

相互類似関係を用いたアソシエーション分析

岩崎 幸子

マーケット・バスケット分析に代表されるアソシエーション分析がさまざまな企業の現場において使用されるようになって久しい。しかしながら、膨大に列挙されるルールのフィルタリングを行い、知りたいルールを発見することは容易ではない。そこで本稿では、アイテム間の相互類似関係を用いてルール抽出する手法を紹介する。

キーワード：購買履歴データ、アソシエーション・ルール、マーケット・バスケット分析

1. はじめに

情報技術の発展に伴い、小売業では購買履歴データをはじめとするシステムが普及している。近年では、製造業や卸売業に購買履歴データを開示して、製・配・販で商品の品揃えから販売施策などを行う「協働マーチャンダイジング」に取り組む小売業も増加している。購買履歴データの分析とその活用用途はさまざまであるが、最も盛んに取り組まれているのはアソシエーション・ルールを用いたマーケット・バスケット分析であろう。その分析結果から、フロア・レイアウトや定番棚割りなどの売場づくり、また関連販売（クロスセル）などへの応用が期待されている。

しかしながら、このマーケット・バスケット分析は厄介である。分析対象とするアイテムを絞り込まずに、計算に時間がかかり列挙されるルールが膨大になる。列挙されたルールを各種評価指標に閾値を設けてフィルタリングすれども、「牛乳⇒もやし」、「納豆⇒バナナ」のように誰もがよく買う購買が多いアイテムのルールが多く残る一方で、興味深いルールではあるものの、そもそもの購買が少なく販売促進施策を実施するには販売コストと業績パフォーマンスが見合わないルールが残りがちである。新しい発見をしたり、また既知のことであっても重要なルールを再確認したいところであるが、列挙された膨大なルールの中からソートしたりフィルタリングしながらそれらを見つけ出すことはとても骨が折れる作業となる。大変な分析作業の割に、得られるものが乏しいことに落胆することも多いのではないだろうか。

そこで本稿では、これらの問題に対する一つの手段

として、膨大に列挙されたアソシエーション・ルールからアイテム間の相互類似関係を測り、それを用いてルールを抽出する手法について紹介する。最初に、本手法を考案することになった経緯と分析手法の解説を行い、また一般的な閾値によるルール抽出手法との違いについて述べる。そして、株式会社マクロミルの消費者購買履歴データ「QPR」を用いた分析事例を紹介する。

2. 手法考案の経緯

本手法は筆者が小売業で勤務していたときに商談室で考案したものである。筆者は商品部で日用品や加工食品などのバイイング業務を担当しており、前述の「協働マーチャンダイジング」ではバスケット分析を頻繁に業務活用していた。メーカーの営業担当者には多忙な時間を割いてバスケット分析を手伝ってもらい、関連販売提案を頂戴していた。鶏卵の担当をしていたとき、「トマトケチャップは卵と同時購買率が高い」、「おでんの素調味料は卵とのリフト値が高い」と、卵売場で関連調味料の関連販売提案を食品メーカーより頂戴した。分析の手間を考えるとすべて採用したい気持ちになるが、限られた売場にあれもこれもと置くと煩雑な売場となる。また提案の多くはメーカー製品の立場から卵を分析したものであり、卵の立場から相応しい商品を分析しているものではない。卵の立場からみて、そのメーカー製品は最も相応しい関連商品なのか、もっとほかにも適切な商品があるのではないかと疑問をもった。実際に、トマトケチャップやおでんの素調味料と卵は一緒に買われやすい関連商品であったが、卵のような購買規模の大きい商品からすればこれらの調味料は「その他大勢」の一つであり、それほど重要で販促パフォーマンスがよい関連商品とはいえなかった。卵と一緒に取り組むべき相手は別の商品であり、またその別の商品にとって卵と組むのはよいことなのか否か

いわさき さちこ

国立情報学研究所情報学プリンシプル研究系

〒101-8430 東京都千代田区一ツ橋 2-1-2

iwasaki@nii.ac.jp

を考慮することが店頭での全体最適な販売だと考えた。

これらのことから、現場におけるバスケット分析においては、アイテム間の「お互いの親しさ」のような相互関係を考慮して分析することが重要ではないかと考えた。またそのアプローチは店頭での全体最適な商品の組合せを発見しようとするだけでなく、膨大なアソシエーション・ルールから何らかの意味をもつルールを取り出し、無意味なルールを削減することにつながるものかと考えた。

では、そのアイテム間の「お互いの親しさ」というものをいかに測り、その相互関係を捉えるといえるだろうか。これに関してはさまざまなアプローチが存在するだろう。当時、高度な分析知識をもち合わせていない筆者は、共起性評価指標（類似度）にランク情報（順序尺度）を用いるという簡単な方法を考えた。それはたとえば、牛乳からみて類似度が高い上位ランキングの相手商品にパンがあった場合、パンからみて上位ランキングの相手商品を見て、お互いに上位であれば牛乳とパンは親しい商品同士とみなす、という方法である。

また、その方法で共起性評価指標（類似度）の値の大小に捉われずにルール抽出が可能になると考えた。共起性の評価指標には、値の大小を比較しても良し悪しが見えられないものがある。たとえば現場でよく使用されている Lift（次節、式(4)）や Jaccard（同、式(3)）は、アイテムの出現頻度に影響を受ける指標である。Lift は出現頻度の少ないアイテムで値が大きくなる傾向があり、Jaccard は出現頻度が大きく異なる組合せでは値が小さくなる傾向がある。両者ともに出現規模が異なるアイテム間で値の大小を比較することは困難である。たとえば、きゅうりと豆腐の Jaccard が 5% であったとする。Jaccard はきゅうりからみても豆腐からみても同じ値であり、方向性をもたない指標である。しかしながら、この 5% という値は双方にとって同じ意味や重みをもっているだろうか。きゅうりからみて豆腐の 5% は最も高い値であるかもしれないが、豆腐からみればきゅうりとの 5% は低い値であり、より高い値を持つ相手は納豆だったりすることがある。評価指標の良し悪しは、アイテムによって相対的なものであるため、評価指標の閾値フィルタリングに加えてこのランク情報を併用することで値の大小にとらわれずにルール抽出が可能になるだろう。

その後、このアイデアと方法を JST CREST のビッグデータに関連するプロジェクトで共同研究メンバーとともに手法として構築した [1, 2]。

3. 分析手法

本手法は一般的なアソシエーション分析の結果を用いてランク情報によるルール選択を行うものである。

3.1 アソシエーション・ルールの列挙

アイテム集合 I についてのトランザクション集合を $D = \{T_1, T_2, \dots, T_n | T_i \subseteq I\}$ とする。そこから得るアソシエーション・ルールは、アイテム集合 $X, Y \subseteq I$ について $X \Rightarrow Y$ と表現する。小売業のバスケット分析の場合、アイテムとトランザクションはそれぞれ商品とレシート単位となる。アソシエーション・ルール集合を作成するにあたり、次式に示すような共起性の評価指標に閾値（下限値）を与え、その閾値を超えて選択されたルールをアソシエーション・ルール集合 R とする。また本手法ではルール解釈の容易性のため、条件部と結論部のアイテム集合のサイズは 1 に限定する ($|X| = |Y| = 1$)。アイテム X を含むトランザクション集合を $occ(X)$ と表す。

$$\text{Support}(X \Rightarrow Y) = \frac{|occ(X) \cap occ(Y)|}{|D|} \quad (1)$$

$$\text{Confidence}(X \Rightarrow Y) = \frac{|occ(X) \cap occ(Y)|}{|occ(X)|} \quad (2)$$

$$\text{Jaccard}(X \Rightarrow Y) = \frac{|occ(X) \cap occ(Y)|}{|occ(X) \cup occ(Y)|} \quad (3)$$

$$\text{Lift}(X \Rightarrow Y) = \frac{\text{Confidence}(X \Rightarrow Y)}{\text{Support}(Y)} \quad (4)$$

$$\text{NPMI}(X, Y) = \frac{\ln \frac{\text{Support}(X \Rightarrow Y)}{\text{Support}(X)\text{Support}(Y)}}{-\ln(\text{Support}(X \Rightarrow Y))} \quad (5)$$

3.2 ランク情報によるルール選択

次に、アソシエーション・ルール集合 R から、式(1)~(5)で与えた共起性評価指標 (sim) のランク情報を用いて部分ルール集合 R' を選択する。ランク情報すなわち条件部に X を持つ部分集合における $X \Rightarrow Y$ の評価順位を式(6)に定義する。

$$\text{rank}_{\text{sim}}(X \Rightarrow Y) = |\{i : \text{sim}(X \Rightarrow Y) \leq \text{sim}(X \Rightarrow \{i\}), \text{sim}(X \Rightarrow \{i\}) \in R\}_{i \in I}| \quad (6)$$

ランク情報によるルール選択方法は二つあり、一つは「両想い (friend) ルール」、もう一つは「片想い (pal) ルール」と呼ぶ。「両想い (friend) ルール」は式(7)、「片想い (pal) ルール」は式(8)の条件を満たすルール $X \Rightarrow Y$ である。

$$\text{rank}_{\text{sim}}(X \Rightarrow Y) \leq k \text{ and } \text{rank}_{\text{sim}}(Y \Rightarrow X) \leq k \quad (7)$$

$$\text{rank}_{\text{sim}}(X \Rightarrow Y) \leq k \text{ and } \text{rank}_{\text{sim}}(Y \Rightarrow X) > k \quad (8)$$

「両想い (friend) ルール」は相互にランク情報の値が任意の条件を満たしており、強い関係をもつルールである。また「片想い (pal) ルール」は X からみると Y は強い関係であるが、 Y からみればそうでもないというルールである。分析目的に応じて、これらのルール選択方法を自由に組み合わせるなどして、データから知見を得る方法を推奨する。

4. 一般的な分析手法との違い

一般的なアソシエーション分析では、アルゴリズムから出力された結果に対して、さらに各種の共起性の評価指標に任意で閾値を与えてルールの選択（絞込み）を行う（以下、一般手法と呼ぶ）。本手法では、アルゴリズムからの出力結果に対してランク情報を付加して、ランク情報を用いてルールの選択を行う。ランク情報を用いたルール選択は一般手法とどのように違うのか。筆者らは、それぞれの手法から選択されたルールの形質的なグラフ比較をする視覚化実験を行った [1]。

表 1 はその実験の一部であり、株式会社マクロミルの消費者購買履歴データ「QPR」を用いて評価指標別に選択されたルールをグラフ視点で比較したものである。まず、両手法とも最小支持度に同じ下限値を用いてルール集合 R を取得し、次に本手法 ($rank$) で評価指標 sim 、ランク k の両想いルール集合 R' を求めている。一方、一般手法 (th) は本手法で得られたルールと同じ数のルールを、評価指標の上位から選んだルール集合 R'' を求めている。ランク $k=1, 2, 3, 4, 5, 10, 20, 50$ とし、評価指標 $sim = \text{Support}, \text{Jaccard}, \text{NPMI}$ の三つにおいて比較したものである。なお、 R に含まれるアイテム数=307 である。

本手法は、いずれの評価指標においてもアイテムカバー率が高く、分析対象となったアイテムが結果により網羅されている。また、最大次数は大きいても k で抑えられるため一般手法よりも小さい値になり、特定のアイテムに接続が極度に集中せずに分散している。そして、連結成分においては指標依存であるが、本手法は一般手法より連結成分が多い傾向であった。両手法ともに、 k の値が高くなるほどアイテム数が増える。一般手法は最大次数が高くなることから枝が複雑になるが、本手法は最大次数が制限されているため枝が複雑になりにくい。つまり、本手法は一般手法よりもアイテムの網羅性が高く、全体を俯瞰しながらも可読性の高いグラフを得やすい。これらのことから、本手法

表 1 結果の比較

sim	k	edge 数 R' , R''	アイテムカバー率		最大次数		連結成分	
			rank	th	rank	th	rank	th
Support	1	7	0.045	0.026	1	7	7	1
	2	13	0.078	0.029	2	8	11	1
	3	22	0.104	0.042	3	12	11	1
	4	30	0.123	0.055	4	15	13	1
	5	41	0.149	0.068	5	19	14	1
	10	122	0.263	0.127	10	38	18	1
	20	397	0.456	0.257	20	74	18	1
50	2096	0.811	0.586	50	176	9	1	
Jaccard	1	52	0.338	0.127	1	10	52	12
	2	110	0.524	0.179	2	17	56	10
	3	153	0.589	0.221	3	22	48	14
	4	188	0.641	0.254	4	23	42	13
	5	236	0.713	0.296	5	26	33	18
	10	426	0.771	0.358	10	34	12	19
	20	901	0.820	0.589	20	50	7	16
50	2860	0.912	0.840	50	88	2	7	
NPMI	1	80	0.521	0.328	1	6	80	31
	2	167	0.726	0.465	2	11	69	30
	3	252	0.840	0.534	3	17	51	31
	4	319	0.863	0.560	4	21	25	30
	5	396	0.876	0.625	5	26	15	29
	10	763	0.938	0.745	10	40	5	14
	20	1472	0.986	0.892	20	66	3	12
50	3623	1.000	1.000	50	97	2	3	

はアソシエーション・ルールを可視化して意味解釈をしようとするときにも便利である。

5. 消費者購買履歴データのバスケット分析事例

5.1 現場におけるバスケット分析

グラフの形質と可読性以上に、より重要なのはルールの意味解釈性のよさであろう。現場におけるバスケット分析の目的は、そこから得た知見から売り方改善のための「打ち手」を考えることにあり、それに使える知見が分析から得られることが必要である。よく、ビールと紙オムツの都市伝説のように、新しくて面白いルールの発見が期待されることが多く、わかりきったこと、当たり前前のルールが軽視されることがある。しかしながら、筆者は誰もが認識しているような既存のルールを再確認することも同じく重要だと考える。なぜならば、当たり前前のことがきちんと店頭実現できているとは限らないからである。そこに大きな販売機会ロスが存在するのではないだろうか。バスケット分析では新しくて面白いルールの発見ができるに越したことはないが、わかりきったことや当たり前前のことを軽視せずに、それらがきちんとできているかという視点からもこの分析を活用していくことが有意義なデータ活用だと考える。このような考え方から、バスケット分析において不要なのは、意味不明なルールである。「牛乳⇒もやし」、「納豆⇒バナナ」などの誰もがよく買うゆえに自ずと共起頻度が高くなるような、一緒に買われることの理由が見いだしにくいルールなどもそれにあたる。

既存のルール確認ができて願わくば新しい発見、そして意味不明なルールが極力少ない分析結果、そのようなアウトプットが得られるような分析が現場において「使える」分析ではないだろうか。本分析もそのようなアウトプットを目標として次のように実施した。

5.2 分析概要

5.2.1 使用データ

本分析は、株式会社マクロミルの消費者購買履歴データ「QPR」を使用している。データ取得期間は、2012年1月1日から2014年6月30日までの2年半、分析対象は全国のスーパー（SM）業態441小売企業である。トランザクション数（レシート（ $|D|$ ））= 5,569,959、アイテム（ $|I|$ ）はJICFS大分類の「食品」と「日用品」内の細分類510カテゴリを使用した。アイテムに関しては「その他」という文字列を含む細分類カテゴリはすべて分析対象外とした。「その他」を分析対象外とする理由は、意味解釈をやすくするためである。「その他」と名がつくカテゴリには分類不明な商品などが含まれていることが多く、一つのカテゴリでありながらも異なる商品が含まれていることが多い。そのため、さまざまなアイテムと無意味な接続をやすく、意味解釈を妨げるため分析対象外とした。

5.2.2 分析方法

前述のデータを用いて、まず最初にアソシエーション・ルールの列挙を行った。ルールの足切りは $\text{Support}=1.8 \times 10^{-6}$ で、共起頻度 = 10 と低めに設定し、できるだけ多くのルールを拾うようにした。次に、ランク情報によるルール選択は「両想い（friend）ルール」のみを使用した。ランク $k=30$ 、評価指標（ sim ）はJaccardとLiftの2種類ですべてAND条件とした。つまり、条件部（ X ）と結論部（ Y ）のJaccardのランクが相互に上位30位かつLiftのランクが相互に上位30位であるルールを選択したものをアソシエーション・ルール集合 R' としている。その結果をより理解しやすくするため、本分析では「Gephi」[3]というグラフを視覚化するツールを用いて描画した。

5.3 分析結果

前述の方法で得られた両想いルール集合 R' の全体図とそれを大まかに意味解釈した内容を加筆したものを図1に示す。 R' に含まれるアイテムは $|V|=392$ 、ルール数は $|E|=2,688$ であり、アイテムカバー率は76.9%であった。連結成分は4であり、大きな連結成分が一つと、小さな連結成分が三つであった。全体を俯瞰すると、「食品群」、「日用雑貨群」、「化粧品群」、そして「医薬品群」の集合が大まかに識別でき、またそれぞれ

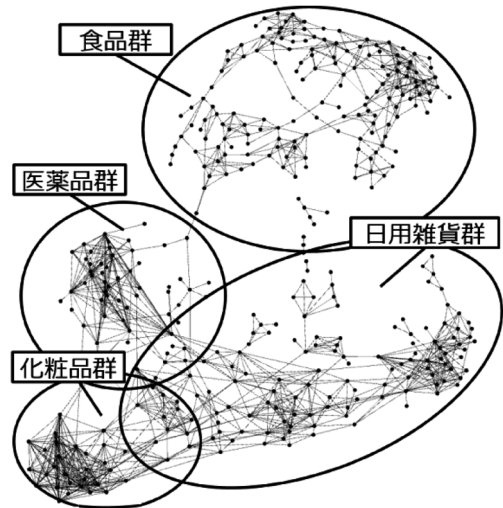


図1 グラフの全体図（枠線と解釈内容は事後的に加筆）

のつながりに意味解釈ができない「意味不明ルール」は少なく、大多数は意味解釈可能なネットワークを形成している。

5.4 ビジネス応用の考察

次に、小売業におけるアソシエーション分析の応用として、購買者が買い物をしやすい売場配置や関連販売など検討する場面を想定し、それらに対する知見が得られるようなバスケット分析のビジネス応用を考察する。本稿では、食品と、日用雑貨と化粧品の関連販売と売場配置について取り上げる。

5.4.1 食品の関連販売

図2, 3, 4は、図1の食品群を部分拡大したものである。まず、図2は、菓子、飲料、パンなどが相互連結している。特に注目したいのは、「スナック」と「カップ麺」の連結である。スナックとは、ポテトチップスなどのスナック菓子類であり、インスタントのカップ麺やコーラなどと連結がある。この結果から、カップ麺は食事よりもおやつ的な買い方がされているように見受けられる。より食事的な調理と食べ方がされるインスタント袋麺がこれら菓子群との連結がない点からもカップ麺と袋麺の消費場面の違いが推測される。スナック菓子は菓子売場、カップ麺は加工食品売場と、大多数の企業では管理部門と売場が異なっており、これらは独立して販売されている。販売促進施策において、スナック菓子だけでなくカップ麺を加えたおやつ企画などは、購買者にとって「おやつ」の選択肢が増えるよい提案となる可能性がある。また、菓子売場と加工食品売場がフロアレイアウトで隣接する場合、スナック菓子とインスタント麺が売場のつながりとして配

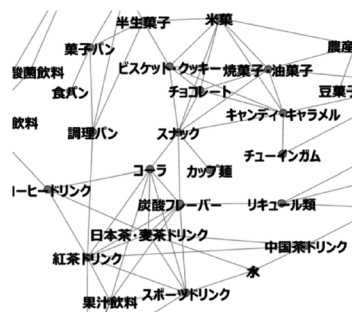


図2 菓子・飲料などの食品群のルール (図1の拡大図)

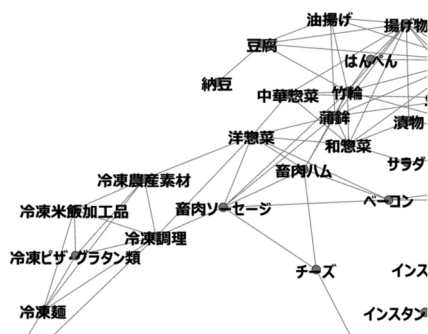


図3 冷食・惣菜・畜肉加工品などの食品群のルール (図1の拡大図)

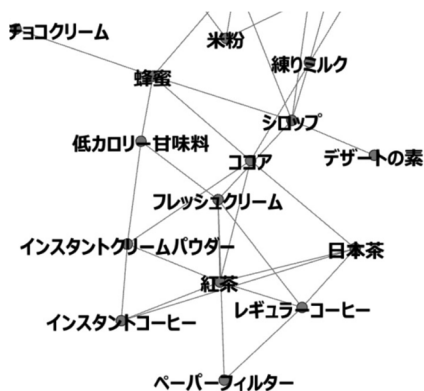


図4 嗜好品などの食品群のルール (図1の拡大図)

置されると購買者にとって便利な買い回りができる可能性がある。本結果からは、このような知見と仮説を得ることができる。

次に、図3は主に、冷凍食品、惣菜、畜肉加工品などが相互連結している部分である。特に、畜肉ソーセージと冷凍調理と洋惣菜が連結しているところが興味深い。冷凍調理はお弁当のおかずなどが多いカテゴリであり、これらは「お弁当」つながりのように見受けら

れる。食肉売場の冷蔵ソーセージと冷凍食品そして惣菜は冷蔵あるいは常温と温度管理が異なる陳列ケースであるため売場を同じにすることは難しい。しかしながら、お弁当企画などの際に、一緒にまとめて購入すると値引きなどの得がある「バンドル企画」などで併買促進をするなどの販売施策の考察が可能である。

また、図4は主に、コーヒーやお茶、デザート関連食材で構成される「嗜好品」と呼ばれる商品群である。ここでは、レギュラーコーヒーとペーパーフィルタの相互連結は重要なルールだと考える。レギュラーコーヒーとは加工食品のコーヒー豆のことで、それをドリップするための雑貨品であるペーパーフィルタと一緒に買うのは容易に想像がつく。これらはそれぞれ食品と雑貨であり、流通チャネルが異なることから実は管理部門も異なり、その結果売場も異なる店舗も多い。ペーパーフィルタがコーヒー売場にもあると、購買者にとって便利で買い忘れしにくい売場となるであろう。

5.4.2 日用雑貨と化粧品売場の配置

もう一つ、パーソナルケア群という日用雑貨と化粧品の売場配置について考察する。小売業のこれらの売場は、大きくパーソナルケア群（化粧品などの身体に使う製品）とホームケア群（衣住居に使う製品群）などに大別されている。パーソナルケア群には、歯磨きやシャンプーなどの日用雑貨品や化粧品がある。これらの商品群の売場配置にはさまざまな捉え方があり、それらの売場配置の考察は悩ましい。たとえばシャンプーの場合、髪に使う「ヘアケア製品」という視点で捉えると、シャンプーは整髪剤やヘアカラーと売場が近いとよいかもしれない。しかしながら、シャンプーを入浴時に使う「インバス製品」と考えると、石鹸や入浴剤などと近いほうがよいかもしれない。また、同じ化粧品といえども男性化粧品と女性化粧品がある。異性の前で化粧品を購入することをためらう購買者がいるため、売場の隣接にも配慮が必要であり売場構成が難しい商品である。図5は、化粧品群と日用雑貨群の部分拡大図であり、パーソナルケア群の商品ネットワークである。図の左側に女性化粧品、中央付近に男性化粧品、そして右側に歯磨きやシャンプーなどの洗面と入浴に関連する商品のネットワークが構成されている。

まず、前述のシャンプーをみると、シャンプーは石鹸やボディシャンプー、入浴剤などと相互連結しており、ヘアスプレーなどの整髪剤とヘアカラーは女性化粧品と男性化粧品の間でこれらの商品と相互連結している。また、リップクリーム、UVケア、制汗剤、と

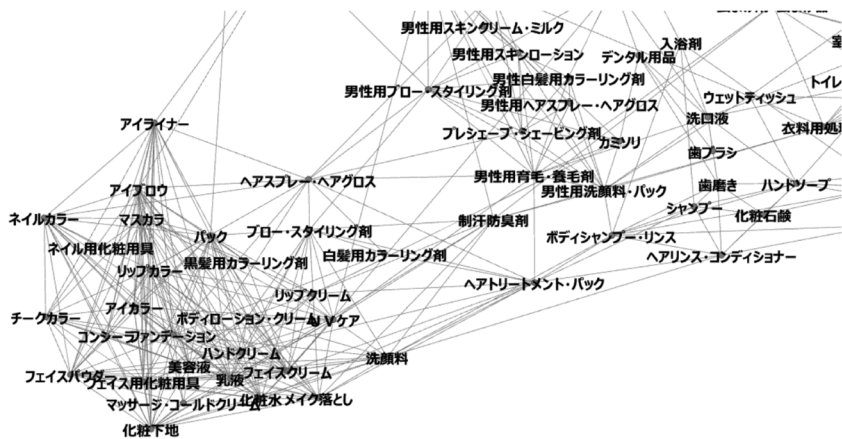


図5 パーソナルケア群 (図1の拡大図)

いったシーズン化粧品が、女性化粧品と男性化粧品の中間で相互連結している。

この結果から売場配置を考察すると、シャンプーは「ヘアケア製品」という視点ではなく、「インバス製品」と捉えて石鹸や入浴剤と近い売場配置することが購買者にとって利便性が高い可能性がある。また、女性化粧品と男性化粧品の間には整髪剤とヘアカラーやシーズン化粧品を配置すると、それらがクッションのような役割となる可能性がある。

筆者は、パーソナルケア群の担当バイヤーだった約7年間、店舗のリニューアル改装や新規店舗の出店の度にこの問題に延々と悩みながらフロアレイアウトの図面を引いていた。時間を遡ることができるならば、この分析結果を当時の筆者に届けたい。

6. おわりに

本稿では、膨大に列挙されたアソシエーション・ルールから、アイテム間の相互類似関係を用いてルール抽出を行う手法について紹介した。本手法は、一般的なアソシエーション分析手法と比べて、アイテムの網羅性が高く、全体を俯瞰しながら可読性の高いルールのネットワーク・グラフを得やすい。またその結果は、購買者の利便性を考慮した売場配置や関連販売などのビジネス応用に対して、さまざまな知見を得ることが可能である。本手法では、商品間の親しさを求めたいという考えから、相互類似関係を測る手段として簡単なランク情報という順序尺度を用いているが、それに対してはさまざまなアプローチがあろう。よりよい方法

を今後も考察していきたい。

また、世のデータ解析アルゴリズムの多くは「レンジでチン」をしたように、すぐに食べられる結果を出してくれるわけではない。データの前処理の重要性はよく知られることであるが、アルゴリズムの出力結果をいかに美味しく食べられる結果にまで仕上げるかの後処理も同じく重要だと考える。データを分析する者や分析結果を活用する者のかゆいところに手を届けるような、アルゴリズムの後処理に関する研究にもより目が向けられ、より活発に研究がなされることを期待したい。

謝辞 本稿の執筆にあたり、株式会社マクロミルより消費者購買履歴データ「QPR」を提供いただきました。また、本研究を学術研究として開始した直後から本稿の執筆に至るまで、中央大学の生田目崇教授に多くの助言をいただきました。深く感謝申し上げます。本研究はJST CRESTの研究助成を受けています。

参考文献

- [1] 岩崎幸子, 中元政一, 中原孝信, 宇野毅明, 羽室行信, “グラフ構造による相関ルールの視覚化ツール: KIZUNA,” 2017年度人工知能学会全国大会 (第31回), No. 2L4-2, 2017.
- [2] 中原孝信, 岩崎幸子, 中元政一, 宇野毅明, 羽室行信, “相互類似関係を用いたグラフ研磨の提案とその評価,” 2017年度人工知能学会全国大会 (第31回), No. 3O2-5, 2017.
- [3] M. Bastian, S. Heymann and M. Jacomy, “Gephi: An open source software for exploring and manipulating networks,” In *Proceedings of the Third International AAAI Conference on Weblogs and Social Media*, pp. 361-362, 2009.