

# インタラクティブにクラスタリングを行い、 構造化を行うツールの開発と実査

伊藤貴一\*1・熊坂賢次\*2

Email: takaichi.ito@gmail.com

\*1: 慶應義塾大学 SFC 研究所

\*2: 慶應義塾大学環境情報学部

## ◎Key Words ツール開発, クラスタリング, 情報可視化

### 1. はじめに

インターネットが社会に浸透するに従い、多くの人が Blog や SNS を使うようになり、人々のライフスタイルを Web サービス上に表明するようになってきている。このような状況下、社会調査も、アンケート調査をするというものだけではなく、Web にある人々の声を拾う、ソーシャルリスニング[1]が一分野になってきている。ソーシャルリスニングのためのツールが求められている。

このようなツールは、あらかじめ明確な答えがないため、探索的アプローチになってしまう。そのため、探索を支援するツールであるべきだ。この探索のためには、機械処理の結果を見せるだけのシステムではいけない。人間の背景知識や分析意図を結果に反映させるようなものでなくてはならない。そのため、インタラクティブ性は重要であり、人間とデータと機械処理が融合するような、知的インタラクティブシステム[2]である必要がある。

このような目的のために、以前に私は、「ひっぱりくん」を開発した[3]。本論では、共起関係のデータを用いてその関係の構造化を行うツールとして、「こうぞうくん」を開発した。

### 2. ツール開発のコンセプト

この論文のツールで扱うのは、バスケット分析の可視化である。商品の購買履歴や、自然言語を形態素解析後のデータを用いた分析である。共起関係に基づきアイテム間の関係を可視化する。これにより、商品購買なら、購買の関係図、自然言語なら、言葉の関係図を作り、データにある構造を読み解くことができるというものである。

このようなデータを分析するために、縦列をクラスタ、横列を頻度のレイヤーの二軸を使い、表形式でインタラクティブに表すツールを作成した。これは、山崎[4]の論文において社会分析に使われた手法のツール化である。

最終的には、Fig.1 のような可視化を行う。

以下、作成したコンセプトを示す。

#### 2.1 べき乗分布とレイヤー

自然言語でのジップの法則[5]、マーケティングで言うロングテール[6]は、べき乗分布ゆえの現象である。べき乗分布を分析するには、ピラミッド状にデータを分割する必要がある。このために、レイヤー分割というを行う。頻度の最頻値と、最低値の対数を取り、それを当分割したものを区切りとして、頻度の層を作る。このようにすると、一番上のレイヤーが、メジャーなもの、一番下のものがマイナーなものとなり、配置により、データの規模

感を把握できるようにしている。

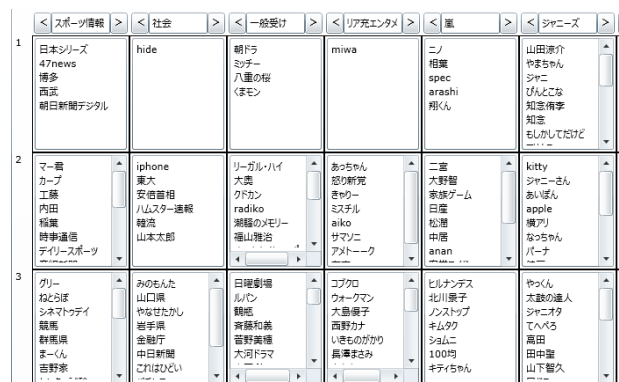


Fig.1

### 2.2 インタラクティブな関係の表示とクラスタリング

関係性の表示のために、アイテムのクリック時、関係が強いアイテムに色を付けるということをしている。これは、グラフにおける、エッジを、インタラクティブに見せることに相当する。そのため、複雑な模様になってしまいがちなネットワーク図と同じ情報をすっきりと見せるようにしている。縦の列で、なるべく関係の強いもので固めるという形で、クラスタリングを行う。教師なしで、関係のみを用いて行うため、これは自己組織化させているともいえる。Fig.2 において、赤は選択したアイテムであり、橙色は、関係しているアイテム。薄い橙は、弱い関係である。ネットワークを可視化しているといえる。ただ、クラスタリングした結果を見るというのではなく、その結果から、分析者の考えによって、違うところにアイテムがあると収まりがいいと思った時、それを移動することができるようにしている。

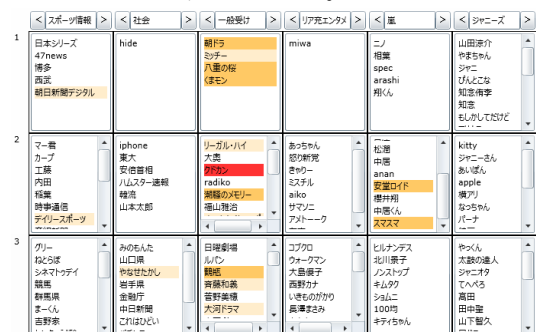


Fig.2

