

<原著論文>

テキストマイニング手法を用いた参与観察データの多角的視点分析

Analyze of participant observation data from the multi-angled viewpoints by text-mining

岩佐 真也¹, 大野 ゆう子²

要 旨

目的: 参与観察データをテキストマイニングで分析することにより、観察者が捉えた現象を多角的にみるための視点を明らかにする。

方法: セネガル共和国セレル民族の文化的価値観およびケアの意味と行為を18の領域に沿って記述的にまとめた参与観察データを対象とした。対象データをテキストマイニング手法により分析し、構成要素を抽出した。全構成要素と18領域の構成要素の高出現頻度語や、高頻度・低頻度に現れる特徴語の関連について検討した。

結果: テキストマイニングにより、609語の全構成要素を得た。その内、出現頻度が10回以上の構成要素数は162で、最も出現頻度の高いものは「子ども」の144回であった。18領域ごとの構成要素の内、最も構成要素が多かったのは《食習慣と栄養についての価値観、信念、タブー》で73であった。全構成要素の中で最も高頻度に出現した「子ども」は、領域5の《世界観》や領域12の《政治的または法的影響要因》、領域18の《仕事》で出現回数が0回であった。全構成要素における高頻度の出現頻度である構成要素が必ずしも各領域での高頻度の特徴語になっているわけではなく、低頻度の特徴語として現れているものもあった。

考察: 低頻度の特徴語は、調査者が各領域のデータを収集する際に見落としがなかったかといった注意喚起を与えてくれる可能性があると考えられる。またこの注意喚起は、その後のデータ収集の際の新たな視点として活用されることで、各領域だけでなく、データ全体像の捉えなおしにもつながると考えられた。

Abstract

Objectives

This study uses text mining to analyze participant observation data and identify the multi-angled viewpoints from which the observers perceive a topic.

Method

The participant observation data cover the cultural values of and the meaning and practice of “care” for the Serer ethnic group in the Republic of Senegal. The study divided the data into 18 descriptive classes and analyzed them by text mining to extract component words. It then examined the relationship between words that occur frequently and characteristic words that occur both frequently and infrequently in the data overall and in the 18 classes.

Results

Text mining produced 609 component words overall, of which 162 words occurred at least ten times, with the word “children” occurring most frequently, at 144 times. Of the 18 classes, that with the largest number of component words was “Values, beliefs, and taboos regarding diet and nutrition,” with 73 instances. “Children,” which occurred most frequently overall, occurred zero times in the classes “5. Worldview,” “12. Political or legal influences,” and “18. Work.” Thus, words that occurred frequently overall were not necessarily frequently occurring characteristic words in every class, but rather appeared as infrequently occurring characteristic words.

1 Maya IWASA 千里金蘭大学看護学部看護学科

受理日：2013年10月15日

2 Yuko OHNO 大阪大学大学院医学系研究科

査読付

Conclusion

Characteristic words that occurred infrequently may serve to caution researchers against oversight when collecting data by class. Exercising this caution is also thought to nurture new perspectives in future data collecting activities and help to interpret the data not only by class but also as a whole.

キーワード：参与観察，テキストマイニング，多角的，特徴語

Participant observation, Text mining, Perspective, Characteristic words

I. はじめに

テキストデータとは、文字コードだけで構成された文字列や文書のデータで、その収集方法はインタビューや自由記載、観察などにより行われる。この質的データには、発言者の特徴や、発言者に共通した普遍的な内容などが含まれていると考えられる¹⁾。一般的にこれらの質的データ分析には、共通性や差異性で分類や整理をし、カテゴリー化するような内容分析や、グランデッド・セオリー・アプローチのような継続的に設定したテーマ間とで比較分析をするものがある²⁾。他にも、文化人類学領域で用いられてきた、民族誌学的なアプローチによる分析などがある³⁾。これらはすべて、概念自体が容易に理解できないもしくは特定の概念を用いて説明ができないような場合に、現象や事象を発見していくための手法である。そのため、大量のデータの中に潜んでいる研究者が思いもよらない意味や個々のデータの関連性を発見しうるには、十分な広がりを持った理解ができたかという主観の適合性に注意を払わなければならない。主観的になりがちなデータ分析の妥当性確保は、今もなお質的研究の課題と考えられる⁴⁾。

そこで近年、注目されているのがテキストマイニングである。テキストマイニングは、データの中から必要な情報を取り出す自然言語処理である。大量の文章や文字に埋没している関係性を発掘するために、多角的にデータを分析し数量化できることが特徴である⁵⁻⁶⁾。

看護学領域におけるテキストマイニングの実際として、質問紙法による自由意見や感想文の分析、面接法による語りの内容分析などがある⁷⁻¹²⁾。これらはいずれも、数値情報や属性情報のような構造化されたデータではない質的データであり、これらのデータからキーワードを抽出しその使用頻度を解析することで、住民のニーズなどを明らかにしている。

このように、質的データをテキストマイニングにより数値化することで、もう一つの角度から現象を捉える試みがなされている。しかし、質的研究の代表的なデータ収集法である参与観察法によるデータを、テキストマイニングにより分析したものは筆者の知る限り見当たらない。

そこで、本研究では参与観察データを取りあげ、テキストマイニング手法を用いて観察者が捉えた現象を多角的にみるための視点を明らかにすることを目的とした。

II. 方法

1. 対象

参与観察データは、2009年3月から2009年6月の間に筆者が観察した、セネガル共和国セレール民族の文化的価値観およびケアの意味と行為を記述的にまとめたデータを用いた。

このデータの収集には、民俗学的参与観察法を用いた。これは、レイニンガーの自然でオープンな発見過程、および主に帰納的に過導き出されたイーミックな様式（人々が現象をどう意識・識別しているかを内側から分析するもの）を中心とするethno nursingを参考とした¹³⁾。

主な観察の領域は、一般的な生活環境、生活環境で使われている技術、世界観、家族の生活様式、宗教的信念と価値観、経済的要因、教育的価値観、民間的ヘルスケアの価値観、専門的なヘルスケアの価値観などの18である（表1）。

2. 分析方法

参与観察で得たデータをテキスト型データ解析ソフトWord Miner version1.15（日本電気計算株式会社）を用いて分析した。

得られたデータを形態素解析により文法的に意味

表1. 18領域

18領域	
1. 言語、コミュニケーション、身振り（自国語、外国語）	
2. 一般的な生活環境（シンボル、物質的および非物質的なサイン）	
3. 衣服と身体的外観	
4. 生活環境で使われている技術	
5. 世界観（世界をどのように見ているか）	
6. 家族の生活様式（価値観、信念、基準）	
7. 一般的な社会相互作用と親族関係の絆	
8. パターン化された日常生活活動	
9. 宗教的（または霊的）信念と価値観	
10. 経済的要因（生活費と所得の概算）	
11. 教育的価値観または信念の要因	
12. 政治的または法律的影響要因	
13. 食習慣と栄養についての価値観、信念、タブー	
14. 民間的（一般的またはその土地固有の）ヘルスケア（キュア）の価値観、信念、実践	
15. 専門的なヘルスケア（キュア）の価値観、信念、実践	
16. 行為の指標となるケアの概念またはパターン	
17. 病气予防の方法と健康維持の方法	
18. 仕事（内容、収入、価値観、信念、養成・研修）	

づけが最小限可能な形態素に分割する分ち書きを行った。次に、分ち書きをした構成要素（語句）について、それ自体には参与観察内容に関連しない品詞や助詞、助動詞、記号類、句読点などを除外した。また、同義語や類義語の統一化を行った。

その後、全構成要素と18領域ごとの構成要素を抽出し、それぞれの高出現頻度語を集計したものから全構成要素と各領域の構成要素の高出現頻度語の関連について分析した。

さらに、18領域ごとの構成要素から、ある領域において他の領域と比較して相対的に高頻度に現れた構成要素を特徴語として集計した。高頻度であっても他の領域でも同様に高頻度である構成要素は一般的な構成要素であり、特徴語とはみなさない。また、低頻度の構成要素も集計し、全構成要素と各領域の高出現頻度語と高頻度もしくは低頻度に現れる特徴語の関連について分析した。

Ⅲ. 結果

1. 構成要素の抽出と出現頻度

参与観察により得たデータの総文字数は41,781文字で、それらを分ち書きした結果、23,110語の形態素に分割された。不要な語句の除去や類義語などを統一し609語の全構成要素を得た。その内、出現頻度が10回以上の構成要素数は162で、20回以上の構成要素数は50であった。最も出現頻度の高いものは「子ども」で144回、次いで「人」が75回、「家」が64回、「挨拶」が60回、「仕事」が56回、「多い」が55回、「村」が54回、「女性」が49回だった（表2）。

全構成要素における上位20位までの構成要素の各領域内の内訳を見ると、「子ども」は18領域中13領

表2. 全構成要素における高頻度構成要素（閾値=20）

構成要素	構成要素数	構成要素	構成要素数
子ども	144	行く	31
人	75	自分	31
家	64	相手	31
挨拶	60	洗う	30
仕事	56	畑	29
多い	55	良く	27
村	54	コーラン	26
女性	49	健康	26
手	47	多く	26
水	47	マンガ	25
食べる	46	悪い	25
食事	46	現代	25
生活	46	収入	25
フランス語	43	話す	25
保健	43	無い	24
体	42	元気	23
家族	40	行って	23
男性	40	者	23
医療	36	小屋	23
学校	36	民族	22
彼	36	持って	21
病气	36	大人	21
言う	35	必要	21
薬	32	小学校	20
お金	31	昼食	20

域でも各領域の高頻度出現上位10位以内に入っていた。「人」では、18領域中6領域で、「家」では5領域、「挨拶」では3領域、「仕事」では2領域、「多い」では4領域、「村」では6領域、「女性」では5領域が各領域の高頻度出現上位10位以内に入っていた（表3）。

2. 全構成要素と領域別特徴語

18領域ごとの高頻度特徴語を見ると、領域5の《世界観》では「貨幣」が全構成要素数6中領域内構成要素数が6と《世界観》の領域の中で最も特徴的な構成要素だった。次いで「世界」で全構成要素数5中領域内構成要素数5、「精神」で全構成要素数7中領域内構成要素数5、「自然」で全構成要素数19中領域内構成要素数が7だった。領域12の《政治的または法律的影響要因》では「政治家」が全構成要素数5中領域内構成要素数5、「政治」では全構

表3. 領域別 高出現頻度語

順番	構成要素	構成要素数	領域																	
			1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
1	子ども	144	8**	7**	15**	3**	0	14**	1	23**	5**	3	20**	0	13**	6**	12**	7**	7**	0
2	人	75	19**	3	5**	0	2	0	9**	5	5**	0	3	0	4	8**	2	1	9**	0
3	家	64	2	16**	1	0	0	7**	4	8**	0	5**	4	0	5	3	0	1	6**	2
4	挨拶	60	26**	0	0	0	2	1	20**	0	0	0	0	0	0	0	0	1	10**	0
5	仕事	56	4	0	2	1	1	1	4	3	0	3	6**	0	0	2	3	0	2	24**
6	多い	55	8	3	5**	0	1	7**	2	8**	1	1	4	1	7**	1	3	1	1	1
7	村	54	9**	7**	1	0	5**	5**	6**	0	2	1	4	5**	2	1	4	1	1	0
8	女性	49	1	4**	19**	1	1	6**	0	1	0	5**	0	0	7**	0	0	0	1	3
9	手	47	8	0	0	4**	1	3	5**	6**	0	0	2	0	7**	2	1	0	7**	1
10	水	47	1	1	0	4**	0	0	0	7**	1	0	0	0	14**	8**	3	0	8**	0
11	食べる	46	1	1	0	0	0	4	0	6**	0	1	0	0	27**	0	1	0	5	0
12	生活	46	4	2	2	2	9**	4	6**	4	2	1	0	1	0	0	3	2	4	0
13	食事	46	1	1	0	0	1	3	0	1	3	1	0	0	26**	0	2	2	5	0
14	フランス語	43	14**	2	0	0	0	0	0	1	0	2	24**	0	0	0	0	0	0	0
15	保健	43	0	0	0	0	0	2	0	0	0	2	0	1	0	1	35**	0	1	1
16	体	42	0	1	1	0	1	0	0	3	0	0	0	0	13**	3	0	0	19**	1
17	男性	40	1	3	9**	0	2	9**	0	6**	0	3	0	0	2	1	0	0	0	4
18	家族	40	5	2	1	0	1	3	6**	2	1	2	0	0	3	0	4	4**	5	1
19	学校	36	1	4**	1	0	0	0	0	3	0	0	24**	0	1	0	0	0	0	2
20	彼	36	9**	0	0	0	2	1	7**	1	2	0	1	1	1	4	0	0	5	2

注釈：*は各領域内での高頻度上位10位内に入っていることを意味する。領域番号は表1に準ずる。

成要素数5中領域内構成要素数が4だった。他にも領域7の《一般的な社会相互作用と親族関係の絆》や領域14の《民間的ヘルスケアの価値観》、領域18の《仕事》の上位にも全構成要素数と領域内構成要素数との差が小さい、高頻度に現れる特徴語があった。

低頻度の特徴語を見ると、全構成要素の中で最も高頻度に出現した「子ども」は、領域5の《世界観》や領域7の《一般的な社会相互作用と親族関係の絆》、領域12の《政治的または法律的影響要因》、領域18の《仕事》で出現回数が0回であった。また「人」は領域6の《家族の生活様式》、領域10の《経済的要因》、領域18の《仕事》で出現回数が0回であった。これらは各領域の構成要素の内、低頻度の特徴語として抽出された(表4)。

全構成要素における高頻度の出現頻度である構成要素が必ずしも各領域での高頻度の特徴語になっているわけではなく、低頻度の特徴語として現れているものもあった。

IV. 考察

参与観察は、特定の社会集団について深い理解を得るために長期にわたり研究者自身がその社会集団に身を置き、そこで暮らす人々の生活やその行動様

式を見聞きするということから社会全体を捉えようとする調査手法である。そのため、収集されたデータの質は、アンケートやレポートなどの研究対象者の意見や価値観が直接明記された質的データとは違い、調査者の能力や調査者のバイアスに左右されやすい¹⁴⁾。

そのため、通常一人で行われるデータ収集の質を向上させるためには、データ収集の途中で幾度となくデータに向き合い主観の適合性について自省することで、十分な広がりをもて理解できたかという多角的視点を常に意識する必要がある¹⁵⁻¹⁶⁾。

そこで本研究では、参与観察データをテキストマイニングすることで、全構成要素における高頻度出現の構成要素が抽出でき、調査者が観察した現象の概観をつかむことに役立つことが考えられた。この抽出は、調査者がどのような言葉を用いて現象を把握しようとしたかが分かると同時に、ほとんど用いられなかった構成要素から、調査者のデータ収集時の領域の個性や調査対象地域の特徴の推察を支援することにつながるのではないかと考える。

一方、全構成要素で高頻度に現れている構成要素が、各領域でも高頻度に用いられている場合であっても、それらの構成要素が必ずしも各領域の特徴的な構成要素とはなっていない。このことは全体として概観を捉えることと、各領域での現象を分析

領域10 (21)				領域11 (34)				領域12 (14)			
順	構成要素	領域内 構成要素数	全構成 要素数	順	構成要素	領域内 構成要素数	全構成 要素数	順	構成要素	領域内 構成要素数	全構成 要素数
1	収入	13	25	1	学校	24	36	1	政治家	5	5
2	副業	6	7	2	フランス語	24	43	2	政治	4	5
高	3 マンゴー	11	25	3	コーラン	19	26	3	政府	4	7
頻	4 支出	5	6	4	教育	13	16	4	代表者	3	3
度	5 農作物	5	6	5	アラビア語	8	10	5	政策	3	4
特	6 現金	7	15	6	小学校	10	20	6	選挙	3	4
徴	7 農耕	7	15	7	幼稚園	5	5	7	ガンビア	2	2
語	8 5人	4	4	8	授業	7	12	8	コミュニテールラ	2	2
	9	24	3	9	親	7	13	9	期待	2	2
	10 歳男	3	3	10	3歳	4	4	10	対して	2	2
<hr/>				<hr/>				<hr/>			
10				10				10			
9				9				9			
低	8			8				8			
頻	7			7				7			
度	6			6	食べる	0	46	6			
特	5			5	食事	0	46	5			
徴	4			4	生活	0	46	4			
語	3 子ども	3	144	3	水	0	47	3			
	2 挨拶	0	60	2	女性	0	49	2			
	1 人	0	75	1	挨拶	0	60	1	子ども	0	144
<hr/>				<hr/>				<hr/>			
領域13 (73)				領域14 (44)				領域15 (40)			
順	構成要素	領域内 構成要素数	全構成 要素数	順	構成要素	領域内 構成要素数	全構成 要素数	順	構成要素	領域内 構成要素数	全構成 要素数
1	食べる	27	46	1	Marabout	14	14	1	保健	35	43
2	米	17	19	2	priere	14	14	2	小屋	19	23
高	3 食事	26	46	3	足	9	10	3	医療	19	36
頻	4 油	8	11	4	葉草	7	8	4	現代	16	25
度	5 砂糖	7	10	5	木	8	14	5	ポスト	12	16
特	6 朝食	9	17	6	痛い	5	6	6	葉	16	32
徴	7 魚	7	11	7	治す	4	4	7	錠剤	6	6
語	8 くず米	4	4	8	怪我	5	8	8	接種	7	10
	9 輸入	4	4	9	布	5	8	9	注射	5	5
	10 コーヒー	7	13	10	見て	4	5	10	予防	8	16
<hr/>				<hr/>				<hr/>			
10	お金	0	31	10				10			
9	行く	0	31	9				9			
低	8 葉	0	32	8				8	学校	0	36
頻	7 医療	0	36	7	家族	0	40	7	彼	0	36
度	6 病気	0	36	6	フランス語	0	43	6	男性	0	40
特	5 フランス語	0	43	5	食べる	0	46	5	体	0	42
徴	4 保健	0	43	4	食事	0	46	4	フランス語	0	43
語	3 生活	0	46	3	生活	0	46	3	女性	0	49
	2 仕事	0	56	2	女性	0	49	2	挨拶	0	60
	1 挨拶	0	60	1	挨拶	0	60	1	家	0	64
<hr/>				<hr/>				<hr/>			
領域16 (13)				領域17 (54)				領域18 (38)			
順	構成要素	領域内 構成要素数	全構成 要素数	順	構成要素	領域内 構成要素数	全構成 要素数	順	構成要素	領域内 構成要素数	全構成 要素数
1	病気	13	36	1	体	19	42	1	仕事	24	56
2	様子	6	13	2	洗う	13	30	2	カシュナツ	8	14
高	3 見る	5	8	3	傷	6	8	3	基礎	5	7
頻	4 関心	3	3	4	蚊帳	6	10	4	保健員	5	7
度	5 見守る	3	3	5	森	5	7	5	職業	6	13
特	6 受診	3	4	6	糞	3	3	6	運転手	4	5
徴	7 過ぎる	2	2	7	掃除	5	10	7	給料	4	5
語	8 我慢	2	2	8	健康	8	26	8	実	4	6
	9 指示	2	2	9	予防	6	16	9	朝食屋	4	6
	10 示す	2	2	10	ごしごし	3	4	10	タンガナ屋	3	3
<hr/>				<hr/>				<hr/>			
10				10				10			
9				9	女性	1	49	9			
低	8			8	行く	0	31	8	食べる	0	46
頻	7			7	子ども	7	144	7	食事	0	46
度	6			6	村	1	54	6	生活	0	46
特	5			5	多い	1	55	5	水	0	47
徴	4			4	医療	0	36	4	村	0	54
語	3			3	学校	0	36	3	挨拶	0	60
	2			2	男性	0	40	2	人	0	75
	1			1	フランス語	0	43	1	子ども	0	144

注釈：()*は構成要素数を表す。 領域番号は表1に準ずる。

することが全く同じではなく、各領域を見る際の特徴語という別の領域からも捉えなおすことの必要性を示唆していると思われる。

特に、低頻度の特徴語は、その領域には全く現れなかったもしくは少ししか現れなかった構成要素であり、逆の意味で特徴的な言葉とみなすことができる。このことから、低頻度の特徴語は、調査者が各領域のデータを収集する際に見落としがなかったかといった注意喚起を与えてくれる可能性があると考えられる。この注意喚起は、その後のデータ収集の際の新たな視点として活用されることで、各領域だけでなく、データ全体像の捉えなおしにもつながると考える。そのためにも、データ収集の途中でテキストマイニングを用いて、全構成要素と高頻度出現構成要素の関連や全構成要素と領域内構成要素から導き出される特徴語の関連を把握することが必要となる。また、途中でデータを振り返ることで、広がりのある多角的な視点を得ながらデータの質を確保し、データの柔軟性を高めることにつながると考える。

V. 結論

参与観察データを多角的にみるための視点とは、高頻度出現構成要素や全構成要素と領域内構成要素から導き出される特徴語の関連からデータを見ることである。またデータ収集の途中でこれらを繰り返し行う事で、新たな高頻度出現構成要素や特徴語が生まれ、データ全体が広がりを持つことにつながると考えられた。

引用文献

- 1) 筒井真優美, 太田有美, 渡邊久美子, 江本リナ, 甲斐恭子, 関根弘子, 中村明子, 日本における研究手法の変遷－量的研究・質的研究・トライアングレーション－, インターナショナルナーシングレビュー, 28(2), 37-51(2005)
- 2) 鈴木敦子, 質的・量的看護研究統合化への模索と質的データの分析, 四日市看護医療大学紀要, 1(1), 1-13(2008)
- 3) 米山俊直, 谷泰編, 文化人類学を学ぶ人のために, 世界思想社, 第9版, 32-45(1996)
- 4) 廣瀬春次, 混合研究法の現在と未来, 山口医学, 61(1・2合併号), 11-16(2012)
- 5) 藤井美和, 小杉考司, 李政元編著, 福祉・心理・看護のテキストマイニング入門, 中央法規, 14-27(2005)
- 6) Ah-Hwee Tan, Text Mining: The state of the art and the challenges, Proceedings of the PAKDD 1999 Workshop on Knowledge Discovery from Advanced Databases, 65-70 (1999)
- 7) 入江拓, 小平朋江, 看護大学生の精神科保護室に対する受け止めおよび視点の変化－テキストマイニングによる非構造型データの分析から－, 聖隷クリストファー大学看護学部紀要, 15, 1-10(2007)
- 8) 大高庸平, 城丸瑞恵, いたうたけひこ, 手術とホルモン療法を受けた乳がん患者の心理－テキストマイニングによる語りの分析から－, 昭和医学会誌, 70(4), 302-314(2010)
- 9) 入江拓, 横井麗子, 比嘉勇人, 精神看護実習をおこなう看護学生の眺める「風景」の視覚化－データマイニングとその活用－, 聖隷クリストファー大学看護学部紀要, 11, 35-48(2003)
- 10) 小林光樹, 齋藤ひろみ, 瀬川香子, 齋藤美華, テキスト解析によるキーワード抽出から見た地域看護活動の分析の試み, 東北大学医学部保健学科紀要, 15(1), 67-72(2006)
- 11) 原修一, 飯干紀代子, 山田弘幸, 天辰雅子, 中山翼, 大森史隆, 笠井新一郎, 言語聴覚士実習生の臨床実習への満足度に影響する要因－テキストマイニングによる検討－, 九州保健福祉大学研究紀要, 12, 149-155(2011)
- 12) 川島大輔, 小山達也, 川野健治, 伊藤弘人, 希死念慮者へのメッセージにみる, 自殺予防に対する意思の説明モデル－テキストマイニングによる分析－, パーソナリティ研究, 17(2), 121-132(2009)
- 13) Leininger, M. M., Culture Care Diversity & Universality: A Theory of Nursing, New York: National League for Nursing, 105-110 (1992)
- 14) A. ハルドン, S. ファン・デル・ヘースト他, 石川信克, 尾崎敬子監訳, 保健と医療の人類学－調査研究の手引き－, 世界思想社, 200-202 (2004)
- 15) 寺島秀明, 篠原徹編, 講座・生態人類学7, エスノ・サイエンス, 京都大学学術出版会, 7-11(2002)

- 16) 平野 (小原) 裕子, 多文化間精神医学的研究における社会調査 - 質的研究法と量的研究法を折衷する -, ころと文化, 6 (2), 108-116 (2007)