} \\
&= \iint \left\{ g^{(j)}(y, \mathbf{x}) \right\}^2 p(y, \mathbf{x}) dy d\mathbf{x} + 2 \iint \partial_j g^{(j)}(y, \mathbf{x}) p(y, \mathbf{x}) dy d\mathbf{x} \\
&\quad + 2 \iint g^{(j)}(y, \mathbf{x}) \left\{ \frac{\partial_j p(\mathbf{x})}{p(\mathbf{x})} \right\} p(y, \mathbf{x}) dy d\mathbf{x} + \iint \left\{ \partial_j \log p(y|\mathbf{x}) \right\}^2 p(y, \mathbf{x}) dy d\mathbf{x}
\end{aligned} \tag{11}$$

最後の等式は、LSLDG と同様に、右辺第 2 項に対して部分積分を実行した。式 (11) の右辺における第 4 項は省略可能、第 1 項と第 2 項はデータ点から容易に推定可能であるが、第 3 項を推定することが容易ではない。なぜならば、被積分関数に、真の確率密度関数の対数微分 $\partial_j \log p(\mathbf{x}) = \frac{\partial_j p(\mathbf{x})}{p(\mathbf{x})}$ が含まれているからである。この問題に対処するために、ここでは、真の確率密度関数の対数微分 $\partial_j \log p(\mathbf{x})$ を第 2 節で紹介した LSLDG による推定結果 $\hat{r}^{(j)}(\mathbf{x})$ で置き換える。これにより、式 (11) の第 3 項は

$$\iint g^{(j)}(y, \mathbf{x}) \left\{ \frac{\partial_j p(\mathbf{x})}{p(\mathbf{x})} \right\} p(y, \mathbf{x}) dy d\mathbf{x} \approx \iint g^{(j)}(y, \mathbf{x}) \hat{r}^{(j)}(\mathbf{x}) p(y, \mathbf{x}) dy d\mathbf{x}$$

と近似されるため、結果として、右辺はデータ点より容易に推定可能である。以上より、LSLDG を用いた近似経験損失関数は、

$$\hat{J}(g^{(j)}) = \frac{1}{n} \sum_{i=1}^n \left\{ g^{(j)}(y_i, \mathbf{x}_i) \right\}^2 + 2 \left\{ \partial_j g^{(j)}(y_i, \mathbf{x}_i) + g^{(j)}(y_i, \mathbf{x}_i) \hat{r}^{(j)}(\mathbf{x}_i) \right\} \tag{12}$$

となる。

次に、LSLDG と同様に、次のような線形モデルを用いる。

$$g^{(j)}(y, \mathbf{x}) = \sum_{i=1}^b \theta_{ij} \underbrace{\exp \left(-\frac{\|\mathbf{x} - \mathbf{x}_i\|^2}{2(\sigma_x^{(j)})^2} - \frac{(y - y_i)^2}{2\sigma_y^2} \right)}_{\varphi_{ij}(y, \mathbf{x})} = \boldsymbol{\theta}_j^\top \boldsymbol{\varphi}_j(y, \mathbf{x})$$

$\sigma_x^{(j)}$ と σ_y はバンド幅パラメータである。線形モデルと l_2 正則化項を近似経験損失関数 (12) に挿入することで、経験近似損失関数は 2 次形式となり、解析的に解を計算で

入力： 出入カデータ $\{(y_i, \mathbf{x}_i)\}_{i=1}^n$.

ステップ 1 LSLCG を用いて，条件付き対数密度勾配 $\nabla_{\mathbf{x}} \log p(y|\mathbf{x})$ を推定.

ステップ 2 推定された勾配 $\hat{\mathbf{g}}(y, \mathbf{x})$ から， $\hat{\Lambda} = \frac{1}{n} \sum_{i=1}^n \hat{\mathbf{g}}(y_i, \mathbf{x}_i) \hat{\mathbf{g}}(y_i, \mathbf{x}_i)^\top$ を計算.

ステップ 3 $\hat{\Lambda}$ に対して固有値分解を実行. 上位 d_x 個の固有値に対応する固有ベクトルを $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_{d_x}$ とする.

出力： $\hat{\mathbf{B}} = (\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_{d_x})$.

図 3 LSGDR アルゴリズム.

きる.

$$\hat{\boldsymbol{\theta}}_j = \arg \min_{\boldsymbol{\theta}_j} [\boldsymbol{\theta}_j^\top \mathbf{G}_j \boldsymbol{\theta}_j + 2\boldsymbol{\theta}_j^\top \mathbf{h}_j + \lambda_j \boldsymbol{\theta}_j^\top \boldsymbol{\theta}_j] = -(\mathbf{G}_j + \lambda_j \mathbf{I})^{-1} \mathbf{h}_j$$

上の式において，

$$\mathbf{G}_j = \frac{1}{n} \sum_{i=1}^n \boldsymbol{\varphi}_j(y_i, \mathbf{x}_i) \boldsymbol{\varphi}_j(y_i, \mathbf{x}_i)^\top, \quad \mathbf{h}_j = \frac{1}{n} \sum_{i=1}^n \partial_j \boldsymbol{\varphi}_j(y_i, \mathbf{x}_i) + \hat{r}^{(j)}(\mathbf{x}_i) \boldsymbol{\varphi}_j(y_i, \mathbf{x}_i)$$

最終的に，条件付き対数密度勾配の推定結果は，

$$\hat{\mathbf{g}}(y, \mathbf{x}) = (\hat{\boldsymbol{\theta}}_1^\top \boldsymbol{\varphi}_1(y, \mathbf{x}), \hat{\boldsymbol{\theta}}_2^\top \boldsymbol{\varphi}_2(y, \mathbf{x}), \dots, \hat{\boldsymbol{\theta}}_{d_x}^\top \boldsymbol{\varphi}_{d_x}(y, \mathbf{x}))^\top$$

で与えられる. この手法を最小 2 乗条件付き対数密度勾配 (least-squares logarithmic conditional density gradients: LSLCG) と呼ぶ.

LSLCG を応用した十分次元削減法のアルゴリズムを図 3 に示す. この十分次元削減法を最小 2 乗勾配次元削減 (least-squares gradients for dimension reduction: LSGDR) と呼ぶ.

3.3 数値実験

ここで，LSGDR の推定性能に関する数値実験を行う. 比較対象として，カーネル次元削減 (KDR) [21] *1 と最小 2 乗次元削減 (LSDR) [22] *2 に加え，dOPG [23]*3 と呼

*1 <http://www.ism.ac.jp/~fukumizu/software.html>

*2 <http://www.ms.k.u-tokyo.ac.jp/software.html#LSDR>

*3 <http://www.stat.nus.edu.sg/~staxyc/dOPG.m>

ばれる条件付き確率密度関数の勾配に基づく手法を用いた。本数値実験において、出力データ y は、モデル

$$y = \sum_{i=1}^{d_z} |x^{(i)}| + 0.3\epsilon$$

によって生成した。ただし、 $x^{(i)}$ と ϵ は標準正規分布に従う確率変数である。このとき、 d_x 行 d_z 列の零行列を $O_{d_x \times d_z}$ と表すと、最適な変換行列 B は、

$$B = \begin{pmatrix} \mathbf{I}_{d_z} \\ O_{(d_x-d_z) \times d_z} \end{pmatrix}$$

となる。LSLCG におけるパラメータ $\sigma_x^{(j)}$ と λ_j は交差検定により決定、 σ_y は $|y_i - y_j|$ の i, j に関する中央値とした。推定誤差は、 $\|BB^T - \hat{B}\hat{B}^T\|_F$ によって評価した。ただし、 \hat{B} は各手法の推定結果、 $\|\cdot\|_F$ はフロベニウスノルムを表す。

部分空間の次元 $d_z = 2, 4, 6$ に対して、LSGDR と KDR の推定誤差が小さいことが分かる (図 4)。一方、dOPG と LSDR の推定誤差は、部分空間の次元が上がるにつれて、急激に大きくなることが分かる。LSDR は、部分空間の次元が上がるほど、 B の初期値の選択が難しく、望ましくない局所最適解が得られている可能性がある。dOPG は、部分空間の次元が高いとき、アルゴリズム内で用いている局所線形平滑化の精度が悪いことが、誤差の増大の原因であると考えられる。KDR の推定誤差は小さいものの、データ点数 n が大きくなるにつれて、計算時間が非常に大きくなることが分かる (図 5)。これは、LSGDR や dOPG とは異なり、KDR が反復最適化を実行するためである。dOPG の局所線形平滑化では、 n^2 個のパラメータを推定する必要があり、データ点数が大きくなると計算時間が LSGDR よりも大きくなる (図 5(b,c))。以上より、推定精度と計算時間の両面において、LSGDR は良い性能を示す手法であることが分かる。

4 確率密度微分に関連した他のトピック

本節では、他の確率密度関数の微分推定法や応用結果について紹介する。

4.1 確率密度関数の高階微分の直接推定法とその応用

LSLDG は確率密度関数の対数の 1 階微分を直接推定する手法であったが、この直接推定法は、確率密度関数の任意の階数の微分を直接推定する手法へと一般化できる [5, 8]。

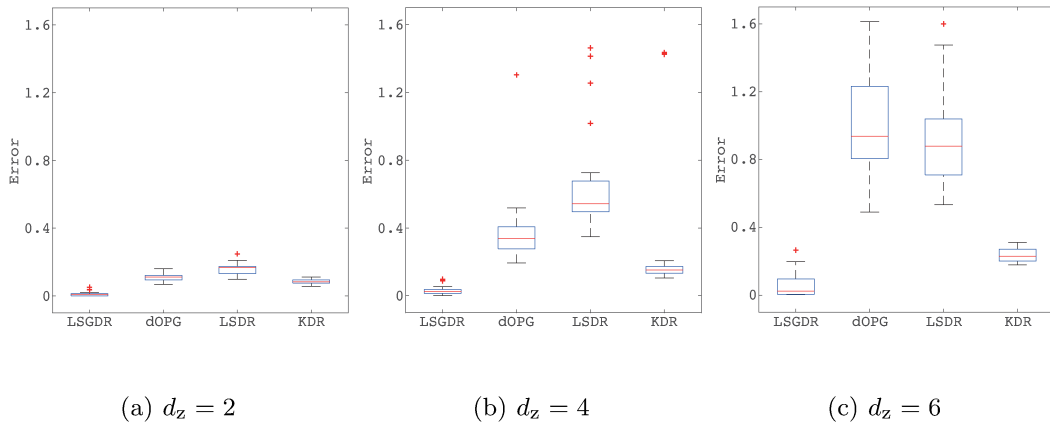


図4 部分空間の次元 d_z に対する各手法の推定誤差。ただし、 $(d_x, n) = (10, 500)$ である。

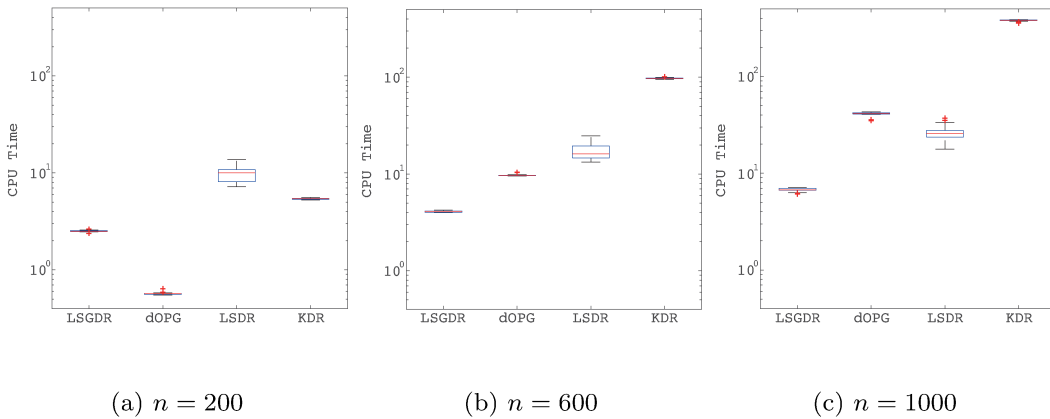


図5 データ点数 n に対する各手法の計算時間。縦軸は、対数目盛で表示された CPU 時間である。ただし、 $(d_x, d_z) = (10, 2)$ である。

例えば、確率密度関数の k 階微分を

$$p_{k,j}(\mathbf{x}) = \frac{\partial^k}{\partial(x^{(1)})^{j_1} \partial(x^{(2)})^{j_2} \dots \partial(x^{(d)})^{j_d}} p(\mathbf{x})$$

と定義し、次の 2 乗損失関数を k 階微分モデル $g_{k,j}$ に関して最小化することで推定できる。

$$J_{k,j}(g_{k,j}) = \int \{g_{k,j}(\mathbf{x}) - p_{k,j}(\mathbf{x})\}^2 d\mathbf{x} \quad (13)$$

ただし、 $\mathbf{j} = (j_1, j_2, \dots, j_d)^\top$, $j_i \in \{0, 1, \dots, d\}$, $j_1 + j_2 + \dots + j_d = k$ である。式 (13) は、真の確率密度関数の k 階を含んでいるため、一見、取り扱いが困難なように見えるが、

LSLDG における部分積分の操作を k 回繰り返すことで、次のような経験損失関数を導出できる。

$$\hat{J}_{k,j}(g_{k,j}) = \int g_{k,j}(\mathbf{x})^2 d\mathbf{x} - \frac{2(-1)^k}{n} \sum_{i=1}^n \frac{\partial^k}{\partial(x^{(1)})^{j_1} \partial(x^{(2)})^{j_2} \dots \partial(x^{(d)})^{j_d}} g_{k,j}(\mathbf{x}_i) \quad (14)$$

LSLDG と LSLCG と同様に、 $g_{k,j}(\mathbf{x})$ に線形モデルを導入することで、解析的に微分を推定できる。この手法は、積分 2 乗誤差密度微分推定法 (integrated squared error for density derivatives: ISED) と呼ばれている。

次に、確率密度関数の 2 階微分が必要な問題として、カルバック・ライブラー情報量推定法のバイアス削減を紹介する。確率密度関数 $p_1(\mathbf{x})$ と確率密度関数 $p_2(\mathbf{x})$ のカルバック・ライブラー情報量は、

$$\text{KL}(p_1 \| p_2) = \int p_1(\mathbf{x}) \log \frac{p_1(\mathbf{x})}{p_2(\mathbf{x})} d\mathbf{x}$$

で定義される。実用的なカルバック・ライブラー情報量推定法の 1 つとして、最近傍法を活用した推定法がある [25]。 $p_1(\mathbf{x})$ と $p_2(\mathbf{x})$ から生成されたデータ点をそれぞれ $\mathcal{X}_1 = \{\mathbf{x}_i\}_{i=1}^{n_1}$ 、 $\mathcal{X}_2 = \{\mathbf{x}_i\}_{i=n_1+1}^{n_1+n_2}$ とすると、

$$\widehat{\text{KL}}(p_1 \| p_2) = \frac{1}{n_1} \sum_{i=1}^{n_1} \log \frac{(n_1 - 1) \text{dist}_1(\mathbf{x}_i)^{d_x}}{n_2 \text{dist}_2(\mathbf{x}_i)^{d_x}}$$

によって、 $\text{KL}(p_1 \| p_2)$ が推定される。上式における $\text{dist}_1(\mathbf{x})$ と $\text{dist}_2(\mathbf{x})$ は \mathbf{x} からそれぞれ \mathcal{X}_1 と \mathcal{X}_2 内の最近傍データ点への距離を示す。しかしながら、 $\widehat{\text{KL}}(p_1 \| p_2)$ は、次のような量に比例したバイアスをもつが知られている [4]。

$$\frac{\text{tr}(\nabla\nabla p_1)}{((n_1 - 1)p_1)^{2/d_x} p_1} - \frac{\text{tr}(\nabla\nabla p_2)}{(n_2 p_2)^{2/d_x} p_2}$$

tr は対角和、 $\nabla\nabla p_1$ は、 p_1 のヘッセ行列を表す。このバイアスを削減するために、[5, 8] では、再スケールしたバイアス行列

$$\tilde{\mathbf{B}} = \frac{1}{(n_1 - 1)^{2/d}} \left(\frac{p_2}{p_1} \right)^{2/d_x + 1} \nabla\nabla p_1 - \frac{1}{n_2^{2/d_x}} \nabla\nabla p_2$$

を推定し、計量学習によって、 $\widehat{\text{KL}}(p_1 \| p_2)$ の推定バイアスを削減する方法が提案されている。 $\tilde{\mathbf{B}}$ を推定する際には、ヘシアン行列 $\nabla\nabla p_1, \nabla\nabla p_2$ は ISED によって、密度比 p_2/p_1 は [26] で提案された密度比推定法を用いている。バイアス削減を施したカルバック・ライブラー情報量推定法は、既存の手法を上回る推定性能をもつことが示されている [5, 8]。

4.2 LSLDG の拡張について

LSLDG の拡張として、マルチタスク学習を適用した拡張がある [27]。マルチタスク学習の目的は、複数のタスクが与えられたとき、タスク間の類似性を活用することで学習精度を向上させることである [28]。LSLDG への応用では、各次元 j に対する $\partial_j \log p(\mathbf{x})$ の推定を1つのタスクとみなすと、 $\partial_j \log p(\mathbf{x})$ は共通の $\log p(\mathbf{x})$ より計算されるため、タスク間に類似性が存在する。次の関係

$$\partial_j \log p(\mathbf{x}) \approx \hat{r}^{(j)}(\mathbf{x}) = \sum_{i=1}^b \hat{\theta}_{ij} \partial_j \phi_i(\mathbf{x}) = \partial_j \sum_{i=1}^b \hat{\theta}_{ij} \phi_i(\mathbf{x}) = \partial_j \hat{\theta}_j^\top \phi(\mathbf{x})$$

より、 $\hat{\theta}_j^\top \phi(\mathbf{x})$ が $\log p(\mathbf{x})$ の近似とみなせるため、 $\hat{\theta}_1 \approx \hat{\theta}_2 \approx \dots \approx \hat{\theta}_d$ となるように、 θ_j を推定することで、タスク間の類似性を活用できる。具体的には、正則化マルチタスク学習 [29, 30] を適用するために、次のような正則化項をさらに加える。

$$\gamma \sum_{j=1}^{d_x} \sum_{j'=1}^{d_x} \|\theta_j - \theta_{j'}\|^2 \quad (15)$$

γ は正則化パラメータである。正則化項 (15) より、 $\gamma \rightarrow \infty$ のとき、 $\theta_1 = \theta_2 = \dots = \theta_d$ となることが分かる。そして、正則化マルチタスク学習によって拡張された LSLDG は、数値実験より、データ点数 n が小さいときに有効であることが示された。

他の LSLDG の拡張では、データがリーマン多様体に属している場合を考える [31]。例えば、画像工学の分野で、動画内の物体の運動を取り扱う際に、物体の特徴点がある多様体に属することがしばしば仮定される [32]。LSLDG では、ユークリッド距離を用いているため、データが非ユークリッド空間内に存在する場合は、必ずしも優れた性能が得られるとは限らない。リーマン多様体上のデータ点を扱うために、[31] では、測地距離を用いて、LSLDG を拡張し、さらに、モード探索クラスリングへ応用している。グラスマン多様体に属するデータを用いた数値実験によって、そのクラスティング法が、LSLDGC を大きく上回る性能を示すことが確認されている。

4.3 非ガウス成分分析への応用

非ガウス成分分析 (non-Gaussian component analysis: NGCA) [9] は、教師無し線形次元削減の枠組みであり、射影後の入力データが非ガウス分布に従うような部分空間を見つけることを目的としている。非ガウス成分分析では、最初に、非ガウス分布に従う

入力： 入力データ $\{\mathbf{x}_i\}_{i=1}^n$.

ステップ 1 \mathbf{x}_i を中心化した後で，白色化.

ステップ 2 LSLDG を用いて，白色化後のデータ $\{\mathbf{y}_i (= \hat{\Sigma}^{-1/2} \mathbf{x}_i)\}_{i=1}^n$ から， $\nabla_{\mathbf{y}} \log p(\mathbf{y})$ を推定.

ステップ 3 推定された勾配 $\hat{\mathbf{r}}(\mathbf{y}_i)$ から， $\hat{\Gamma} = \frac{1}{n} \sum_{i=1}^n \{\hat{\mathbf{r}}(\mathbf{y}_i) + \mathbf{y}_i\} \{\hat{\mathbf{r}}(\mathbf{y}_i) + \mathbf{y}_i\}^T$ を計算.

ステップ 4 $\hat{\Gamma}$ に対して固有値分解を実行. 上位 d_s 個の固有値に対応する固有ベクトルより部分空間 $\hat{\mathcal{L}}$ を構築.

出力： $\hat{\mathcal{L}} = \hat{\Sigma}^{-1/2} \hat{\mathcal{L}}$.

図 6 LSNCA のアルゴリズム.

確率変数 $s^{(j)}$ を含むベクトル $\mathbf{s} = (s^{(1)}, s^{(2)}, \dots, s^{(d_s)})^T$ と d_x 行 d_s 列の行列 \mathbf{A} を用いて，入力データ \mathbf{x} がモデル

$$\mathbf{x} = \mathbf{A}\mathbf{s} + \mathbf{n}$$

によって生成されると仮定する. \mathbf{n} は平均 $\mathbf{0}$ ，共分散行列 \mathbf{C} のガウスノイズである. この仮定の下で，確率密度関数 $p(\mathbf{x})$ が次のようなセミパラメトリックモデルで表現されることが証明されている [9].

$$p(\mathbf{x}) = f_{\mathbf{x}}(\mathbf{W}^T \mathbf{x}) \mathcal{N}(\mathbf{0}, \mathbf{C}) \quad (16)$$

式 (16) における \mathbf{W} は d_x 行 d_s 列の行列， $f_{\mathbf{x}}$ は非負の関数である. 式 (16) では， \mathbf{W} , $f_{\mathbf{x}}$, \mathbf{C} を一意に決定できないが，次のような部分空間は同定可能である [33].

$$\mathcal{L} = \text{Ker}(\mathbf{W}^T)^\perp = \text{Range}(\mathbf{W}) \quad (17)$$

ここで， Ker と Range は零空間と値域を表し， $^\perp$ は直交補空間を表す. 上式における \mathcal{L} は非ガウス部分空間と呼ばれる部分空間である. 非ガウス成分分析における問題は， $\{\mathbf{x}_i\}_{i=1}^n$ から \mathcal{L} を推定することである.

非ガウス成分分析においても LSLDG を活用した計算効率の良いアルゴリズムを導出することができる [10]. 入力データを白色化することで，セミパラメトリックモデル (16) は次のような簡略化した形式となる [34].

$$p(\mathbf{y}) = f_{\mathbf{y}}(\mathbf{W}'^T \mathbf{y}) \mathcal{N}(\mathbf{0}, \mathbf{I}_{d_x}) \quad (18)$$

$\mathbf{y} = \Sigma^{-1/2}\mathbf{x}$ は白色化後のデータ, $\Sigma = E\{\mathbf{x}\mathbf{x}^\top\}$, $f_{\mathbf{y}}$ は非負の関数, \mathbf{W}' は $\mathbf{W}'^\top\mathbf{W}' = \mathbf{I}_{d_s}$ を満たす d_x 行 d_s 列の行列である. 元のセミパラメトリックモデル (16) との最大の違いは, \mathcal{N} 内の共分散行列が単位行列 \mathbf{I}_{d_x} へと変換されたことである. 簡略化されたセミパラメトリックモデル (18) の下で, 推定したい部分空間は, $\mathcal{L} = \text{Range}(\mathbf{W}) = \Sigma^{-1/2}\text{Range}(\mathbf{W}')$ となる. Σ は, $\hat{\Sigma} = \frac{1}{n}\sum_{i=1}^n \mathbf{x}_i\mathbf{x}_i^\top$ によって推定可能なので, $\text{Range}(\mathbf{W}')$ を何らかの方法で推定する必要がある.

$\text{Range}(\mathbf{W}')$ を推定するために, 簡略化されたセミパラメトリックモデル (18) から次の関係が得られる.

$$\nabla_{\mathbf{y}} [\log p(\mathbf{y}) - \log \mathcal{N}(\mathbf{0}, \mathbf{I}_{d_x})] = \mathbf{W}' \nabla_{\mathbf{W}'^\top \mathbf{y}} \log f_{\mathbf{y}}(\mathbf{W}'^\top \mathbf{y}) \quad (19)$$

式 (19) は, 左辺 $\nabla_{\mathbf{y}} [\log p(\mathbf{y}) - \log \mathcal{N}(\mathbf{0}, \mathbf{I}_{d_x})] = \nabla_{\mathbf{y}} \log p(\mathbf{y}) + \mathbf{y}$ が, $\text{Range}(\mathbf{W}')$ に存在することを意味する. したがって, $\boldsymbol{\eta} = \nabla_{\mathbf{y}} \log p(\mathbf{y}) + \mathbf{y}$ とすると, 主成分分析と同様に, $\text{Range}(\mathbf{W}')$ は, $E\{\boldsymbol{\eta}\boldsymbol{\eta}^\top\}$ の上位 d_s 個の固有値に対応した固有ベクトルが張る部分空間として推定することができる.

この考えに基づいた手法である最小 2 乗非ガウス成分分析 (least-squares NGCA: LSNGCA) [10] のアルゴリズムが, 図 6 に示されている. LSNGCA は, 数値実験によって, 既存の NGCA の手法と比較して, 推定精度と計算効率性が良いことも示されている [10]. しかしながら, 推定された共分散行列 $\hat{\Sigma}$ の条件数が大きいとき, LSNGCA における白色化ステップは大きな問題となる. この問題に対して, [35] では, 元のセミパラメトリックモデル (16) から出発し, LSNGCA を白色化ステップを含まない手法へ発展させ, 更に非ガウス部分空間の推定精度を向上させている.

5 まとめ

本稿では, 確率密度関数の微分を直接推定するいくつかの手法について解説・紹介した. これら手法は, 多くの機械学習問題に適用可能であるため, 今後のさらなる発展が期待できる. 提案手法の理論的な解析は興味深い研究課題であり, 今後の重要課題の 1 つである.

謝辞

本稿の研究において, 佐々木博昭は, JSPS 科研費 15H06103 の助成を受けた. 杉山将は, JST CREST の助成を受けた.

参考文献

- [1] K. Fukunaga and L. Hostetler. The estimation of the gradient of a density function, with applications in pattern recognition. *IEEE Transactions on Information Theory*, 21(1):32–40, 1975.
- [2] Y. Cheng. Mean shift, mode seeking, and clustering. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 17(8):790–799, 1995.
- [3] D. Comaniciu and P. Meer. Mean shift: A robust approach toward feature space analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 24(5):603–619, 2002.
- [4] Y. K. Noh, M. Sugiyama, S. Liu, M. C. du Plessis, F. C. Park, and D. D. Lee. Bias reduction and metric learning for nearest-neighbor estimation of Kullback-Leibler divergence. In *Proceedings of the 17th International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 669–677, 2014.
- [5] H. Sasaki, Y. K. Noh, and M. Sugiyama. Direct density-derivative estimation and its application in KL-divergence approximation. In *Proceedings of the 18th International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 809–818, 2015.
- [6] B.W. Silverman. *Density Estimation for Statistics and Data Analysis*. CRC press, 1986.
- [7] M. P. Wand and M. C. Jones. Multivariate plug-in bandwidth selection. *Computational Statistics*, 9(2):97–116, 1994.
- [8] S. Sasaki, Y. K. Noh, G. Niu, and M. Sugiyama. Direct density derivative estimation. *Neural Computation*. to appear.
- [9] G. Blanchard, M. Kawanabe, M. Sugiyama, V. Spokoiny, and K.R. Müller. In search of non-Gaussian components of a high-dimensional distribution. *Journal of Machine Learning Research*, 7:247–282, 2006.
- [10] H. Sasaki, G. Niu, and M. Sugiyama. Non-Gaussian component analysis with log-density gradient estimation. In *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2016. to appear.
- [11] K.C. Li. Sliced inverse regression for dimension reduction. *Journal of the American Statistical Association*, 86(414):316–327, 1991.

- [12] H. Sasaki, V. Tangkaratt, and M. Sugiyama. Sufficient dimension reduction via direct estimation of the gradients of logarithmic conditional densities. In *Proceedings of the 7th Asian Conference on Machine Learning (ACML)*, volume 45, pages 33–48, 2015.
- [13] R.S. Singh. Applications of estimators of a density and its derivatives to certain statistical problems. *Journal of the Royal Statistical Society. Series B*, 39(3):357–363, 1977.
- [14] H. Sasaki, A. Hyvärinen, and M. Sugiyama. Clustering via mode seeking by direct estimation of the gradient of a log-density. In *Machine Learning and Knowledge Discovery in Databases Part III- European Conference, ECML/PKDD 2014*, volume 8726, pages 19–34, 2014.
- [15] D. D. Cox. A penalty method for nonparametric estimation of the logarithmic derivative of a density function. *Annals of the Institute of Statistical Mathematics*, 37(1):271–288, 1985.
- [16] C.M. Stein. Estimation of the mean of a multivariate normal distribution. *The Annals of Statistics*, pages 1135–1151, 1981.
- [17] A. Hyvärinen. Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6:695–709, 2005.
- [18] L. Hubert and P. Arabie. Comparing partitions. *Journal of Classification*, 2(1):193–218, 1985.
- [19] R. D. Cook. *Regression Graphics: Ideas for Studying Regressions Through Graphics*. John Wiley & Sons, 1998.
- [20] F. Chiaromonte, R. D. Cook, and B. Li. Sufficient dimension reduction in regressions with categorical predictors. *The Annals of Statistics*, pages 475–497, 2002.
- [21] K. Fukumizu, F. R. Bach, and M. I. Jordan. Dimensionality reduction for supervised learning with reproducing kernel hilbert spaces. *Journal of Machine Learning Research*, 5:73–99, 2004.
- [22] T. Suzuki and M. Sugiyama. Sufficient dimension reduction via squared-loss mutual information estimation. *Neural Computation*, 25(3):725–758, 2013.
- [23] Y. Xia. A constructive approach to the estimation of dimension reduction directions. *The Annals of Statistics*, 35(6):2654–2690, 2007.
- [24] Jianqing Fan, Qiwei Yao, and Howell Tong. Estimation of conditional densities

- and sensitivity measures in nonlinear dynamical systems. *Biometrika*, 83(1):189–206, 1996.
- [25] Q. Wang, S. R. Kulkarni, and S. Verdu. A nearest-neighbor approach to estimating divergence between continuous random vectors. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 55(5):2392–2405, 2006.
- [26] T. Kanamori, S. Hido, and M. Sugiyama. A least-squares approach to direct importance estimation. *Journal of Machine Learning Research*, 10:1391–1445, 2009.
- [27] I. Yamane, H. Sasaki, and M. Sugiyama. Regularized multi-task learning for multi-dimensional log-density gradient estimation. *Neural Computation*. to appear.
- [28] R. Caruana. Multitask learning. *Machine learning*, 28(1):41–75, 1997.
- [29] T. Evgeniou and M. Pontil. Regularized multi-task learning. In *Proceedings of the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 109–117. ACM, 2004.
- [30] A. Argyriou, T. Evgeniou, and M. Pontil. Multi-task feature learning. In *Advances in Neural Information Processing Systems (NIPS)*, pages 41–48, Cambridge, MA, 2007. MIT Press.
- [31] M. Ashizawa, H. Sasaki, Sakai T., and M. Sugiyama. Least-squares log-density gradient clustering for riemannian manifolds. *IEICE Tech. Rep.*, 115(511):17–24, 2016.
- [32] V.M. Govindu. Lie-algebraic averaging for globally consistent motion estimation. In *IEEE Computer Society Conference on Computer Vision and Pattern Recognition 2004. CVPR 2004.*, volume 1, pages 684–691. IEEE, 2004.
- [33] F.J. Theis and M. Kawanabe. Uniqueness of non-Gaussian subspace analysis. In *Proceedings of the sixth international conference on Independent Component Analysis and Blind Signal Separation*, volume 3889, pages 917–925. 2006.
- [34] M. Sugiyama, M. Kawanabe, G. Blanchard, and K.R. Müller. Approximating the best linear unbiased estimator of non-Gaussian signals with Gaussian noise. *IEICE transactions on information and systems*, 91(5):1577–1580, 2008.
- [35] H. Shiino, H. Sasaki, G. Niu, and M. Sugiyama. Whitening-free least-squares non-Gaussian component analysis. *arXiv preprint arXiv:1603.01029*, 2016.