

密度比推定によるビッグデータ解析

Big Data Analysis by Density Ratio Estimation

杉山 将

東京工業大学 計算工学専攻

sugi@cs.titech.ac.jp

<http://sugiyama-www.cs.titech.ac.jp/~sugi>

概要

統計的データ解析のあらゆるタスクは、データの背後に潜む確率分布を推定することにより解決できる。しかし、確率分布の推定は一般に困難であることが知られているため、これを回避しつつ所望のデータ解析を行うことが望ましい。本稿では、確率分布でなく確率密度の比を通して様々なデータ解析を行う密度比推定の枠組みを紹介する。この枠組には、非定常環境適応学習、異常値検出、二標本検定、独立性検定、特徴選択、独立成分分析、条件付き確率密度推定、確率的パターン認識など多くの重要なデータ解析パラダイムが含まれる。

キーワード

機械学習, 確率分布, 密度比, 変化検知

1 はじめに

ビッグデータ時代には、データサイエンティストは多種多様なデータに対して、様々なデータ解析を迅速に行う必要がある。その際、データ解析タスクごとにアルゴリズムの開発・実装を行うのは効率が悪い。このような背景のもと、様々なデータ解析タスクを統一的に解決できる密度比推定 [1] とよばれるアプローチが近年提案された。

密度比推定の枠組みには、非定常環境適応学習、異常値検出、二標本検定、独立性検定、特徴選択、独立成分分析、条件付き確率密度推定、確率的パターン認識など、多くの重要なデータ解析パラダイムが含まれる。そして、これらの多様なデータ解析タスクを、確率密度比関数の推定を通して統一的かつ精度良く解決する。

本稿では、まず第2節で様々な密度比推定の方法を概説し、第3節で密度比推定により解決できる様々な機械学習タスクの例を紹介する。そして、第4節では、密度比推定に基づく教師なし変化検知手法を詳しく説明し、最後に第5節で今後の展望を述べる。

2 密度比推定

確率密度 $p_{\text{nu}}^*(\mathbf{x})$ を持つ確率分布に独立に従う標本 $\{\mathbf{x}_i^{\text{nu}}\}_{i=1}^{n_{\text{nu}}}$ と、確率密度 $p_{\text{de}}^*(\mathbf{x})$ を持つ確率分布に独立に従う標本 $\{\mathbf{x}_j^{\text{de}}\}_{j=1}^{n_{\text{de}}}$ から、密度比

$$r^*(\mathbf{x}) = \frac{p_{\text{nu}}^*(\mathbf{x})}{p_{\text{de}}^*(\mathbf{x})}$$

を推定する問題を考える．‘nu’ と ‘de’ は，分子 (numerator) と分母 (denominator) の頭文字である．

標本 $\{\mathbf{x}_i^{\text{nu}}\}_{i=1}^{n_{\text{nu}}}$ から確率密度 $p_{\text{nu}}^*(\mathbf{x})$ を，標本 $\{\mathbf{x}_j^{\text{de}}\}_{j=1}^{n_{\text{de}}}$ から確率密度 $p_{\text{de}}^*(\mathbf{x})$ をそれぞれ推定すれば，密度比 $r^*(\mathbf{x}) = p_{\text{nu}}^*(\mathbf{x})/p_{\text{de}}^*(\mathbf{x})$ を推定できる．しかし，推定した確率密度の比を取るにより，確率密度の推定誤差が増幅されてしまう恐れがあるため，このような二段階の方式は推定が精度が良くない (図1)．そこで本節では，各確率密度 $p_{\text{nu}}^*(\mathbf{x})$ や $p_{\text{de}}^*(\mathbf{x})$ の推定を経由せずに，密度比 $r^*(\mathbf{x}) = p_{\text{nu}}^*(\mathbf{x})/p_{\text{de}}^*(\mathbf{x})$ を直接推定できる手法を紹介する [1]．

2.1 確率的分類法

確率的分類法では， $p_{\text{nu}}^*(\mathbf{x})$ と $p_{\text{de}}^*(\mathbf{x})$ から生成された標本に，ラベル $y = \text{‘nu’}$ と ‘de’ をそれぞれ割り当てる．このとき， $p_{\text{nu}}^*(\mathbf{x})$ と $p_{\text{de}}^*(\mathbf{x})$ を

$$\begin{aligned} p_{\text{nu}}^*(\mathbf{x}) &= p^*(\mathbf{x}|y = \text{‘nu’}) \\ p_{\text{de}}^*(\mathbf{x}) &= p^*(\mathbf{x}|y = \text{‘de’}) \end{aligned}$$

と表せるため，ベイズの定理により，密度比を

$$r^*(\mathbf{x}) = \frac{p^*(y = \text{‘de’}) p^*(y = \text{‘nu’}|\mathbf{x})}{p^*(y = \text{‘nu’}) p^*(y = \text{‘de’}|\mathbf{x})}$$

と表現できる．ここで，ラベルの事前確率 $p^*(y)$ の比を標本数の比で近似し，ラベルの事後確率 $p^*(y|\mathbf{x})$ を $\{\mathbf{x}_i^{\text{nu}}\}_{i=1}^{n_{\text{nu}}}$ と $\{\mathbf{x}_j^{\text{de}}\}_{j=1}^{n_{\text{de}}}$ に対する確率的分類器 $\hat{p}(y|\mathbf{x})$ (例えば，ロジスティック回帰や最小二乗確率的分類により求める) で近似すれば，密度比の近似 $\hat{r}(\mathbf{x})$ を次式で求められる．

$$\hat{r}(\mathbf{x}) = \frac{n_{\text{de}} \hat{p}(y = \text{‘nu’}|\mathbf{x})}{n_{\text{nu}} \hat{p}(y = \text{‘de’}|\mathbf{x})}$$

2.2 積率適合法

積率適合法では，密度比のモデル $r(\mathbf{x})$ を用いて， $r(\mathbf{x})p_{\text{de}}^*(\mathbf{x})$ の積率を $p_{\text{nu}}^*(\mathbf{x})$ の積率に最小二乗適合させる．例えば一次の積率 (すなわち期待値) を適合させる場合は，次式を解く．

$$\min_r \left\| \mathbb{E}_{p_{\text{de}}^*}[\mathbf{x}r(\mathbf{x})] - \mathbb{E}_{p_{\text{nu}}^*}[\mathbf{x}] \right\|^2$$

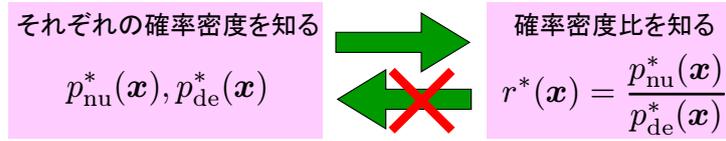


図 1: 密度比推定. 分子と分母の密度 $p_{\text{nu}}^*(\mathbf{x})$, $p_{\text{de}}^*(\mathbf{x})$ がわかればそれらの比 $r^*(\mathbf{x})$ もわかるが, 密度比 $r^*(\mathbf{x})$ がわかったとしてもそれぞれの密度はわからない. 従って, 分子と分母の密度を個別に推定するよりも, 密度比を直接推定の方が易しいと考えられる.

ただし, $\|\cdot\|$ はユークリッドノルム, \mathbb{E} は期待値を表す. 真の密度比を正しく求めるためには全ての次数の積率を適合させる必要がある. ガウスカーネルのような普遍再生核関数 $K(\mathbf{x}, \mathbf{x}')$ を用いれば, これを効率良く実現できる.

$$\min_r \left\| \mathbb{E}_{p_{\text{de}}^*} [K(\mathbf{x}, \cdot) r(\mathbf{x})] - \mathbb{E}_{p_{\text{nu}}^*} [K(\mathbf{x}, \cdot)] \right\|_{\mathcal{H}}^2$$

ただし, $\|\cdot\|_{\mathcal{H}}$ は $K(\mathbf{x}, \mathbf{x}')$ が属するヒルベルト空間のノルムを表す. 実際には, 期待値を標本平均で近似した規準を最小化することにより解を求める.

2.3 密度適合法

密度適合法では, 一般化カルバック距離のもとで $r(\mathbf{x})p_{\text{de}}^*(\mathbf{x})$ を $p_{\text{nu}}^*(\mathbf{x})$ に適合させる.

$$\min_r \mathbb{E}_{p_{\text{nu}}^*} \left[\log \frac{p_{\text{nu}}^*(\mathbf{x})}{r(\mathbf{x})p_{\text{de}}^*(\mathbf{x})} \right] + \mathbb{E}_{p_{\text{de}}^*} [r(\mathbf{x})]$$

ただし, 実際の推定には期待値を標本平均で近似した規準を用いる. $r(\mathbf{x})$ として, 線形モデル, 対数線形モデル, 混合モデルを用いた手法が提案されている.

2.4 密度比適合法

密度比適合法では, 密度比モデル $r(\mathbf{x})$ を真の密度比 $r^*(\mathbf{x})$ に最小二乗適合させる.

$$\min_r \mathbb{E}_{p_{\text{de}}^*} [(r(\mathbf{x}) - r^*(\mathbf{x}))^2]$$

ただし, 実際の推定には期待値を標本平均で近似した規準を用いる. $r(\mathbf{x})$ として線形モデルを用いれば, 密度比適合法の解は解析的に求められる.

2.5 統一的枠組み

上記の最小二乗密度比適合法を一般化し，ブregマン距離のもとで $r(\mathbf{x})$ を $r^*(\mathbf{x})$ に適合させる．

$$\min_r \mathbb{E}_{p_{\text{de}}^*} [f(r^*(\mathbf{x})) - f(r(\mathbf{x})) - f'(r(\mathbf{x}))(r^*(\mathbf{x}) - r(\mathbf{x}))]$$

ただし， $f(t)$ は微分可能な強凸関数であり， $f'(t)$ はその微分を表す． $f(t)$ を変えることにより，様々な密度比推定法が表現できる．

- ロジスティック回帰： $t \log t - (1+t) \log(1+t)$
- 再生核積率適合： $(t-1)^2/2$
- カルバック密度適合： $t \log t - t$
- 最小二乗密度比適合： $(t-1)^2/2$
- ロバスト密度比適合： $(t^{1+\alpha} - t)/\alpha$, ($\alpha > 0$)

2.6 次元削減付き密度比推定

\mathbf{x} を線形射影により $\begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix}$ と分解したときに， \mathbf{v} 成分が $p_{\text{nu}}^*(\mathbf{x})$ と $p_{\text{de}}^*(\mathbf{x})$ で共通，すなわち，ある共通の $p^*(\mathbf{v}|\mathbf{u})$ を用いて $p_{\text{nu}}^*(\mathbf{x})$ と $p_{\text{de}}^*(\mathbf{x})$ が

$$p_{\text{nu}}^*(\mathbf{x}) = p^*(\mathbf{v}|\mathbf{u})p_{\text{nu}}^*(\mathbf{u}), \quad p_{\text{de}}^*(\mathbf{x}) = p^*(\mathbf{v}|\mathbf{u})p_{\text{de}}^*(\mathbf{u})$$

と表現できるならば，密度比 $r^*(\mathbf{x})$ を $p_{\text{nu}}^*(\mathbf{u})/p_{\text{de}}^*(\mathbf{u})$ と簡略化できる．従って， \mathbf{u} が属する部分空間（異分布部分空間とよぶ）を特定すれば，高次元の密度比推定問題を低次元の問題に還元できる．異分布部分空間の探索は，教師付き次元削減手法により $\{\mathbf{x}_i^{\text{nu}}\}_{i=1}^{n_{\text{nu}}}$ と $\{\mathbf{x}_j^{\text{de}}\}_{j=1}^{n_{\text{de}}}$ を最も良く分離する部分空間を求める，あるいは， $p_{\text{nu}}^*(\mathbf{u})$ から $p_{\text{de}}^*(\mathbf{u})$ へのピアソン距離

$$\mathbb{E}_{p_{\text{de}}^*} [(p_{\text{nu}}^*(\mathbf{u})/p_{\text{de}}^*(\mathbf{u}) - 1)^2]$$

を最大にする部分空間を求めることにより行う．

3 密度比に基づく機械学習

本節では，密度比推定により解決できる機械学習タスクの例を示す．

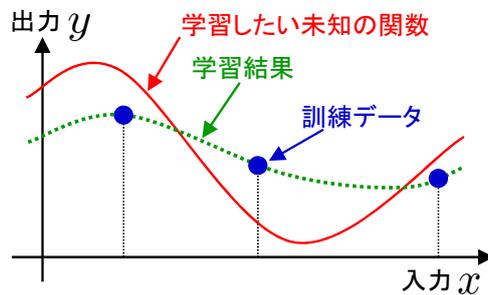


図 2: 教師付き学習. 入力と出力が対になった訓練データから, その背後に潜んでいる入出力関係を推定する.

3.1 重点標本化

入力 \mathbf{x} から出力 y への変換規則を学習する教師付き学習 (図 2) において, 訓練標本とテスト標本の入力分布が $p_{\text{tr}}^*(\mathbf{x})$ から $p_{\text{te}}^*(\mathbf{x})$ に変化するが, 入出力関係 $p^*(y|\mathbf{x})$ は変化しない状況を共変量シフトとよぶ [2]. 共変量シフト下では最尤推定などの学習法はバイアスを持つが, このバイアスは損失関数を重要度 $p_{\text{te}}^*(\mathbf{x})/p_{\text{tr}}^*(\mathbf{x})$ に従って重み付けすることにより打ち消される. すなわち, 損失関数 $\ell(\mathbf{x})$ の $p_{\text{te}}^*(\mathbf{x})$ に関する期待値は, 損失関数の $p_{\text{tr}}^*(\mathbf{x})$ に関する重要度重み付き期待値により計算できる (図 3).

$$\begin{aligned} \mathbb{E}_{p_{\text{te}}^*}[\ell(\mathbf{x})] &= \int \ell(\mathbf{x})p_{\text{te}}^*(\mathbf{x})d\mathbf{x} \\ &= \int \ell(\mathbf{x})\frac{p_{\text{te}}^*(\mathbf{x})}{p_{\text{tr}}^*(\mathbf{x})}p_{\text{tr}}^*(\mathbf{x})d\mathbf{x} = \mathbb{E}_{p_{\text{tr}}^*} \left[\ell(\mathbf{x})\frac{p_{\text{te}}^*(\mathbf{x})}{p_{\text{tr}}^*(\mathbf{x})} \right] \end{aligned}$$

交差確認などのモデル選択法も共変量シフト下では不偏性を失うが, 同様に重要度重み付けを行うことにより不偏性が回復できる.

3.2 確率分布比較

正常標本集合に基づいて, 評価標本集合に含まれる異常値を検出する問題を考える. これら二つの標本集合に対する密度比を考えれば, 正常値に対する密度比の値は 1 に近く, 異常値に対する密度比の値は 1 から大きく離れる. 従って, 密度比の値を評価基準とすることにより異常値を検出できる (図 4).

また, 密度比 $r^*(\mathbf{x}) = p_{\text{nu}}^*(\mathbf{x})/p_{\text{de}}^*(\mathbf{x})$ の推定量 $\hat{r}(\mathbf{x})$ を用いることにより, 分布 $p_{\text{nu}}^*(\mathbf{x}), p_{\text{de}}^*(\mathbf{x})$ 間の距離を精度良く推定できる [3].

- カルバック距離: $\frac{1}{n_{\text{nu}}} \sum_{i=1}^{n_{\text{nu}}} \log \hat{r}(\mathbf{x}_i^{\text{nu}})$
- ピアソン距離: $\frac{2}{n_{\text{nu}}} \sum_{i=1}^{n_{\text{nu}}} \hat{r}(\mathbf{x}_i^{\text{nu}}) - \frac{1}{n_{\text{de}}} \sum_{j=1}^{n_{\text{de}}} \hat{r}(\mathbf{x}_j^{\text{de}})^2 - 1$

これらの距離推定量を用いれば, 並べ替え検定により二つの分布の同一性を検定できる.

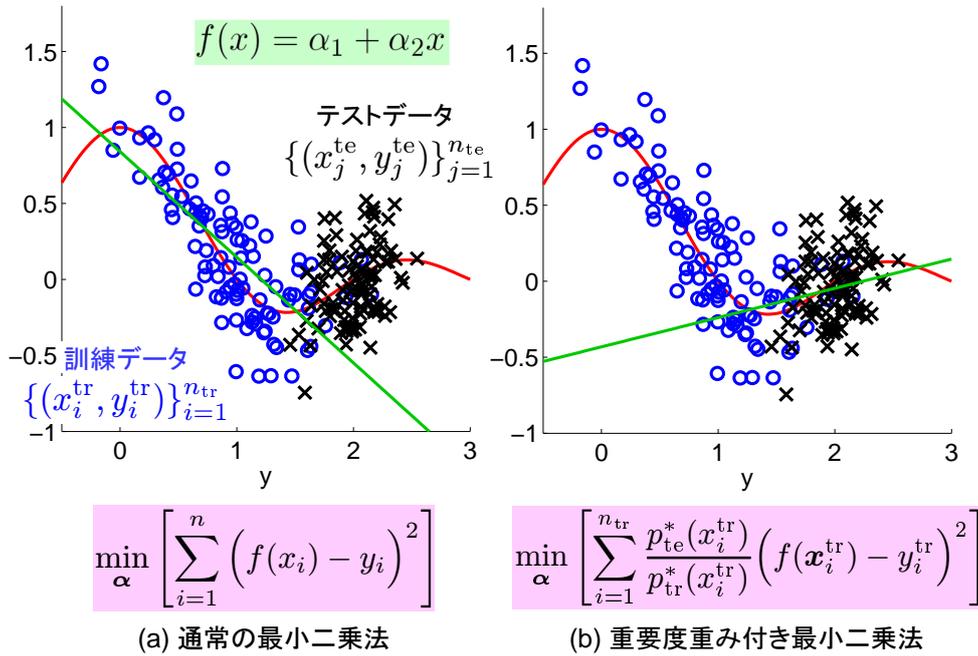


図 3: 重要度重み付き最小二乗法による共変量シフト適応. (a) 通常の最小二乗法ではモデルを訓練データに適合させるため, テストデータが訓練データと異なる分布に従う場合はテストデータをうまく予想できない. (b) テストデータに近い訓練データに強い重みをつけることにより, テスト出力に適合させる.

3.3 相互情報量推定

確率密度 $p_{x,y}^*(\mathbf{x}, \mathbf{y})$ を持つ分布に独立に従う n 個の標本 $\{(\mathbf{x}_k, \mathbf{y}_k)\}_{k=1}^n$ から, \mathbf{x} と \mathbf{y} の相互情報量

$$\mathbb{E}_{p_{x,y}^*} \left[\log \frac{p_{x,y}^*(\mathbf{x}, \mathbf{y})}{p_x^*(\mathbf{x})p_y^*(\mathbf{y})} \right]$$

を推定する問題を考える. ただし, $p_x^*(\mathbf{x})$ と $p_y^*(\mathbf{y})$ は \mathbf{x} と \mathbf{y} の周辺密度である. $\{(\mathbf{x}_k, \mathbf{y}_k)\}_{k=1}^n$ を分子の確率分布からの標本とみなし, $\{(\mathbf{x}_k, \mathbf{y}_{k'})\}_{k,k'=1}^n$ を分母の確率分布からの標本とみなせば, 密度比推定により相互情報量が推定できる. 同様に, 二乗損失版の相互情報量

$$\mathbb{E}_{p_x^*(\mathbf{x})} \mathbb{E}_{p_y^*(\mathbf{y})} \left[\left(\frac{p_{x,y}^*(\mathbf{x}, \mathbf{y})}{p_x^*(\mathbf{x})p_y^*(\mathbf{y})} - 1 \right)^2 \right]$$

も推定できる. 相互情報量は確率変数間の独立性を表す指標であり, その推定量は, 独立性検定, 特徴選択, 特徴抽出, クラスタリング, 独立成分分析, オブジェクト適合, 因果推定など, 様々な機械学習タスクに応用できる [4] (図5).

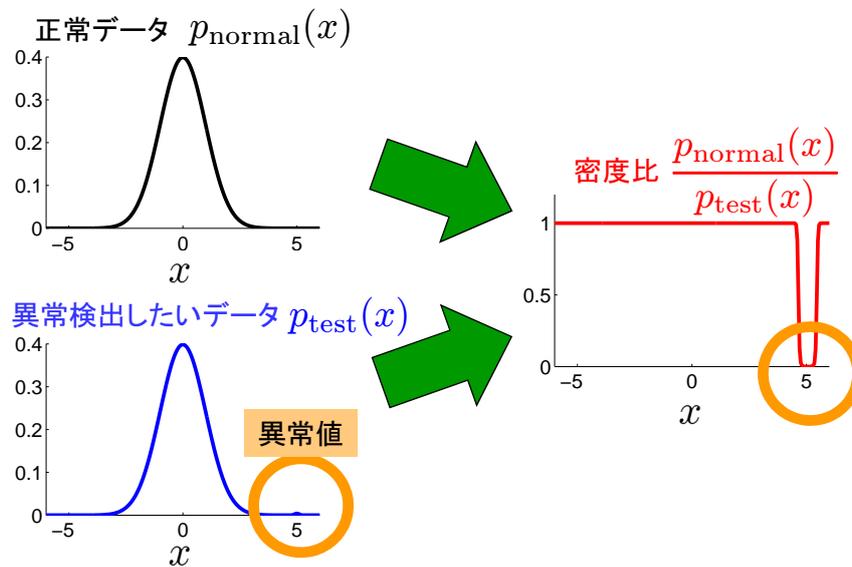


図 4: 密度比に基づく異常値検出. $p(x)$ の異常値を直接検出するのは困難だが (左下), 正常データの密度 (左上) との比を取ることで, 異常値が強調され検出が容易になる (右).

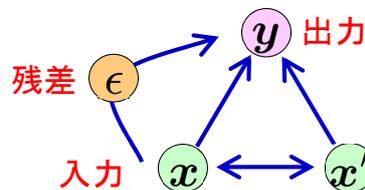


図 5: 相互情報量に基づく独立性判定. 入出力間の独立性判定により, 特徴選択, 特徴抽出, クラスタリングなどが行え, 入力間の独立性判定により, 独立成分分析やオブジェクト適合などが行える. また, 入力と残差の間の独立性判定により, 因果推論が行える.

3.4 条件付き確率推定

確率密度 $p_{x,y}^*(\mathbf{x}, \mathbf{y})$ を持つ分布に独立に従う n 個の標本 $\{(\mathbf{x}_k, \mathbf{y}_k)\}_{k=1}^n$ から, 条件付き確率

$$p_{y|x}^*(\mathbf{y}|\mathbf{x}) = \frac{p_{x,y}^*(\mathbf{x}, \mathbf{y})}{p_x^*(\mathbf{x})}$$

を推定する問題を考える. $\{(\mathbf{x}_k, \mathbf{y}_k)\}_{k=1}^n$ を分子の確率分布からの標本とみなし, $\{\mathbf{x}_k\}_{k=1}^n$ を分母の確率分布からの標本とみなせば, 密度比推定により条件付き確率が推定できる (図 6). \mathbf{y} が連続変数の場合, これは条件付き密度推定に対応し, \mathbf{y} がカテゴリ変数の場合は確率的パターン認識となる (図 6).

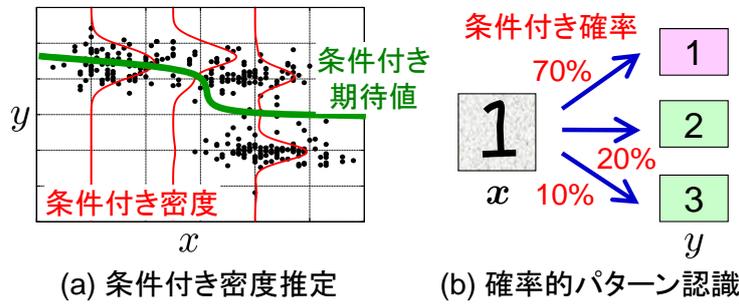


図 6: 条件付き確率推定. (a) 出力変数 y が連続値を取るとき, 条件付き密度の推定に対応する. これは, 条件付き期待値を推定する回帰分析の一般化になっており, 出力の条件付き分布が多峰性や非対称性を持つときに有用である. (b) 出力変数 y がカテゴリ値を取るとき, 確率的パターン認識とよばれ, カテゴリの予測だけでなく予測の信頼度も同時に得られる.



図 7: 二次のマルコフネットワーク. $k = 1, \dots, d$ が頂点で, $\|\theta_{k,k'}\| > 0$ のときに頂点 k, k' 間に枝があるグラフとして, 確率分布が表現される. マルコフネットワークの変化は, グラフの枝の有無の変化に対応する.

4 密度比推定に基づく教師なし変化検知

本節では, 確率分布比較の一例である変化検知について詳しく説明する.

教師なし変化検知の目的は, 確率密度 $p^A(\mathbf{x})$ を持つ確率分布に独立に従う標本 $\{\mathbf{x}_i^A\}_{i=1}^{n_A}$ と確率密度 $p^B(\mathbf{x})$ を持つ確率分布に独立に従う標本 $\{\mathbf{x}_j^B\}_{j=1}^{n_B}$ から与えられたとき, $p^A(\mathbf{x})$ と $p^B(\mathbf{x})$ が等しいかどうかを判定することである.

第 3.2 節で述べた方法で $p^A(\mathbf{x})$ と $p^B(\mathbf{x})$ の距離を推定すれば, $p^A(\mathbf{x})$ と $p^B(\mathbf{x})$ が等しいかどうかを判定できる. ここではさらに, $p^A(\mathbf{x})$ と $p^B(\mathbf{x})$ が二次のマルコフネットワーク

$$q(\mathbf{x}; \theta) \propto \exp \left(\sum_{k \geq k'} \theta_{k,k'}^\top \mathbf{f}(x^{(k)}, x^{(k')}) \right)$$

によりモデル化できると仮定する. ただし, $\mathbf{f}(x, x')$ は特徴ベクトルである.

このモデルのパラメータ θ の変化から, 多次元ベクトル $\mathbf{x} = (x^{(1)}, \dots, x^{(d)})^\top$ の要素 $x^{(k)}, x^{(k')}$ 間の相互作用の変化を捉えられる (図 7). $\mathbf{f}(x, x') = xx'$ とおけば, 二次のマルコフネットワークは正規モデルと一致し, $x^{(k)}, x^{(k')}$ 間の共分散をモデル化することに対

応する. $\mathbf{f}(x, x') = [x^t, x^{t-1}x', \dots, x, x', 1]^\top$ のような高次特徴を考えると, $x^{(k)}, x^{(k')}$ 間の高次の共分散を考慮したモデル化が行える.

モデル $q(\mathbf{x}; \boldsymbol{\theta})$ を標本 $\{\mathbf{x}_i^A\}_{i=1}^{n_A}$ と $\{\mathbf{x}_j^B\}_{j=1}^{n_B}$ にそれぞれ適合させ, 推定したパラメータ $\hat{\boldsymbol{\theta}}^A, \hat{\boldsymbol{\theta}}^B$ の差 $\hat{\boldsymbol{\theta}}^A - \hat{\boldsymbol{\theta}}^B$ を求めることにより, 変数間の相互作用の変化を捉えられる. 例えば, スパース最尤推定を用いると,

$$\begin{aligned} \max_{\boldsymbol{\theta}^A} \sum_{i=1}^{n_A} \log q(\mathbf{x}_i^A; \boldsymbol{\theta}^A) - \lambda_A \sum_{k \geq k'} \|\boldsymbol{\theta}_{k,k'}^A\|, \quad \lambda_A \geq 0 \\ \max_{\boldsymbol{\theta}^B} \sum_{j=1}^{n_B} \log q(\mathbf{x}_j^B; \boldsymbol{\theta}^B) - \lambda_B \sum_{k \geq k'} \|\boldsymbol{\theta}_{k,k'}^B\|, \quad \lambda_B \geq 0 \end{aligned}$$

しかし, 我々が知りたいのはパラメータが変化したかどうかだけであり, 個々のパラメータの値 $\hat{\boldsymbol{\theta}}^A, \hat{\boldsymbol{\theta}}^B$ を求めることが目的ではない.

そこで, 二つのパラメータの差 $\boldsymbol{\theta}^A - \boldsymbol{\theta}^B$ を直接スパース化する学習法が提案された.

$$\begin{aligned} \max_{\boldsymbol{\theta}^A, \boldsymbol{\theta}^B} \sum_{i=1}^{n_A} \log q(\mathbf{x}_i^A; \boldsymbol{\theta}^A) + \sum_{j=1}^{n_B} \log q(\mathbf{x}_j^B; \boldsymbol{\theta}^B) \\ - \gamma \sum_{k \geq k'} \|\boldsymbol{\theta}_{k,k'}^A - \boldsymbol{\theta}_{k,k'}^B\|, \quad \gamma \geq 0 \end{aligned}$$

しかし, この方法でもまだ二つのパラメータ $\boldsymbol{\theta}^A, \boldsymbol{\theta}^B$ を明示的に推定している. また, マルコフネットワークの正規化項 $\int \exp(\sum_{k \geq k'} \boldsymbol{\theta}_{k,k'}^\top \mathbf{f}(x^{(k)}, x^{(k')})) d\boldsymbol{\theta}$ は, 正規モデルやノンパラ正規モデル [5] など, 限られたモデルに対してしか効率良く計算することができない.

これらの問題点は, パラメータ $\boldsymbol{\theta}^A, \boldsymbol{\theta}^B$ を明示的に推定せず, それらの差 $\boldsymbol{\alpha} = \boldsymbol{\theta}^A - \boldsymbol{\theta}^B$ を直接推定することにより解決できる [6]. パラメータの差 $\boldsymbol{\alpha}$ の推定は, 対数線形密度比モデルの学習に対応する.

$$r(\mathbf{x}; \boldsymbol{\alpha}) = \frac{q(\mathbf{x}; \boldsymbol{\theta}^A)}{q(\mathbf{x}; \boldsymbol{\theta}^B)} \propto \exp\left(\sum_{k \geq k'} \boldsymbol{\alpha}_{k,k'}^\top \mathbf{f}(x^{(k)}, x^{(k')})\right)$$

対数線形密度比モデルは, 第 2.3 節で紹介した一般化カルバック距離のもとでの密度適合法により, 効率良く学習できる.

$$\min_{\boldsymbol{\alpha}} \int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{p'(\mathbf{x})r(\mathbf{x}; \boldsymbol{\alpha})} d\mathbf{x} + \gamma \sum_{k \geq k'} \|\boldsymbol{\alpha}'_{k,k'}\|$$

5 まとめ

様々な機械学習タスクを統一的に解決できる密度比推定の枠組みを紹介した. 密度比推定の精度や計算効率を向上させれば, 密度比推定に基づく全ての機械学習アルゴリズムの

性能を改善できるため、密度比推定手法の更なる発展が望まれる。近年、密度比よりも安定した量である相対密度比 [7]

$$r_{\beta}^*(\mathbf{x}) = \frac{p_{\text{nu}}^*(\mathbf{x})}{\beta p_{\text{nu}}^*(\mathbf{x}) + (1 - \beta)p_{\text{de}}^*(\mathbf{x})}, \quad 0 \leq \beta < 1$$

や密度差 [8]

$$d(\mathbf{x}) = p_{\text{nu}}^*(\mathbf{x}) - p_{\text{de}}^*(\mathbf{x})$$

の直接推定法も研究されている。また、連続出力の隠れマルコフモデルの学習 [9] や教師なしラベル付け [10] など、密度比推定や密度差推定により解決できる新たな機械学習タスクを開拓することも今後の重要な課題である。

参考文献

- [1] M. Sugiyama, T. Suzuki, and T. Kanamori, *Density Ratio Estimation in Machine Learning*, Cambridge University Press, Cambridge, UK, 2012.
- [2] M. Sugiyama and M. Kawanabe, *Machine Learning in Non-Stationary Environments: Introduction to Covariate Shift Adaptation*, MIT Press, Cambridge, Massachusetts, USA, 2012.
- [3] M. Sugiyama, S. Liu, M.C. duPlessis, M. Yamanaka, M. Yamada, T. Suzuki, and T. Kanamori, “Direct divergence approximation between probability distributions and its applications in machine learning,” *Journal of Computing Science and Engineering*, vol.7, no.2, pp.99–111, 2013.
- [4] M. Sugiyama, “Machine learning with squared-loss mutual information,” *Entropy*, vol.15, no.1, pp.80–112, 2013.
- [5] H. Liu, J. Lafferty, and L. Wasserman, “The nonparanormal: Semiparametric estimation of high dimensional undirected graphs,” *The Journal of Machine Learning Research*, vol.10, pp.2295–2328, 2009.
- [6] S. Liu, J. Quinn, M.U. Gutmann, and M. Sugiyama, “Direct learning of sparse changes in Markov networks by density ratio estimation,” *Machine Learning and Knowledge Discovery, Part II*, eds. by H. Blockeel, K. Kersting, S. Nijssen, and F. Železný, vol.8189, pp.596–611, *Lecture Notes in Computer Science*, Springer, Berlin, Sep. 23–27 2013.
- [7] M. Yamada, T. Suzuki, T. Kanamori, H. Hachiya, and M. Sugiyama, “Relative density-ratio estimation for robust distribution comparison,” *Neural Computation*, vol.25, no.5, pp.1324–1370, 2013.

- [8] M. Sugiyama, T. Suzuki, T. Kanamori, M.C. duPlessis, S. Liu, and I. Takeuchi, “Density-difference estimation,” *Neural Computation*, vol.25, no.10, pp.2734–2775, 2013.
- [9] J.A. Quinn and M. Sugiyama, “Density ratio hidden Markov models,” Technical Report 1302.3700, arXiv, 2013.
- [10] M.C. duPlessis and M. Sugiyama, “Clustering unclustered data: Unsupervised binary labeling of two datasets having different class balances,” Technical Report 1305.0103, arXiv, 2013.