

最尤共クラスタリングによる二相對応関係の確率モデル化

群馬大学理工学府 海老原 諒介, 関 庸一

1 はじめに

顧客群とブランド群の間の購買関係データなどを理解しようとするとき、二相の対応関係のモデル化が必要となる。観測されるのが顧客の購買の有無のみなら、2相の関係は単純な二部グラフ構造をもち、共クラスタリング [?] の問題となる。

しかし、継続した複数回の購買行動が与えられるような場合は購買頻度等の重み付き二部グラフとなり複雑な関係となる。このような二相間の関係を観測した頻度行列に対して、MDL [?] 基準でさらに簡潔なモデル表現を選択し、共クラスタリングするアルゴリズムを提案する。

2 提案手法

顧客 $i (i = 1, \dots, I)$ とサービス $j (j = 1, \dots, J)$ が、それぞれ N 個と M 個のセグメント $\mathbf{n} \in \{1, \dots, N\}^I$ 、 $\mathbf{m} \in \{1, \dots, M\}^J$ に整理され、観測データ x_{ij} は各セグメント組合せごとに同一分布に従うものとする。たとえば、顧客セグメント n_i は商品セグメント m_j を確率 $p_{n_i m_j}$ で購入するとすると、ベルヌーイ分布 $x_{ij} \sim \mathcal{B}(p_{n_i m_j})$ となる。生起頻度を $\mathbf{X} = \{x_{ij}\}_{j=1, \dots, J}^{i=1, \dots, I}$ 、分布パラメータを $\mathbf{P} = \{p_{nm}\}_{m=1, \dots, M}^{n=1, \dots, N}$ とする。

ここで、すべてのセグメント間に関係があれば、完全二部グラフ $K_{N, M}$ で表現されることとなる。しかし、偶然のサービス利用のように閾値 p_{th} 以下の無視できるような弱い関係が含まれており、このような関係は整理して、より単純なグラフで表現できることが期待される場合も考えられる。そこで、閾値以下のパラメータは、所属観測データを統合して求めることで弱い関係を単純化して評価する。このとき、モデルパラメータ数は $l = n \times m - \#\{\hat{p}_{nm} < p_{th}\} + 1$ となる。

簡潔なモデル表現を選択するために、 $f(x|p)$ をデータの確率関数として、(1) 式の MDL 基準を用いる。

$$\text{MDL}(\mathbf{P}) = - \sum_{i,j} \log f(x_{ij}|p_{n_i m_j}) + \frac{l}{2} \log IJ \quad (1)$$

第一項は対数尤度であり、モデルの当てはまりの良さを表している。第二項はモデルの複雑さを表している。

アルゴリズムとして図 1 に示す MLBICL を提案する。これは最尤基準に基づき EM アルゴリズムを用いた共クラスタリングである。これを用いて局所

```
1: procedure MLBICL( $\mathbf{X}, \mathbf{n}, \mathbf{m}, K, C$ )
2:    $\hat{\mathbf{P}} = \text{Est}(\mathbf{X}, \mathbf{n}, \mathbf{m})$ 
3:    $\mathbf{P} = \mathbf{0}$ 
4:   for  $k = 1$  to  $K$  do
5:     for  $i = 1$  to  $I$  do
6:        $\hat{n}_i = \arg \max_n \text{PartLrow}(\mathbf{X}, \hat{\mathbf{P}}, n_i, \mathbf{m})$ 
7:     end for
8:     for  $j = 1$  to  $J$  do
9:        $\hat{m}_j = \arg \max_m \text{PartLcol}(\mathbf{X}, \hat{\mathbf{P}}, \mathbf{n}, m_j)$ 
10:    end for
11:     $\hat{\mathbf{P}} = \text{Est}(\mathbf{X}, \hat{\mathbf{n}}, \hat{\mathbf{m}})$ 
12:     $(\hat{\mathbf{P}}, \hat{\mathbf{n}}, \hat{\mathbf{m}}) = \text{Shrink}(\hat{\mathbf{P}}, \hat{\mathbf{n}}, \hat{\mathbf{m}})$ 
13:     $p_{nx} = \min\{\min_{nm} \{p_{nm} \mid p_{nm} > p_{th}\}, \max_{nm} p_{nm}\}$ 
14:     $p_{bf} = \max\{\max_{nm} \{p_{nm} \mid p_{nm} < p_{th}\}, \min_{nm} p_{nm}\}$ 
15:     $p_{th} = \arg \max_{p_{th} \in \{p_{nx}, p_{th}, p_{bf}\}} \text{MDL}(\mathbf{X}, \text{Cut}(\hat{\mathbf{P}}, p_{th}), \hat{\mathbf{n}}, \hat{\mathbf{m}})$ 
16:     $\hat{\mathbf{P}} = \text{Cut}(\hat{\mathbf{P}}, p_{th})$ 
17:    if  $\max_{nm} |\hat{P}_{nm} - P_{nm}| < C$  then
18:      Break
19:    end if
20:     $\mathbf{P} = \hat{\mathbf{P}}$ 
21:  end for
22:  return  $(\hat{\mathbf{P}}, \hat{\mathbf{n}}, \hat{\mathbf{m}}, p_{th})$ 
23: end procedure
```

- K : 繰り返し数
- \mathbf{n}, \mathbf{m} : 所属セグメント
- C : 収束判定の閾値
- $\text{Est}(\mathbf{X}, \mathbf{n}, \mathbf{m})$: \mathbf{P} を推定する関数
- $\text{PartLrow}(\mathbf{X}, \mathbf{P}, n, \mathbf{m}), \text{PartLcol}(\mathbf{X}, \mathbf{P}, \mathbf{n}, m)$: 顧客 n 、サービス m の部分尤度を計算する関数。
- $\text{Cut}(\mathbf{P}, p)$: 閾値 p 以下のセグメント組合せについてプーリングしてパラメータを求め置換する関数。
- $\text{Shrink}(\mathbf{P}, \mathbf{n}, \mathbf{m})$: 所属のないセグメントを削除する関数。

図 1: MLBICL アルゴリズム

解をできるだけ避けるため十分な回数の推定を行い、最良 MDL を与える解を選ぶ。

実データとして、MovieLens データ [?] を用いた評価を行った結果を報告する。

参考文献

- [1] M.Deodhar and J.Ghosh, SCOAL: A Framework for Simultaneous Co-Clustering and Learning from Complex Data, ACM, Trans. on Knowledge Discovery from Data, 4(3), 2010.
- [2] Rissanen, J., A Universal Prior for Integers and Estimation by Minimum Description Length, *The Annals of Statistics*, 11(2), 416–431, 1983.
- [3] F.M.Harper, and J.A.Konstan, The movielens datasets: History and context., ACM Trans. on Interact. Intell. Syst., 5(4), 2016