

密度比推定

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>

統計的機械学習のほとんど全てのタスクは、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処理を実現することが望ましい。ここでは、確率分布でなく確率密度の比を通して様々なデータ処理タスクを解決する密度比推定の枠組みを紹介する [1].

1 密度比推定手法

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、確率密度比

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

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。

1.1 確率的分類法 (Probabilistic Classification)

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

$$p_{\text{nu}}^*(\mathbf{x}) = p^*(\mathbf{x} | y = \text{‘nu’}), \quad p_{\text{de}}^*(\mathbf{x}) = p^*(\mathbf{x} | y = \text{‘de’})$$

と表すことができ、ベイズの定理より、密度比を

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

と表現することができる。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 $\hat{r}(\mathbf{x})$ を次式で求めることができる：

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

1.2 積率適合法 (Moment Matching)

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：

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

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。普遍再生核 $K(\mathbf{x}, \mathbf{x}')$ を用いれば、これを効率よく実現することができる：

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

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。

1.3 密度適合法 (Density Fitting)

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

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

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。

1.4 密度比適合法 (Density-Ratio Fitting)

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。更に非負拘束と ℓ_1 正則化項を加えた場合は、全ての正則化パラメータに対する解が効率よく計算できる。

1.5 統一的枠組み

上記の最小二乗密度比適合法を一般化し、ブレッグマン距離のもとで $r(\mathbf{x})$ 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 : $(t^{1+\alpha} - t)/\alpha, (\alpha > 0)$

1.6 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 $p_{\text{nu}}^*(\mathbf{u})/p_{\text{de}}^*(\mathbf{u})$ と簡略化することができる. 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, 局所フィッシャー判別分析などの教師付き次元削減手法により $\{\mathbf{x}_i^{\text{nu}}\}_{i=1}^{n_{\text{nu}}}$ と $\{\mathbf{x}_j^{\text{de}}\}_{j=1}^{n_{\text{de}}}$ を最もよ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.

2 密度比に基づく機械学習

2.1 重点標本化 (Importance Sampling)

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すことができる. 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:

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

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.

2.2 分布比較 (Distribution Comparison)

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することができる。

また、密度比推定量 $\hat{r}(\mathbf{x})$ を用いることにより、分布間の距離を精度よく推定することができる：

- 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することができる。

2.3 相互情報量推定 (Mutual Information)

確率密度 $p_{\mathbf{x},\mathbf{y}}^*(\mathbf{x}, \mathbf{y})$ 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

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

を推定する問題を考える。ただし、 $p_{\mathbf{x}}^*(\mathbf{x})$ と $p_{\mathbf{y}}^*(\mathbf{y})$ 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

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

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することができる。

2.4 条件付き確率 (Conditional Probability) 推定

確率密度 $p_{\mathbf{x},\mathbf{y}}^*(\mathbf{x}, \mathbf{y})$ 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

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

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。 \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。

3 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。また、密度比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重要である。

References

- [1] M. Sugiyama, T. Suzuki & T. Kanamori: *Density Ratio Estimation in Machine Learning*, Cambridge University Press, 2012.
- [2] M. Sugiyama & M. Kawanabe: *Machine Learning in Non-Stationary Environments*, MIT Press, 2012.