

はじめに 私はなぜ IBM SPSS Decision Trees なのか？

IBM SPSS Statistics のアドオン・モジュールである決定木(decision trees)分析を使った言語研究の方法を紹介します。決定木分析は、音韻、語彙、格助詞、文理解、敬語、ポライトネス、方言、言語習得、日本語教育など多様なテーマに活用することができます。IBM SPSS Decision Trees が、言語研究に有効な理由は以下の6つです。

第1に、決定木分析は、予測する目的変数(従属変数)が質的変数(名義尺度)である場合には分類木分析(classification trees)、量的変数(スケール)である場合は回帰木分析(regression trees)を自動的に選んで実行してくれます。

第2に、決定木分析は、量的変数と質的変数の両方を、目的変数を予測する説明変数(独立変数)として使用することができます。

第3に、決定木分析は、説明変数の予測力を統計検定によって評価してくれます。その際、目的変数に対して、説明変数の値が統計的に等質であると判断された場合は結合し、異質であると判断された場合は保持して、説明変数を最適に分類してくれます。

第4に、決定木分析は、CHAID またはその修正版である Exhaustive CHAID を使用することで、結果の樹形図を多分岐で描くことができます。

第5に、決定木分析は、説明変数で目的変数を予測した結果を、視覚的にわかり易い樹形図で描いてくれます。この樹形図のおかげで、難しい統計解析の手法がわからなくても、直感的に結果を理解することができます。

第6に、決定木分析は、1層目の説明変数を指定することができます。もちろん、有意でない説明変数を指定することは意味がないのですが、理論的に検証したい説明変数を1層目に置くことができるのは、研究目的の全貌を把握するのに非常に有効です。

この本の執筆にあたり、一橋大学の早川杏子先生がテキストの内容を確認してくださり、また宮崎大学の張婧禪先生が章のイラストなどのデザインをしてくださいました。心より感謝申し上げます。

2023年4月26日 名古屋大学にて
玉岡 賀津雄



| | |
|--------------------------------|----|
| はじめに | i |
| 第 1 章 決定木分析 | |
| —主語の有生性と動詞の自他性の関係..... | 1 |
| 1 決定木分析の目的 | 2 |
| 2 主語の有生性を予測する決定木分析..... | 4 |
| 3 樹形図の見方 | 10 |
| 4 決定木分析の種類 | 13 |
| 第 2 章 分類木分析 | |
| —マレーシア人日本語学習者のスピーチレベルシフト | 17 |
| 1 言語研究における決定木分析の目的..... | 18 |
| 2 会話での丁寧体と普通体の使用..... | 18 |
| 3 普通体と丁寧体の頻度記録とデータセット..... | 21 |
| 4 SPSS へのデータセットの読み込み..... | 22 |
| 5 ケースの重み付け | 26 |
| 6 目的変数と説明変数の決定 | 27 |
| 7 樹形図の枝の成長手法 | 30 |
| 8 樹形図の成長の制限 | 32 |
| 9 樹形図の描き方の選択 | 36 |
| 10 その他の出力の選択 | 37 |
| 11 モデルの要約, 相対リスク, 正解の割合 | 40 |
| 12 樹形図の描画..... | 42 |
| 13 樹形図の解釈 | 46 |
| 14 分類木分析で「最初の変数を適用」する信憑性..... | 52 |

| | |
|-------------------------------|----|
| 15 スピーチレベルシフト, 個人差, 男女差 | 55 |
| 16 森をみてから個々の木々をみる | 56 |
| 17 『12人の優しい日本人』の丁寧体使用 | 57 |

第3章 回帰木分析

| | |
|-----------------------------|----|
| — 行為要求表現における丁寧度の変化 | 61 |
| 1 行為要求表現と丁寧度の関係 | 62 |
| 2 行為要求表現の丁寧度の測定 | 62 |
| 3 丁寧度のデータセットの作成 | 65 |
| 4 データセットの読み込みと変数の尺度設定 | 69 |
| 5 回帰木分析の変数の設定 | 70 |
| 6 分岐基準の設定 | 72 |
| 7 樹形図の成長の制限 | 73 |
| 8 樹形図のターミナルノードと予測値の保存 | 74 |
| 9 回帰木分析のモデルの要約 | 76 |
| 10 樹形図の概要 | 78 |
| 11 樹形図による予測と相対リスク | 81 |
| 12 樹形図の決定係数を算出する2つの方法 | 83 |
| 13 樹形図の解釈 | 87 |
| 14 回帰木分析によって証明されたこと | 88 |

第4章 回帰木分析

| | |
|-------------------------------|-----|
| — 中国人日本語学習者による間接発話の理解 | 91 |
| 1 慣習的および非慣習的な間接発話 | 92 |
| 2 間接発話理解の先行研究からみいだせる問題点 | 93 |
| 3 日本語習熟度と間接発話の理解の測定 | 95 |
| 4 間接発話の理解を測定するテスト | 98 |
| 5 テスト項目の内容 | 101 |
| 6 日本語習熟度別の間接発話理解テストの結果 | 102 |

| | |
|--------------------------|-----|
| 7 間接発話の理解を予測する回帰木分析 | 103 |
| 8 回帰木分析の結果 | 107 |
| 9 樹形図の推定値 | 109 |
| 10 回帰木分析で「最初の変数を適用」する信憑性 | 111 |
| 11 樹形図の決定係数 | 113 |
| 12 回帰木分析からわかること | 114 |

第 5 章 分類木分析

| | |
|-------------------------|-----|
| —山口方言話者のアクセントにおける世代間の変化 | 119 |
| 1 アクセント核を担う単位と方言 | 120 |
| 2 山口方言と調査方法 | 121 |
| 3 撥音のアクセント核の頻度と独立性の検定 | 122 |
| 4 長音にアクセント核が置かれた頻度 | 130 |
| 5 促音にアクセント核が置かれた頻度 | 132 |
| 6 二重母音にアクセント核が置かれた頻度 | 133 |
| 7 アクセント核を予測する分類木分析 | 134 |
| 8 樹形図の解釈 | 139 |
| 9 分類木分析からわかること | 142 |

第 6 章 分類木分析

| | |
|------------------------|-----|
| —中国人および韓国人日本語学習者の連濁の習得 | 145 |
| 1 日本語の連濁とは | 146 |
| 2 連濁と語彙層 | 146 |
| 3 ライマンの法則 | 150 |
| 4 なぜ中国人と韓国人の日本語学習者なのか | 152 |
| 5 2つのグループの特性を統制する方法 | 155 |
| 6 性別 | 156 |
| 7 年齢（月齢） | 157 |
| 8 日本語学習期間 | 164 |

| | |
|------------------------|-----|
| 9 日本滞在期間 | 165 |
| 10 日本語の語彙知識 | 165 |
| 11 日本語の文法知識 | 170 |
| 12 分類木分析のデータセット | 173 |
| 13 分類木分析の実行 | 176 |
| 14 樹形図から窺えるライマンの法則の普遍性 | 182 |

第 7 章 分類木分析

| | |
|-------------------------|-----|
| — 絵本にみる疑問詞の習得順序 | 189 |
| 1 絵本の冊数をデータとすることの意味 | 190 |
| 2 疑問詞がみられる対象年齢別の絵本数 | 191 |
| 3 分類木分析のためのデータセット | 193 |
| 4 絵本の冊数頻度を使った分類木分析 | 197 |
| 5 樹形図からみえてくる幼児の疑問詞の習得順序 | 203 |

第 8 章 分類木分析

| | |
|--------------------------|-----|
| — 副詞と共起する接続助詞の文中・文末の出現頻度 | 207 |
| 1 コーパス研究とは | 208 |
| 2 副詞の共起する接続助詞の文中・文末での位置 | 210 |
| 3 コーパスの選択と頻度データ収集 | 212 |
| 4 共起頻度を分析するためのデータセット | 215 |
| 5 共起頻度の分類木分析 | 217 |
| 6 樹形図からみえてくる接続助詞の文中での位置 | 219 |

第 9 章 分類木分析と回帰木分析

| | |
|------------------------------------|-----|
| — 中国人日本語学習者による精神動詞および物理動詞の二格とヲ格の付与 | 225 |
| 1 動詞が付与する意味役割と格 | 226 |
| 2 3種類の活動動詞と記憶のテンプレート | 227 |
| 3 日本語学習者を対象とした二格とヲ格の習得 | 229 |

| | |
|----------------------|-----|
| 4 早川・玉岡・初(2015)研究の仮説 | 231 |
| 5 動詞の選択と条件の統制 | 232 |
| 6 格助詞テスト | 237 |
| 7 読解テストによる日本語能力の群分け | 238 |
| 8 ダミー変数による回帰木分析 | 240 |
| 9 ダミー変数による分類木分析 | 246 |
| 10 分類木分析の樹形図の結果 | 246 |
| 11 分類木分析からみえてくるもの | 248 |

第 10 章 回帰木分析

| | |
|---------------------------|-----|
| — 中国人日本語学習者の助言の難しさを決める諸要因 | 255 |
| 1 言うべきか, 言わざるべきか | 256 |
| 2 助言の定義と助言行動を決める諸要因 | 256 |
| 3 助言を巡る日中の社会文化差 | 259 |
| 4 中国語と日本語での助言の基準 | 261 |
| 5 助言場面の設定と質問内容 | 263 |
| 6 助言の難しさを予測する回帰木分析 | 265 |
| 7 回帰木分析からみえてくるもの | 275 |
| おわりに | 283 |
| 索引 | 286 |



第1章

決定木分析

主語の有生性と動詞の自他性の関係

決定木分析(けっていぎぶんせき、けっていほくぶんせき)は、ある目的の事柄(目的変数)があり、それを複数の背景要因(説明変数)で説明したり、あるいは構成する背景要因の構造を把握したりするための解析法です。分析の結果が樹形図で描かれるので、統計の知識がじゅうぶんではなくても、視覚的に結果を把握することができます。とても実用的な手法です。もともと、マーケットリサーチのために開発された解析法ですが、近年、医学、生物学、心理学などさまざまな分野で広く使われています。言語研究にも大いに活用できるはずです。この章では、主語の有生性を予測する研究を例に、決定木分析で描かれる樹形図の見方と名称を紹介します。そして、決定木分析が言語研究に効果的に応用できることを実感していただきます。決定木分析には、予測したい事柄のデータの違いによって、分類木分析と回帰木分析の2種類があります。最後に、これら2種類の決定木分析の特徴を説明します。

分析練習用のデータのダウンロード

<https://tamaoka.org/download/index.html> の jyogen.xls です。

注：このデータは分析の練習用に作成したものです。オリジナルのデータとは異なりますが、この章の分析結果と同じになります。



1 決定木分析の目的

決定木分析は、英語では decision tree analysis といわれます。英語の decision が「決定」「決断」「解決」「判断」などと訳されることからわかるように、なにかを決めることを目的とした分析法です。決定木分析は、もともとマーケティングリサーチの分野で開発されました。企業は、顧客の要求に応える商品やサービスを提供します。そして、それらの情報を顧客に効率的に流し、顧客を満足させることで利益を得ます。その際、どのような商品やサービスが市場で求められているか、どのような人達が商品やサービスを購入する見込みがあるか、さらに商品やサービスに対して顧客がどの程度満足しているのか、これらのことを効率よく把握しなくてはなりません。こうした一連のマーケティングリサーチにおける問いに答えるために決定木分析が活用されてきました。決定木分析は、結果を樹形図(dendrogram)で視覚的にわかり易く描いてくれます。そのため、統計の知識がなくても分析の結果が直感的に理解できるきわめて実用的な解析法です。

決定木分析は、予測あるいは説明したい事柄について、それを決める背景となる諸要因を解明することを目的としています。そのため、言語研究でも大いに活用することができるはずです。まず、簡単な例で考えてみましょう。日本語では、人間・動物が主語である場合には「いる」、それ以外には「ある」が用いられ、有生か無生かが区別されています。有生(animate)か無生(inanimate)かの区別を有生性(animacy)と呼びます。「先生が(生徒に)たずねる」の主語の「先生」、「犬が吠える」の主語の「犬」は有生です。一方、「山が聳(そび)える」の主語の「山」は無生の主語です。有生と無生の代わりに、有生物と無生物という言い方をすることもあります。また、これらはそれぞれ有情物と非情物、活動体と不活動体などという言い方をすることもあります。ヨーロッパ言語では、無生物が主語になることが多いといわれます。一方、日本語では、主語になるのは有生の名詞のほうが無生の名詞より多いというのが一般的な理解ではないでしょうか。

私達が典型的だと考えるのは、有生の誰かがなにかを行うという他動詞文でしょう。たとえば、図1-1のように、「ケンジとマサルがカレーを作った。」という文です。主語は「ケンジとマサル」で、主語を示す格助詞であ



分類木分析

マレーシア人日本語学習者の スピーチレベルシフト

分類木分析(ぶんるいぎぶんせき、ぶんるいぼくぶんせき)は、決定木分析の一種です。ある複数の**説明変数**(独立変数)で名義尺度の目的変数(従属変数)を予測する分析法です。カイ2乗検定の有意水準で樹形図の分岐が決められます。樹形図の分岐をたどることで目的変数を予測する諸要因(説明変数)を視覚的に把握することができます。この章では、マレーシア人日本語学習者が初対面の目上の人と話す際の、日本に来る前と日本に来てからの普通体と丁寧体の使用頻度を調べた研究(Jamila, 2008)を紹介します。統計分析ソフト IBM SPSS Decision Trees による普通体と丁寧体を予測する分類木分析の手順を説明します。

分析練習用のデータのダウンロード

<https://tamaoka.org/download/index.html> の speech.xls

注：このデータセットは、Jamila(2006, 2007, 2008)に印刷されている頻度表から作成しました。

<https://tamaoka.org/download/index.html> の baishin.xls

注：このデータセットは、日高・伊藤(2007)に印刷されている頻度表から作成しました。



1 言語研究における決定木分析の目的

目的変数が名義尺度である場合に分類木分析を使います。説明変数としては、名義尺度、順序尺度、量的尺度のいずれも使えます。第1章では、主語が有生または無生のいずれをとるかという仮説を分類木分析で検証しました。「上がる」「変わる」「始まる」の3つの自動詞では、毎日新聞18年の記事を調べたところ、出現頻度(以下、頻度)で有生の名詞が主語になることはあまりなく、無生の名詞が主語になる場合がほとんどでした。一方、対応する「上げる」「変える」「始める」の3つの他動詞では有生と無生の主語の頻度に違いはみられませんでした。ただし、第1章で紹介した6つの動詞はあくまで分類木分析を説明するために使用した頻度データです。有生性と自他性との関係については、64種類の動詞を対象に研究した玉岡・張・牧岡(2019)を参照してください。なお、新聞のサンプル数は約3億語で、非常に大きいので、玉岡・張・牧岡(2019)の結果の予測力はじゅうぶんであると思われます。そのため、今後書かれる新聞記事においても、これらの動詞について類似した主語の有生性の頻度傾向を予測することができると考えられます。

すでに、決定木分析の全体的なイメージは、第1章でつかめたと思います。そこで、第2章では、IBM SPSS Statisticsのアドオン・ソフトであるDecision Treesを使って分類木分析を実際に行ってみましょう。以下では、(1)Excelを使ってどのようにデータを配置するか、(2)SPSSを使ってどのような手順で分析するか、(3)結果をどう解釈するか、を説明します。



2 会話での丁寧体と普通体の使用

分類木分析を行う前に、どのような仮説を証明しようとしているかを把握しておくことが大切です。そこで、まずデータ収集の理論的な背景について説明します。

日本語には体系的な敬語表現が存在します(菊池, 1997)。日本語の会話では、同じ内容であっても、対人関係、場面、状況、目的などに応じて表現を変えることがあります。「見る」と表現するにしても、尊敬語の「ご覧になる」、謙譲語の「拝見する」、丁寧語の「見ます」と、いろいろな表現があり

第3章

回帰木分析

行為要求表現における丁寧度の変化

回帰木分析(かいきぎぶんせき, かいきぼくぶんせき)は決定木分析の一種で、複数の説明変数(独立変数)で量的尺度の目的変数(従属変数)を予測する解析法です。回帰木分析では、 F 検定で平均の差を比較して、樹形図を分岐します。日本語の行為要求表現は、平叙形の肯定の「～てくれる」よりも否定の「～てくれない」のほうが、さらに疑問形でも肯定の「～てもらえますか」より否定の「～てもらえませんか」のほうがより丁寧だといわれています。この章では、肯定形か否定形か、平叙形か疑問形か、という言語特性によって丁寧度がどう変わるかを主観的判断で検討した研究(林・玉岡・宮岡, 2011)を例に、SPSSのDecision Treesを使って回帰木分析の手順を解説します。

分析練習用のデータのダウンロード

<https://tamaoka.org/download/index.html> の polite.xls

注：このデータは分析練習用に作成したものです。オリジナルのデータとは異なりますが、この章の分析結果と同じになります。



1 行為要求表現と丁寧度の関係

はじめに、この章の回帰木分析に使うデータの背景について説明します。相手になんらかの行動を起こさせることを目的とする**行為要求表現**(国立国語研究所, 1992)は、「話して」「話しなさい」と直接に要求する形式から、「話してくれる」「話していただけませんか」といった疑問形で間接的に要求する形式まで多様です。表現の形式としては、授受の補助動詞「くれる／ください」や「もらう／いただく」を含んだ表現が頻繁に使用されます(岡本, 2007)。これらは、「くれる」<「もらう」<「くださる」<「いただく」と相手の許可を求める際の丁寧度も変化すると予想されます(‘<’は丁寧度の順位を示す)。これに加えて、ストレートに要求するのではなく、「～てもらえますか」と相手に可能かどうかとたずねることで、より丁寧さが増すとされています(岡本, 1986, 2000; Okamoto, 1992)。さらに、日本語の要求表現は、「～いただけますか」の肯定の疑問形よりも、「～いただけませんか」の否定の疑問形のほうが丁寧であるといわれています。その理由としては、依頼の相手に対して選択の余地を与えることで、心理的負担を軽減し、より丁寧な印象を与えるからだと説明されています(菊池, 1997; 岡本, 2007)。これらの先行研究(菊池, 1997; 岡本, 1986, 2000; Okamoto, 1992)から、日本語の行為要求表現において、否定や疑問などの言語特性が丁寧度を上げる要因となっているかどうかを実証的に検討したのが、林・玉岡・宮岡(2011)の研究です。この章では、林ほか(2011)の研究で使われたデータを基に、回帰木分析の手順を説明します。



2 行為要求表現の丁寧度の測定

依頼される側である聞き手が感じる丁寧度は、どのように測定するのでしょうか。表3-1に示したように、林ほか(2011)は、平叙形の「～してくれる」「～もらえる」から、疑問形の「～てもらえますか」「～くださいますか」「～いただけますか」という丁寧度が異なると予想される行為要求表現を設定しています。さらに、これらの5つの表現を「～てくれない」「～てもらえない」「～てもらえませんか」「～くださいませんか」「～ていた