

多項ロジットモデル及び主成分分析を用いた統計的マッチング手法の提案

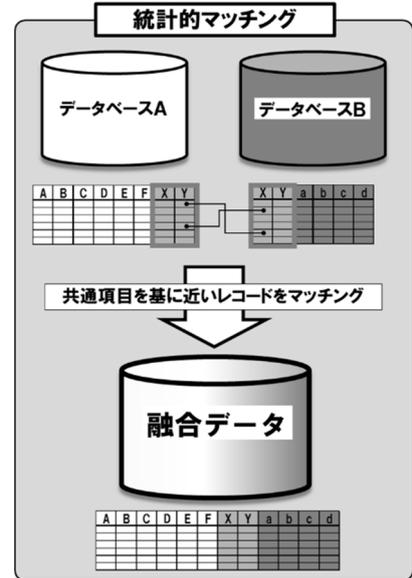
総合研究大学院大学 高部 勲
統計数理研究所教授 山下 智志

1. 統計的マッチングの概要と課題

データリンケージは、異なるデータベースを結合し、豊富な情報を持つ単一のデータベースを構築する技術である。データリンケージにより、新たな調査やデータ収集を行うことなく、有用なデータの作成が可能となることから、近年、様々な分野で利用されるようになってきている。

データリンケージを行う際に、各レコードを識別できる照合キー（共通一連番号、名称、所在地など）が存在する場合は完全照合（Exact Matching）が可能だが、そのようなキーが存在しない、あるいは利用できない場合には、レコード間の類似度を表す距離関数を定義し、近いレコード同士を結合する統計的マッチング（Statistical Matching）の手法が用いられる。

マッチングにはウエイト付き距離が用いられることが多いが、ウエイトの決定方法が恣意的であるとの指摘がある。また企業データは一般にレコード間の差が大きいため距離の定義が難しく、さらにデータ量も多く計算量も多くなるため、統計的マッチングの適用例はあまり多くない。



【ウエイト付き距離関数の例】 $d_{ij} = \sum_{p=1}^P \beta_p |X_{ip} - X_{jp}|$

2. 提案手法と結果

本報告では、多項ロジットモデル及び主成分分析を用いた統計的マッチング手法を提案する。2種類の企業データの共通フィールドから算出したレコード間の距離を説明変数として多項ロジットモデルを構築することにより、ウエイトの合理的な決定及びマッチング精度の確率的な評価が可能となる。

ここで各データのレコード数が増加した場合、距離計算の対象となるレコードの組合せが飛躍的に増大し、計算が困難となることから、それらをうまく削減する必要がある。本報告では主成分分析によりデータを層化し、近隣のレコードのみを距離・尤度計算の対象とすることで計算の効率化を図っている。

提案手法を商用データ及び経済センサスのマイクロデータに適用した結果、誤分類率などの点で、従前のデータマッチングの手法（Nearest Neighbor Methodなど）よりも優れていることが示された。

(1) 多項ロジットモデルの概要

- 複数の選択肢から選択対象を確率的に決定するモデル。
- 交通機関の選択に関する意思決定の分析などに利用。

【交通機関・手段の選択問題の例】

(2) 多項ロジットモデルに基づく統計的マッチング

【上記のモデルの枠組みを統計的マッチングの問題に適用】

マッチング「元」レコード(企業 i)

マッチング「先」レコード(企業 j)

マッチング確率

レコード(企業)間の距離

$D_{ij} = \beta_1 |X_{i,乗車時間} - X_{j,乗車時間}| + \beta_2 |X_{i,乗車時間} - X_{j,乗車時間}| + \dots$

- マッチング「元」のデータベースのレコード(企業)を選択主体とみなす。
- マッチング「先」のデータベースのレコード(企業)を選択肢とみなす。
- レコード間のウエイト付き距離を各選択肢の効用とみなす。
- ⇒多項ロジットモデルの枠組を統計的マッチングの問題に適用することが可能
- ⇒マッチングの精度 について、確率 の形で表現することが可能。(レコード間の距離が大きいくほど、マッチング確率が低くなる考えられる。)
- 距離関数のウエイトをパラメータとして扱い、最尤法により推定。
- ⇒距離関数のウエイトを 統計的・合理的に決定することが可能。
- ⇒カテゴリ変数 の距離のウエイトについても、同じ枠組みで扱うことが可能。

参考文献：

[1] D’Orazio, M., M. Di Zio & M. Scanu (2006), *Statistical Matching: Theory and Practice*, Wiley
[2] Rässler, S. (2002), *Statistical Matching*, Springer 本研究は科研費（16H02013及び15H03390）の助成を受けている。