

決定木を用いた質的データからのグラフ構造学習

川崎 隆史† 鈴木 大輔† 杉山 歩† Dam Hieu Chi†

†北陸先端科学技術大学院大学 知識科学研究科

1 はじめに

近年、グラフやネットワークから知識発見を行おうという試みが増加している。SNSからのコミュニティ抽出や、タンパク質の相互作用リンク予測などはその一例である。

しかしながら、グラフを用いた知識発見には、分析の目的に合わせてデータからグラフ構造を決定する必要がある。グラフ構造を如何にして決定するかという問題は非常に重要である。その中でも、近年増加している多次元の複雑なデータに対応するため、重要な関係性のみを取り出してグラフ構造を決定するという変数選択の問題を考慮した手法が求められている。

このような問題を解決する手法としてN.Meinshausen-P.Bühlmannの手法(MB法)[1]が提案された。これは回帰分析を全変数に網羅的に行うことによってグラフ構造を学習する手法であり、変数選択にはLasso正則化[2]を用いている。これにより、多次元の複雑なデータから重要な関係性のみを抽出したグラフ構造を学習することに成功した。

しかし、MB法は回帰分析を用いるために量的なデータを分析するには適切であるが、質的データの分析、とりわけ二値データに対する分析には不適切である。この問題に対し、ロジスティック回帰を用いて対応した既存研究が存在する[3]。しかし、目的変数に対する説明変数間の関係性を吟味出来ないなどの問題があり[4]、あらゆる場合において適切な手法であるとはいえないことがわかっている。

2 研究目的

本研究の目的は、質的データから説明変数間の関係性を考慮したグラフ構造を学習することである。そのために、決定木を用い、変数選択に変数減少法を採用する。グラフ構造学習のアルゴリズムにMB法の発想を利用する。これにより、MB法の優位点を継承した質的データを扱う手法の開発を目指す。

3 研究方法

3.1 決定木によるグラフ構造学習

MB法によるグラフ構造の学習は、データの持つN個の変数 (x_1, x_2, \dots, x_N) について各変数 x_i を目的変数とし、その他を説明変数とした回帰分析をN回行い、その結果を統合することにより行われる。本研究はこの方法を利用する。

決定木による構造学習の流れを下記に示す。なお、決定木のアルゴリズムはCART、変数重要度の算出に関しては、Gini-Importanceを採用した。モデル推定には10-fold-Cross Validation(CV)を用い、各分割データにおいて変数減少法を用いたモデル構築を行う。その結果から、最良なモデルの構築を行う。

Algorithm 1 決定木によるグラフ構造学習

- 1: Input: データ $X = (x_1, x_2, \dots, x_N)$
- 2: Output: 変数重要度の行列
- 3: **for** $i = 1$ to N **do**
- 4: Set $f(x_i) = (x_1, x_2, \dots, x_{N-1})$
- 5: $f(x_i)$: モデル推定
- 6: $f(x_i)$: モデル構築
- 7: $f(x_i)$: 変数重要度の算出
- 8: **end for**

3.2 決定木とグラフ構造の対応

決定木モデルとグラフ構造の対応を図1に示す。変数A,B,Cが存在する時に、Aが目的変数だった場合の決定木とグラフを表している。エッジの方向はB,CがAを説明することを意味し、実線は決定木の第一層、破線は第二層以降に用いられた変数で、それ以前の層による条件付きの関係性であることを表している。

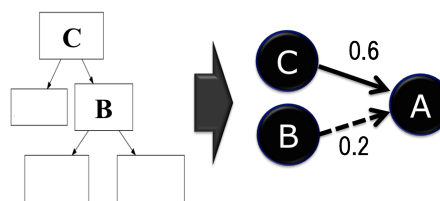


図1: 決定木モデルとグラフ構造の対応

Graph Structure Learning from Categorical data by Decision Tree
 †Takafumi KAWASAKI †Daisuke SUZUKI †Ayumu SUGIYAMA
 †Dam Hieu Chi
 †School of Knowledge Science, Japan Advanced Institute of Science and Technology

