

**IBM SPSS Modeler Text
Analytics 18.1.1 ユーザーズ・
ガイド**

IBM

注記

本書および本書で紹介する製品をご使用になる前に、237 ページの『特記事項』に記載されている情報をお読みください。

本書は、IBM SPSS Modeler Text Analytics バージョン 18.1.1 リリース 0 モディフィケーション 0 および新しい版で明記されない限り、以降のすべてのリリースおよびモディフィケーションに適用されます。

お客様の環境によっては、資料中の円記号がバックスラッシュと表示されたり、バックスラッシュが円記号と表示されたりする場合があります。

原典： IBM SPSS Modeler Text Analytics 18.1.1 User's Guide

発行： 日本アイ・ビー・エム株式会社

担当： トランスレーション・サービス・センター

目次

前書き	vii
IBM Business Analytics について	vii
技術サポート	vii
第 1 章 IBM SPSS Modeler Text Analytics について	1
IBM SPSS Modeler Text Analytics バージョン	
18.1.1 へのアップグレード	1
テキスト マイニングについて	2
抽出の方法	5
カテゴリー化の方法	7
IBM SPSS Modeler Text Analytics ノード	8
アプリケーション	9
第 2 章 ソース・テキストの読み取り	11
ファイル リスト ノード	11
ファイル リスト ノード:「設定」タブ	12
ファイル リスト ノード:その他のタブ	13
テキスト マイニングでのファイル リスト ノードの使用	13
Web フィード ノード	13
Web フィード ノード:「入力」タブ	14
Web フィード ノード:「レコード」タブ	15
Web フィード ノード:「コンテンツ フィルター」タブ	16
テキスト マイニングでの Web フィード ノードの使用	17
言語ノード	17
言語ノード:「設定」タブ	18
第 3 章 コンセプトおよびカテゴリーのマイニング	19
テキスト マイニング・モデル作成ノード	20
テキスト マイニング・ノード:「フィールド」タブ	21
テキスト マイニング・ノード:「モデル」タブ	24
テキスト マイニング・ノード:「エキスパート」タブ	28
時間を削減する上流のサンプリング	30
ストリーム内のテキスト マイニング・ノードの使用	30
テキスト マイニング モデル ナゲット:コンセプト・モデル	31
コンセプト・モデル:「モデル」タブ	32
コンセプト モデル:「設定」タブ	34
コンセプト モデル:「フィールド」タブ	35
コンセプト・モデル:「要約」タブ	36
ストリームでのコンセプト モデル ナゲットの使用	36

テキスト マイニング モデル ナゲット:カテゴリー・モデル	40
カテゴリー モデル ナゲット:「モデル」タブ	41
カテゴリー モデル ナゲット:「設定」タブ	42
カテゴリー モデル ナゲット:その他のタブ	43
ストリームでのカテゴリー モデル ナゲットの使用	44
第 4 章 テキスト・リンクのマイニング	47
テキスト リンク分析ノード	47
「テキスト リンク分析ノード」:「フィールド」タブ	48
テキスト リンク分析ノード:「エキスパート」タブ	49
TLA ノード出力	51
TLA 結果のキャッシュ	51
ストリーム内のテキスト リンク分析ノードの使用	52
第 5 章 外部ソース・テキストの参照	55
ファイル ビューアー ノード	55
ファイル ビューアー ノード設定	55
ファイル ビューアー ノードの使用	56
第 6 章 スクリプト用のノードのプロパティ	59
ファイル リスト ノード:filelistnode	59
Web フィード ノード:webfeednode	59
言語ノード: languageidentifier	60
テキスト マイニング・ノード	
:TextMiningWorkbench	61
テキスト マイニング モデル ナゲット	
:TMWBModelApplier	63
テキスト リンク分析ノード: textlinkanalysis	64
第 7 章 インタラクティブ・ワークベンチ・モード	67
カテゴリーとコンセプト・ビュー	68
クラスター・ビュー	70
テキスト リンク分析ビュー	72
リソース・エディター・ビュー	74
オプションの設定	76
オプション:「セッション」タブ	76
オプション:「表示」タブ	76
オプション:「サウンド」タブ	77
Microsoft Internet Explorer ヘルプの設定	77
モデル ナゲットおよびモデル作成ノードの生成	78
モデル作成ノードの更新および保存	78
セッションの終了	78
キーボード・アクセシビリティ	79
ダイアログ・ボックスのショートカット	80

第 8 章 コンセプトとタイプの抽出	81
抽出結果: コンセプトとタイプ	81
データの抽出	82
抽出結果のフィルタリング	85
コンセプト・マップの検証	86
コンセプト・マップ・インデックスの作成	89
抽出結果の調整	89
類義語の追加	90
コンセプトのタイプへの追加	91
コンセプトの抽出からの除外	92
単語を抽出に強制投入	93
第 9 章 テキストデータの Kategorisierung	95
カテゴリー・ペイン	96
カテゴリー作成の方法と戦略	98
カテゴリー作成の方法	98
カテゴリー作成の方略	99
カテゴリー作成のヒント	99
最適な記述子の選択	100
カテゴリーとは	103
カテゴリーのプロパティ	104
データ・ペイン	104
カテゴリーの関連性	105
カテゴリーの作成	106
言語学的手法の詳細設定	108
言語学的手法について	110
出現頻度に基づく手法の詳細設定	115
カテゴリーの拡張	116
手作業でのカテゴリーの作成	119
カテゴリーの新規作成または名前の変更	119
ドラッグ&ドロップによるカテゴリーの作成	120
カテゴリー規則の使用	121
カテゴリー規則シンタックス	121
カテゴリー規則内の TLA パターンの使用	123
カテゴリー規則におけるワイルドカードの使用	125
カテゴリー規則の例	127
カテゴリー規則の作成	129
規則の編集および削除	130
定義済みのインポートおよびエクスポート	130
定義済みカテゴリーのインポート	131
カテゴリーのエクスポート	135
テキスト分析パッケージの使用	136
テキスト分析パッケージの作成	136
テキスト分析パッケージの読み込み	137
テキスト分析パッケージの更新	138
カテゴリーの編集および調整	139
記述子のカテゴリーへの追加	139
カテゴリー記述子の編集	139
カテゴリーの移動	140
カテゴリーのフラット化	140
カテゴリーの結合・組み合わせ	141
カテゴリーの削除	141
第 10 章 クラスターの分析	143
クラスターの作成	144
類似度リンク値の計算	145

クラスターの検証	146
クラスター定義	147
第 11 章 テキスト リンク分析の検証	149
TLA パターン結果の抽出	150
タイプ・パターンおよびコンセプト・パターン	151
TLA 結果のフィルタリング	152
データ パネル	153
第 12 章 グラフの視覚化	155
カテゴリー・グラフおよび図表	155
カテゴリー棒グラフ	156
カテゴリー Web グラフ	156
カテゴリー Web テーブル	157
クラスター・グラフ	157
コンセプト Web グラフ	157
クラスター Web グラフ	158
テキスト リンク分析のグラフ	158
コンセプト Web グラフ	159
タイプ Web グラフ	159
グラフのツールバーおよびパレットの使用	159
第 13 章 セッション・リソース・エディター	163
リソース・エディターを使用したリソースの編集	163
テンプレートの作成および更新	164
リソース・テンプレートの切り替え	165
第 14 章 テンプレートとリソース	167
テンプレート エディターとリソース エディターの比較	168
エディターのインターフェース	168
テンプレートを開く	172
テンプレートの保存	173
読み込み後のノード・リソースの更新	173
テンプレートの管理	174
テンプレートのインポートおよびエクスポート	175
テンプレート・エディター の終了	175
リソースのバックアップ	176
リソース・ファイルのインポート	176
第 15 章 ライブラリーの使用	179
付属ライブラリー	179
ライブラリーの作成	180
パブリック・ライブラリーの追加	181
キーワードおよびタイプの検索	181
ライブラリーの表示	182
ローカル・ライブラリーの管理	182
ローカル・ライブラリーの名前の変更	182
ローカル・ライブラリーを使用不可に	183
ローカル・ライブラリーの削除	183
パブリック・ライブラリーの管理	183
ライブラリーの共有	184
ライブラリーの公開	186
ライブラリーの更新	186
競合の解決	187

第 16 章 ライブラリー辞書について	189
キーワード辞書	189
ビルトインのタイプ	190
キーワード辞書の作成	191
キーワードの追加	192
キーワードの強制	195
キーワード辞書の名前変更	195
キーワード辞書の移動	196
キーワード辞書の無効化および削除	196
類義語辞書	197
類義語の定義	198
オプションの要素の定義	199
類義語の無効化および削除	199
不要語辞書	200
第 17 章 拡張リソースについて	203
検索	204
置換	204
リソースの対象言語	205
Fuzzy Grouping	205
固有表現	206
正規表現の定義	207
正規化	209
構成	210
言語処理	211
抽出パターン	211
強制定義	213
省略形	214

第 18 章 テキスト リンク規則について	215
テキスト リンク規則を扱う場所	215
作業の開始	216
規則の編集または作成が必要な場合	216
テキスト リンク分析結果のシミュレーション	217
シミュレーションのデータ定義	217
シミュレーション結果の理解	218
ツリー内の規則およびマクロのナビゲート	219
マクロの作業	220
マクロの作成および編集	221
マクロの無効化および削除	221
エラーのチェック、保存およびキャンセル	222
特殊マクロ：mTopic、mNonLingEntities、SEP	223
テキスト リンク規則の使用	223
条件規則の作成および編集	226
条件規則の無効化および削除	227
エラーのチェック、保存およびキャンセル	227
条件規則の処理順序	228
ルール・セットの使用 (多段階処理)	229
条件規則およびマクロにサポートされている要素	230
入力モードでの表示および作業	232
特記事項	237
商標	238
索引	239

前書き

IBM® SPSS® Modeler Text Analytics は強力なテキスト分析機能を提供するものであり、高度な言語テクノロジーと自然言語処理 (Natural Language Processing, NLP) を使用して、さまざまな構造のないテキスト データを高速で処理し、このテキストから重要なコンセプトを抽出および整理します。さらに、IBM SPSS Modeler Text Analytics はこれらのコンセプトをカテゴリーにグループ化できます。

組織内に保持されるおよそ 80% のデータは、テキスト ドキュメントの形式です (例: レポート、Web ページ、電子メール、コール センターのメモ)。テキストは、組織が顧客の動向をより良く理解するための重要な要素です。NLP を組み込むシステムは、複合句などのコンセプトを効率的に抽出できます。さらに、基底となる言語の情報を使用して、コンセプトを製品、組織、人物など、意味や状況に応じて関連グループに分類できます。その結果、情報のニーズに対する関連性を迅速に確認できます。これらの抽出されたコンセプトとカテゴリーは、人口統計など既存の構造化されたデータと組み合わせることができ、さらに IBM SPSS Modeler の完全なデータ マイニング ツールを使ったモデル作成に適用することにより、より適切で焦点を絞った決定を行うことができます。

言語学的なシステムは、知識に依存します。つまり、辞書に含まれている情報が多いほど、より高い品質の結果が得られます。IBM SPSS Modeler Text Analytics には、キーワードや類義語の辞書、ライブラリー、およびテンプレートなど、一連の言語リソースが付属しています。またこの製品を使用すると、状況に合わせてこれらの言語リソースを開発および調整できます。言語リソースの調整はインタラクティブなプロセスで、正確なコンセプトの取得とカテゴリー化に必要です。CRM およびゲノムなど、特定のドメインのカスタム テンプレート、ライブラリー、辞書も含まれています。

IBM Business Analytics について

IBM Business Analytics ソフトウェアは、意思決定者がビジネス・パフォーマンスの改善のために使用可能な完全で整合性があり、正確な情報を提供します。ビジネス・インテリジェンス、予測分析、財務実績および戦略管理、分析アプリケーション の包括的なポートフォリオを利用することによって、現在の実績を明確、迅速に理解し、将来の結果を予測することができます。豊富な産業用ソリューション、証明された実践法、それに専門家によるサービスを組み合わせることにより、あらゆる規模の会社組織が、最高の生産性を推進し、信頼できる意志決定を自動化し、そして、よりよい結果を実現させることができます。

このポートフォリオの一部として、IBM SPSS Predictive Analytics ソフトウェアを使用する組織は、将来のイベントを予測し、その洞察に基づいて積極的に行動し、より優れた業績を実現することができます。全世界の企業、政府、学術分野のお客様が IBM SPSS の技術を活用し、不正行為を減少させ、リスクを軽減させながら、顧客の獲得、保持、成長において、競争優位を高めることができます。IBM SPSS ソフトウェアを日々の業務に取り入れることによって、組織は業務目標を達成し、大きな競争的優位を獲得することができるよう、意思決定を方向付け、自動化することができるようになります。詳細な情報、または営業担当者へのお問い合わせ方法については、<http://www.ibm.com/spss> を参照してください。

技術サポート

お客様はテクニカル・サポートをご利用いただけます。IBM Corp. 製品の使用方法、または対応するハードウェア環境へのインストールについてサポートが必要な場合は、テクニカル・サポートにご連絡ください。テクニカル・サポートの詳細は、IBM Corp. Web ページ <http://www.ibm.com/support> を参照してください。ご依頼の際には、氏名、組織名、およびサポート同意契約をご用意ください。

第 1 章 IBM SPSS Modeler Text Analytics について

IBM SPSS Modeler Text Analytics は強力なテキスト分析機能を提供するものであり、高度な言語テクノロジーと自然言語処理 (Natural Language Processing, NLP) を使用して、さまざまな構造のないテキストデータを高速で処理し、このテキストから重要なコンセプトを抽出および整理します。さらに、IBM SPSS Modeler Text Analytics はこれらのコンセプトをカテゴリーにグループ化できます。

組織内に保持されるおよそ 80% のデータは、テキスト ドキュメントの形式です (例: レポート、Web ページ、電子メール、コール センターのメモ)。テキストは、組織が顧客の動向をより良く理解するための重要な要素です。NLP を組み込むシステムは、複合句などのコンセプトを効率的に抽出できます。さらに、基底となる言語の情報を使用して、コンセプトを製品、組織、人物など、意味や状況に応じて関連グループに分類できます。その結果、情報のニーズに対する関連性を迅速に確認できます。これらの抽出されたコンセプトとカテゴリーは、人口統計など既存の構造化されたデータと組み合わせることができ、さらに IBM SPSS Modeler の完全なデータ マイニング ツールを使ったモデル作成に適用することにより、より適切で焦点を絞った決定を行うことができます。

言語学的なシステムは、知識に依存します。つまり、辞書に含まれている情報が多いほど、より高い品質の結果が得られます。IBM SPSS Modeler Text Analytics には、キーワードや類義語の辞書、ライブラリー、およびテンプレートなど、一連の言語リソースが付属しています。またこの製品を使用すると、状況に合わせてこれらの言語リソースを開発および調整できます。言語リソースの調整はインタラクティブなプロセスで、正確なコンセプトの取得とカテゴリー化に必要です。CRM およびゲノムなど、特定のドメインのカスタム テンプレート、ライブラリー、辞書も含まれています。

展開: 非構造化データのリアルタイムのスコアリングに IBM SPSS Modeler Solution Publisher を使用して、テキスト マイニング・ストリームを展開できます。これらのストリームを展開する機能により、正常でクローズドループのテキスト マイニングの実装を実現します。例えば、組織は、予測モデルを適用してマーケティング メッセージの精度をリアルタイムに向上させることにより、受信者または発信者からのメモ帳のメモを分析することができます。

IBM SPSS Modeler Text Analytics を IBM SPSS Modeler Solution Publisher とともに実行するには、ディレクトリー <install_directory>/ext/bin/spss.TMWBServer を \$LD_LIBRARY_PATH 環境変数に追加します。

注: IBM SPSS Modeler Text Analytics の日本語アダプターは、バージョン 18.1 以降は推奨されていません。

IBM SPSS Modeler Text Analytics バージョン 18.1.1 へのアップグレード

旧バージョンの PASW Text Analytics または Text Mining for Clementine からのアップグレード。

IBM SPSS Modeler Text Analytics バージョン 18.1.1 をインストールする前に、新しいバージョンで使用する現在のバージョンの TAP、テンプレート、ライブラリーを保存してエクスポートする必要があります。これらのファイルは、最新バージョンをインストールしても削除も上書きもされないディレクトリーに保存することをお勧めします。

最新バージョンの IBM SPSS Modeler Text Analytics をインストールした後、保存した TAP ファイルの読み込み、保存したライブラリーの追加、保存したテンプレートのインポートおよび読み込みを行って、最新バージョンで使用することができます。

重要: 最初に必要なファイルを保存してエクスポートせずに現在のバージョンをアンインストールすると、旧バージョンで実行していた TAP、テンプレート、パブリック ライブラリーの作業が失われ、IBM SPSS Modeler Text Analytics バージョン 18.1.1 で使用できなくなります。

テキスト マイニングについて

現在、顧客の電子メール、コール センターのメモ、自由記述式のアンケート回答、ニュース フィード、Web フォームなど、非構造化または半構造の形式で保持された情報量が増加しています。この情報過多によって、多くの組織に「この情報をどのように収集、検証そして活用するのか」という問題をもたらします。

テキスト マイニングとは、テキスト形式の素材のコレクションを分析するプロセスで、作者がこれらのコンセプトの表現に使用した正確な単語またはキーワードを知らなくても、主要なコンセプトやテーマをキャプチャーし、隠れた関連性や傾向を明らかにします。テキスト マイニングと情報検索は全く異なりますが、これらが混同される場合があります。情報の正確な検索および保存は大きな課題ですが、情報に含まれる高品質な内容、用語、および関連性の抽出および管理は非常に重要なプロセスです。

テキスト マイニングおよびデータ・マイニング

テキストの各項目について、言語学的テキスト マイニングによりコンセプトのインデックス、およびこれらのコンセプトについての情報を返します。この抜き出された、構造化された情報は、その他のデータ・ソースと組み合わせて、次のような質問を処理することができます。

- 一緒に出現するのはどのコンセプトですか？
- コンセプトが他に何かリンクしているものがありますか？
- 抽出した情報から作成できる高レベルのカテゴリーは何ですか？
- コンセプトまたはカテゴリーから予測するのは何ですか？
- コンセプトまたはカテゴリーからどのように動作を予測しますか？

テキスト マイニングとデータ・マイニングを組み合わせると、構造化データまたは非構造化データだけで行うよりも、すぐれた洞察が可能です。この処理には通常、次のステップが含まれます。

1. マイニングするテキストを特定する。 マイニングするテキストを準備します。テキストが複数のファイルにある場合、ファイルを 1 つの場所に保存します。データベースについては、テキストが含まれているフィールドを決定します。
2. テキストをマイニングして構造化データを抽出する。 テキスト マイニング アルゴリズムをソース・テキストに適用します。
3. コンセプト モデルおよびカテゴリー モデルを作成する。 主要なコンセプトを特定し、カテゴリーを作成します (あるいは、そのいずれか)。非構造化データから返されるコンセプト数は通常、非常に多くなります。スコアリングに最適なコンセプトおよびカテゴリーを特定します。
4. 構造化データを分析する。 クラスターリング、分類、予測モデル作成など、従来のデータ マイニング手法を採用して、コンセプト間の関連性を検出します。抽出されたコンセプトを他の構造化データに結合し、コンセプトに基づいて今後の動作を予測します。

テキスト分析およびカテゴリー化

定性的分析の形式であるテキスト分析では、テキストからの役立つ情報を抽出し、このテキスト内の主要なアイデアまたはコンセプトを適切な数のカテゴリーにグループ化します。テキスト分析はすべての種類および長さのテキストに実行できますが、分析へのアプローチは若干異なります。

比較的短いレコードまたはドキュメントは、それほど複雑でなく、通常不明確な単語や回答があまり含まれていないため、最も容易にカテゴリー化されます。例えば、短い自由記述式のアンケートで好きな休日の過ごし方を 3 つ挙げるよう質問した場合、ビーチに行く、国立公園に行く、または何もしない などの多くの短い回答が見られることが予想される場合があります。一方、比較的長い自由記述式のアンケートの回答は、特に回答者が高学歴で意欲があり、またアンケートを記入するのに十分な時間がある場合、非常に複雑で長くなる場合があります。アンケートで政治に関する考えを尋ねたり、または政治に関するブログ フィードがあったりする場合、あらゆる種類の問題および立場について、長いコメントがいくつかあると予想されることがあります。

非常に短い時間で長いテキスト・ソースから主要キーワードを抽出して洞察に満ちたカテゴリーを作成する機能は、IBM SPSS Modeler Text Analytics を使用するうえでの重要な利点です。この利点は、自動化された言語学的手法と統計的手法を組み合わせ得られるもので、テキスト分析プロセスの段階ごとに最も信頼できる結果を生成します。

言語処理および NLP

すべての構造のないテキスト・データの管理における主な問題は、コンピューターが理解できるようなテキストを作成するための標準的な規則がないという点です。言語、すなわち意味はすべてのドキュメントおよびすべてのテキストの部分で異なります。そのような非構造化データを正確に取得し構成する唯一の方法は、言語を分析してその意味を明らかにすることです。非構造化情報からコンセプトを抽出するには、いくつかの異なる自動化されたアプローチがあります。これらのアプローチは、言語学的アプローチと非言語学的アプローチの 2 種類に分けられます。

いくつかの組織が、統計およびニューラル・ネットワークに基づく自動化された非言語学的ソリューションを採用しようとしてきました。これらのソリューションでは、コンピューター技術を駆使して、人間が読み込むよりはるかに迅速に主要キーワードをスキャンおよびカテゴリー化できます。しかし、こうしたソリューションの精度は非常に低くなります。多くの統計的システムでは、単語が出現する回数をただカウントし、関連するコンセプトへの統計的近接性を計算するだけです。これにより関連性の低い多くの結果、すなわちノイズを生み出し、見つけるべき結果や無視すべき結果を見逃したりすることになります。

限られた精度を補うために、いくつかのソリューションで複雑な非言語学的規則を組み込み、関連性のある結果および関連性のない結果とを区別します。これを、規則に基づくテキスト マイニングといいます。

一方、言語学に基づくテキスト マイニングでは、人間の言語をコンピューターによる支援で分析する自然言語処理 (NLP) の原則をテキストの単語、句、構文、または構造に適用します。NLP を組み込むシステムは、複合句などのコンセプトを効率的に抽出できます。さらに、基底となる言語の情報を使用して、コンセプトを製品、組織、人物など、意味や状況に応じて関連グループに分類できます。

言語学に基づくテキスト マイニングでは、さまざまな単語の形式が類似した意味を持っていることを認識し、文の構造を分析してテキストを理解するための枠組みを提供することによって、人間と同じようにテキストの意味を検出します。このアプローチでは、統計的システムの速度およびコストの効率の点を利用し、人間の手をほとんど必要とせず、精度がはるかに高くなります。

抽出プロセス時における統計的アプローチと言語学的アプローチとの違いを説明するために、reproduction of documents (ドキュメントの複製) についての質問に対する回答について考えてみましょう。統計的ソリ

ューションおよび言語学的ソリューションのいずれも、reproduction (複製)という単語を展開して、copy (コピー)やduplication (重複)などの類義語を含めるようにする必要があります。展開しない場合、関連情報が見落とされてしまいます。ただし、統計的ソリューションによって、こうした種類の類義語集、同じ意味を持つ他のキーワードを検索使用する場合、birth (誕生)というキーワードも加わり、関連しない多くの結果を生成する場合があります。言語の理解により、テキストの曖昧さが無くなり、本質的に、言語学に基づくテキストマイニングをより信頼できるアプローチにします。

抽出プロセスがどのように機能するのかを理解しておく、言語リソース (ライブラリー、タイプ、類義語など) を微調整する際に主要な決定を下すのに役立ちます。抽出プロセスのステップには以下のものがあります。

- ソース・データの標準フォーマットへの変換
- 候補となる用語の特定
- 類義語の等価クラスおよび統合の特定
- タイプの割り当て
- 二次分析によるインデックスの付与、および必要に応じてパターン・マッチ

手順 1: ソース・データの標準フォーマットへの変換

最初のステップでは、後続の分析に利用できるように、インポートしたデータを決まった形式に変換します。この変換は内部的に実行され、元のデータは変更されません。

手順 2: 候補となる用語の特定

言語学的抽出において、候補となるキーワードを特定する際の言語リソースの役割を理解しておくのは大切なことです。言語リソースは、抽出が実行されるごとに使用されます。言語リソースは、テンプレート、ライブラリー、およびコンパイル済みリソースの形式で保存されています。ライブラリーには、語のリスト、関係性、また抽出の特定や調整に使用されるその他の情報が含まれています。基幹辞書は表示・編集ができません。ただし、残りのリソースをテンプレート・エディターで、またはインタラクティブワークベンチセッションの場合はリソース・エディターで編集できます。

コンパイル済み辞書は、IBM SPSS Modeler Text Analytics の抽出エンジンの主要な、内部コンポーネントです。これらのリソースには、品詞コード (名詞、動詞、形容詞など) を含む基本形のリストを取めた一般辞書が含まれています。

これらコンパイル済み辞書のほか、製品にはいくつかのライブラリーが付属し、それらを使用して、コンパイル済み辞書のタイプ定義およびコンセプト定義を補い、また類義語を提供することができます。これらのライブラリー、および作成したユーザー指定のライブラリーは、いくつかの辞書で構成されています。これらには、キーワード辞書、類義語辞書、および不要語辞書が含まれています。

データがインポートおよび変換されると、抽出エンジンは抽出の候補のキーワードの特定を開始します。候補となるキーワードとは、テキスト内の概念を特定するのに使用される語や、語の集まりのことです。テキストを処理しているとき、単語 (ユニターム) および複合語 (マルチターム) は、品詞パターン抽出を使用して特定されます。そして、候補の感性キーワードは、感性テキスト・リンク分析を使用して特定されます。

注: 前述のコンパイル済み一般辞書にあるキーワードは、ユニタームとして重要でないまたは言語学的にあいまいであるすべての単語を示します。これらの単語は、ユニタームを特定するときに抽出から除外されます。ただし、それらは、品詞を決定またはより長い候補の複合語 (マルチターム) を参照している場合に再評価されます。

手順 3: 類義語の等価クラスおよび統合の特定

候補のユニタームおよびマルチタームが特定された後、ソフトウェアは正規化辞書を使用して、等価クラスを特定します。等価クラスとは、ある語句の基本形式、つまり同じ語句の 2 つの異なる表現を 1 つの形式で表わしたものです。等価クラスについてどのコンセプトを使用するかを判断するために、抽出エンジンは、次の規則を上から順に適用します。

- ライブラリーのユーザー指定の形式。
- コンパイル済みリソースで定義されている最も頻度の高い形式。

手順 4: タイプの割り当て

次に、抽出されたコンセプトにタイプを割り当てます。タイプは、コンセプトの意味上のグループ化です。基幹辞書ならびにライブラリーの両方がこのステップで使用されます。タイプには、上位レベルのコンセプト、肯定的な単語および否定的な単語、人名、地名、組織名などが含まれます。詳しくは、189 ページの『キーワード辞書』のトピックを参照してください。

言語学的なシステムは、知識に依存します。つまり、辞書に含まれている情報が多いほど、より高い品質の結果が得られます。類義語の定義など、辞書の内容の変更は、そのまま結果の改善につながります。これは、通常、対話的な処理で、正確なコンセプトの検索に不可欠です。NLP は IBM SPSS Modeler Text Analytics の主要な要素です。

抽出の方法

回答の主要キーワードの抽出時、IBM SPSS Modeler Text Analytics は言語学に基づくテキスト分析に依存します。このアプローチを用いると統計に基づくシステムがもたらすようなスピードと費用対効果が得られます。また人の手を介することがほとんどないので、極めて高い精度が得られます。言語学に基づくテキスト分析は、自然言語処理、あるいは計量言語学と呼ばれる研究分野に基づいています。

抽出プロセスがどのように機能するのかを理解しておく、言語リソース (ライブラリー、タイプ、類義語など) を微調整する際に主要な決定を下すのに役立ちます。抽出プロセスのステップには以下のものがあります。

- ソース・データの標準フォーマットへの変換
- 候補となる用語の特定
- 類義語の等価クラスおよび統合の特定
- タイプの割り当て
- インデックスの付与
- パターンおよびイベント抽出のマッチング

手順 1: ソース・データの標準フォーマットへの変換

最初のステップでは、後続の分析に利用できるように、インポートしたデータを決まった形式に変換します。この変換は内部的に実行され、元のデータは変更されません。

手順 2: 候補となる用語の特定

言語学的抽出において、候補となるキーワードを特定する際の言語リソースの役割を理解しておくのは大切なことです。言語リソースは、抽出が実行されるごとに使用されます。言語リソースは、テンプレート、ライブラリー、およびコンパイル済みリソースの形式で保存されています。ライブラリーには、語のリスト、関係性、また抽出の特定や調整に使用されるその他の情報が含まれています。基幹辞書は表示・編集ができ

ません。ただし、残りのリソース (テンプレート) を テンプレート・エディター で、またはインタラクティブ・ワークベンチ・セッションの場合は リソース・エディター で編集できます。

コンパイル済み辞書は、IBM SPSS Modeler Text Analytics の抽出エンジンの主要な、内部コンポーネントです。これらのリソースには、品詞コード (名詞、動詞、形容詞、副詞、分詞、限定詞、接続詞、前置詞) を含む基本形のリストを取めた一般辞書が含まれています。また、リソースには、<地名>、<組織名>、または<人名> のタイプに多くの抽出されたキーワードを割り当てるために使用する、予約済みのビルトインのタイプも含まれています。詳しくは、190 ページの『ビルトインのタイプ』のトピックを参照してください。

これらコンパイル済み辞書のほか、製品にはいくつかのライブラリーが付属し、それらを使用して、コンパイル済み辞書のタイプ定義およびコンセプト定義を補い、またその他のタイプや類義語を提供することができます。これらのライブラリー、および作成したユーザー指定のライブラリーは、いくつかの辞書で構成されています。これらには、キーワード辞書、類義語辞書 (類義語およびオプションの要素)、および不要語辞書が含まれています。詳しくは、179 ページの『第 15 章 ライブラリーの使用』のトピックを参照してください。

データがインポートおよび変換されると、抽出エンジンは抽出の候補のキーワードの特定を開始します。候補となるキーワードとは、テキスト内の概念を特定するのに使用される語や、語の集まりのことです。テキストの処理中、コンパイル済み辞書にない単語 (ユニターム) は、抽出の候補のキーワードとして見なされます。候補の複合語 (マルチターム) は、品詞パターン抽出を使用して特定されます。例えば、品詞パターンが「形容詞、名詞」のマルチターム `sports car` (スポーツ カー) は、2 つの部分に分けられます。品詞パターンが「形容詞、形容詞、名詞」のマルチターム `fast sports car` (高速スポーツ カー) は、3 つの部分に分けられます。

注: 前述のコンパイル済み一般辞書にあるキーワードは、ユニタームとして重要でないまたは言語学的にあまり重要でないすべての単語を示します。これらの単語は、ユニタームを特定するときに抽出から除外されます。ただし、それらは、品詞を決定またはより長い候補の複合語 (マルチターム) を参照している場合に再評価されます。

最後に、特殊なアルゴリズムを使用して、役職などの大文字の文字列を処理し、これらの特殊なパターンを抽出できるようにします。

手順 3: 類義語の等価クラスおよび統合の特定

候補のユニタームおよびマルチタームが特定された後、一連のアルゴリズムを使用して、ユニタームやマルチタームを比較し、等価クラスを特定します。等価クラスは、ある語句の基本形、すなわち同じ語句の 2 つの表現を 1 つの形で表わしたものです。句を等価クラスに割り当てる目的は、例えば、`president of the company` (会社の社長) および `company president` (会社社長) を別のコンセプトとして扱わないようにすることです。等価クラスのどのコンセプトを使用するか、つまり、`president of the company` (会社の社長) または `company president` (会社社長) のどちらを主要キーワードとして使用するかを判断するために、抽出エンジンは、次の規則を順に適用します。

- ライブラリーのユーザー指定の形式。
- テキスト全体で最も出現頻度の高い形式。
- テキスト全体で最も短い形式 (通常、基本型に該当)。

手順 4: タイプの割り当て

次に、抽出されたコンセプトにタイプを割り当てます。タイプは、コンセプトの意味上のグループ化です。基幹辞書ならびにライブラリーの両方がこのステップで使用されます。タイプには、上位レベルのコンセプト

ト、肯定的な単語および否定的な単語、人名、地名、組織名などが含まれます。ユーザーがタイプを定義して追加することもできます。詳しくは、189 ページの『キーワード辞書』のトピックを参照してください。

手順 5: インデックスの付与

レコードまたはドキュメントのセット全体に、テキストの位置と各等価クラスの代表キーワードの間にポインタを確定してインデックスを付けます。候補のコンセプトの活用形インスタンスはすべて、候補の基本型としてインデックスが付けられます。基本形ごとに全体の出現頻度が計算されます。

手順 6: パターンおよびイベント抽出のマッチング

IBM SPSS Modeler Text Analytics は、タイプやコンセプトだけでなく、それらの関係性も見つけることができます。この製品ではいくつかのアルゴリズムおよびライブラリーを使用でき、またタイプおよびコンセプトの間の関係性パターンを抽出する機能が用意されています。製品に対する反応などの特定の意見、または政治的グループやゲノムのリンクなど、人々またはオブジェクトの間の関係性リンクを探す場合に特に役立ちます。

カテゴリー化の方法

IBM SPSS Modeler Text Analytics でカテゴリーモデルを作成する場合、いくつかの手法から選択して、カテゴリーを作成できます。すべてのデータセットが一意であるため、手法の数やそれらを適用する順序は、変わる場合があります。結果の解釈が、他の人とは異なる場合があるため、テキスト・データにとってどの手法が最良の結果を生み出すか、それぞれの手法を検証する必要があります。IBM SPSS Modeler Text Analytics では、カテゴリーをさらに検証し、調整できるワークベンチ・セッションでカテゴリー・モデルを作成できます。

このガイドの場合、カテゴリーの作成は、カテゴリー定義の生成および、1 つまたは複数のビルトインの手法を使用した分類を指し、またカテゴリー化は、スコアリング、またはラベル付け、一意の識別子 (名前/ID/値) を各レコードまたはドキュメントのカテゴリー定義に割り当てるプロセスのことを指します。

カテゴリー作成時、抽出されたコンセプトおよびタイプはカテゴリーの構築ブロックとして使用されます。カテゴリーを作成すると、カテゴリー定義の要素に一致するテキストが含まれる場合、レコードおよびドキュメントが自動的にカテゴリーに割り当てられます。

IBM SPSS Modeler Text Analytics には、自動カテゴリー作成手法がいくつか用意されており、ドキュメントまたはレコードを迅速にカテゴリー化することができます。

グループ化手法

使用できるそれぞれの手法は、特定の種類のデータおよび状況に適していますが、同じ分析で手法を組み合わせて、全範囲のドキュメントまたはレコードをキャプチャーすると役に立つ場合があります。複数のカテゴリーのコンセプトを表示したり、重複するカテゴリーを見つけることができます。

派生関係のコンセプトの語幹: コンセプト・コンポーネントが形態的に関連するか、または語幹を共有するかどうかを分析するとき、コンセプトを取得し、そのコンセプトに関連するその他のコンセプトを検索することによって、カテゴリーを作成します。この手法は、生成された各カテゴリーのコンセプトが類義語または意味の上で密接に関連しているため、類義語の複合語コンセプトを特定するのに非常に役立ちます。長さの異なるデータを処理し、コンパクトなカテゴリーをより少なく生成します。例えば、コンセプト `opportunities to advance` は、コンセプト `opportunity for advancement` および `advancement opportunity` とグループ化されます。詳しくは、111 ページの『派生関係のコンセプトの語幹』のトピックを参照してください。

セマンティック・ネットワーク: 各コンセプトの考えられる意味を、単語の関係の拡張インデックスから特定することによって開始し、関連するコンセプトをグループ化することによってカテゴリを作成します。この手法は、コンセプトがセマンティック・ネットワークに認識され、あまり曖昧でない場合に最も適しています。テキストに、ネットワークが認識していない特殊な用語または専門用語が含まれている場合はあまり役に立ちません。例えば、コンセプト `granny smith apple` は、`granny smith` と横の関係があるため、`gala apple` および `winesap apple` とグループ化されます。また別の例では、コンセプト `animal` は、その下位語である `cat` および `kangaroo` とグループ化されます。このリリースでは、英語テキストにのみ使用できます。詳しくは、113 ページの『セマンティック・ネットワーク』のトピックを参照してください。

内包関係のコンセプト: この手法では、一方の共通の文字列である単語を含むかどうかに基づき、マルチタームのコンセプト (複合語) をグループ化することによってカテゴリを作成します。例えば、コンセプト `seat` (シート) は、コンセプト `safety seat` (セーフティ シート)、`seat belt` (シート・ベルト)、および `seat belt buckle` (シート・ベルトのバックル) とグループ化されます。詳しくは、112 ページの『内包関係のコンセプト』のトピックを参照してください。

共起: この手法では、テキスト内の共起関係のコンセプトからカテゴリを作成します。ドキュメントおよびレコードでコンセプトまたはコンセプト・パターンがいっしょに出現することが多いとき、共起関係のコンセプトはおそらくカテゴリ定義の値のものである基底となる関連を反映します。単語が頻繁に共起する場合、共起規則が作成され、新しいサブカテゴリのカテゴリの記述子として使用できます。例えば、多くのレコードに単語 `price` (価格) および `availability` (有効性) が含まれている場合 (ただし、一方を含み、もう一方を含まないレコードはほとんどない)、これらのコンセプトを共起規則 (`price & available`) にグループ化し、例えばカテゴリ `price` のサブカテゴリに割り当てることができます。詳しくは、114 ページの『共起規則』のトピックを参照してください。

最小 ドキュメント。ドキュメント数: 共起関係のコンセプトの重要性を判断できるようにするため、カテゴリの記述子として使用されるよう、指定の共起関係のコンセプトを含む必要のあるドキュメントまたはレコードの最小数を定義します。

IBM SPSS Modeler Text Analytics ノード

IBM SPSS Modeler に付属する多くの標準ノードとともに、テキスト マイニング・ノードを使用して、テキスト分析の精度をストリームに組み込むこともできます。IBM SPSS Modeler Text Analytics では、それを実行するためのテキスト マイニング・ノードがいくつか用意されています。これらのノードは、ノード・パレットの IBM SPSS Modeler Text Analytics タブに保存されています。

次のノードが含まれます。

- **ファイル・リスト入力ノード**で、テキスト マイニング・プロセスへの入力として、ドキュメント名のリストを生成します。これは、テキストがデータベースや他の構造化ファイルではなく、外部ドキュメントに存在する場合に役立ちます。ノードは、リストされた各ドキュメントあるいはフォルダの 1 レコードを持つたった 1 つのフィールドを出力します。これは、後続のテキスト マイニング・ノードの入力として選択できます。詳しくは、11 ページの『ファイル リスト ノード』のトピックを参照してください。
- **Web フィード入力ノード**を使用して、RSS または HTML 形式のブログまたはニュース・フィードなど、Web フィードからテキストを読み取り、このデータをテキスト マイニング・プロセスで使用できます。Web フィード入力ノードは、フィードの各レコードに 1 つまたは複数のフィールドを出力しますが、後続のテキスト マイニング・ノードでは、これを入力として選択できます。詳しくは、13 ページの『Web フィード ノード』のトピックを参照してください。

- 言語の識別子ノードはプロセス ノードの一種であり、ソース テキストをスキャンして、書かれている人間の言語を判別し、言語のマークを新しいフィールドに書き込みます。主に大量のデータとともに使用するよう設計されています。このノードは、データ ソースに複数の言語が存在するが、1 つの言語のみを処理したい場合に特に役立ちます。詳しくは、17 ページの『言語ノード』のトピックを参照してください。
- テキスト マイニング・ノードでは、言語学的手法を使用して、主要なコンセプトをテキストから抽出します。これらのコンセプトおよびそのほかのデータを使用してカテゴリーを作成することができ、既知のパターンに基づいてコンセプト間の関係および関連を特製する機能 (テキスト リンク分析) を用意しています。ノードを使用して、テキスト・データの内容を検討、またはコンセプト・モデルまたはカテゴリー・モデルのいずれかを作成できます。コンセプトおよびカテゴリーは、人口統計などの既存の構造化されたデータを組み合わせることができ、モデル作成に適用することができます。詳しくは、20 ページの『テキスト マイニング・モデル作成ノード』のトピックを参照してください。
- テキスト リンク分析ノードは、コンセプトを抽出し、またテキスト内の既知のパターンに基づいて、コンセプト間の関係を特定します。パターンを抽出して、これらのコンセプトに関連付けられた意見または識別子のほか、コンセプト間の関係を見出すことができます。テキスト リンク分析ノードを使用して、より直接的にテキストからパターンを特定および抽出し、パターンの結果をストリーム内のデータセットに追加できます。ただし、テキスト マイニングモデル作成ノードのインタラクティブ・ワークベンチ・セッションを使用して TLA を実行することもできます。詳しくは、47 ページの『テキスト リンク分析ノード』のトピックを参照してください。
- 外部ドキュメントからテキストをマイニングする場合、テキスト マイニング出力ノードが、コンセプトが抽出されたドキュメントへのリンクを含む HTML ページを生成するのに使用できます。詳しくは、55 ページの『ファイル ビューアー ノード』のトピックを参照してください。

アプリケーション

一般的に、定期的に多くの分量のドキュメントを確認して、より詳細に検討するために主要な要素を特定する必要がある人々は、IBM SPSS Modeler Text Analytics を使用すると多くの利点があります。

特定のアプリケーションには、次のような機能があります。

- 科学的小および医学的研究: 特許報告書、雑誌の記事、計画書の発行物など、二次リサーチの使用を検証します。以前は知られていなかった (例えば特定の製品に関連した医者など) 関連性を識別します。薬品の開発プロセスにかかる時間を最小化します。遺伝子調査の補助として使用します。
- 投資リサーチ: 毎日のアナリスト・レポート、ニュース記事、企業のプレス・リリースを確認して、主要な戦略ポイントまたは市場シフトを特定します。こうした情報のトレンド分析により、一定の期間にわたって、企業または業界の緊急の問題または機会について明らかにします。
- 不正検出: 銀行および医療費の不正を使用して、以上を検出し、多数のテキストから警告を検出します。
- 市場リサーチ: 市場リサーチの段階で使用し、自由記述式アンケートの回答の主要なトピックを特定します。
- ブログおよび Web フィード分析: 新しいフィード、ブログなどの主要なキーワードを使用して、モデルを検証および作成します。
- CRM. 電子メール、取引、調査など、すべての顧客との接点からのデータを使用し、モデルを作成します。

第 2 章 ソース・テキストの読み取り

テキストマイニングのデータは、データベースや、データを行と列で表現する他の「長方形」の形式などの IBM SPSS Modeler で使用される標準的な形式、またはこの構造に準拠しない Microsoft Word、Adobe PDF、HTML などのドキュメント形式です。

- Microsoft Word、Microsoft Excel、Microsoft PowerPoint のほか、Adobe PDF、XML、HTML など、標準のデータ構造に従っていないドキュメントのテキストを読み取るために、ファイル リスト ノードを使用して、ドキュメントまたはフォルダのリストをテキストマイニングへの入力として生成できます。詳しくは、『ファイル リスト ノード』を参照してください。
- RSS または HTML 形式のブログまたはニュース・フィードなど、Web フィードからテキストを読み取るために、Web フィード ノードを使用して Web フィード・データをテキストマイニング・プロセスの入力用に書式設定できます。詳しくは、13 ページの『Web フィード ノード』を参照してください。
- 顧客のコメント用の 1 つ以上のテキスト フィールドを含むデータベースなど、SPSS Modeler で使用する標準データ形式のテキストを読み取るために、任意の SPSS Modeler 入力ノードを使用できます。詳しくは、SPSS Modeler ノードの資料を参照してください。
- 大量のデータを処理するときに、テキストに複数の言語が存在する可能性がある場合は、言語ノードを使用して、特定のフィールドで使用されている言語を識別してください。詳しくは、17 ページの『言語ノード』を参照してください。

ファイル リスト ノード

Microsoft Word、Microsoft Excel、Microsoft PowerPoint、さらに Adobe PDF、XML、HTMLなどの形式で保存された、構造のないドキュメントのテキストを読み取るために、ファイル リスト ノードを使用して、ドキュメントまたはフォルダのリストをテキストマイニング プロセスへの入力として生成できます。構造のないテキスト・ドキュメントは、IBM SPSS Modeler で使用される他のデータのようにフィールドやレコード (行および列) で表すことができないため、この方法が必要になります。

ファイル リスト ノードは、入力ノードとして機能します。

このノードは、IBM SPSS Modeler ウィンドウの下部にあるノード・パレットの IBM SPSS Modeler Text Analytics タブにあります。詳しくは、8 ページの『IBM SPSS Modeler Text Analytics ノード』のトピックを参照してください。

重要: マシンのローカル エンコードに含まれる文字を含むディレクトリ名およびファイル名はサポートされません。ファイル リスト ノードを含むストリームを実行しようとする、これらの文字を含むファイル名またはディレクトリ名により、ストリームの実行が失敗します。この問題は、外国語のディレクトリ名またはファイル名を使用する場合 (フランス語のロケールでドイツ語のファイル名を使用する場合など) に発生する可能性があります。

ローカル データのサポート。リモートの IBM SPSS Modeler Text Analytics Server に接続し、ファイル リスト ノードを含むストリームがある場合、データが IBM SPSS Modeler Text Analytics Server と同じマシン上にあるか、ファイル リスト ノードのソース データが格納されているフォルダへのアクセス権限がサーバー マシンに割り当てられている必要があります。

注: ファイル リスト ノードを IBM SPSS Collaboration and Deployment Services - Scoring 構成内のスコアリングに使用することはできません。

ファイル リスト ノード: 「設定」 タブ

このタブで、このノードのディレクトリー、ファイルの拡張子、入力を定義します。

注: テキスト マイニング抽出は、Microsoft Windows 以外のプラットフォームにある Microsoft Office と Adobe PDF ファイルを処理できません。ただし、XML、HTML またはテキスト・ファイルは常に処理可能です。

マシンのローカル エンコードに含まれる文字を含むディレクトリ名およびファイル名はサポートされません。ファイル リスト ノードを含むストリームを実行しようとする、これらの文字を含むファイル名またはディレクトリ名により、ストリームの実行が失敗します。この問題は、外国語のディレクトリ名またはファイル名を使用する場合 (フランス語のロケールでドイツ語のファイル名を使用する場合など) に発生する可能性があります。

ディレクトリー 一覧表示するドキュメントを含むルート フォルダーを指定します。

- サブディレクトリーを含める サブディレクトリーもスキャンする場合に指定します。

リストに含めるファイル形式: 使用するファイル形式および拡張子を選択または選択解除できます。ファイルの拡張子の選択を解除すると、その拡張子を持つファイルは無視されます。次の拡張子でフィルタリングできます。

表 1. ファイル拡張子でフィルターされるファイル・タイプ :

• .rtf、.doc、.docx、.docm	• .xls、.xlsx、.xlsm	• .ppt、.pptx、.pptm	• .txt、.text
• .htm、.html、.shtml	• .xml(X)	• .pdf	• .

注: 詳しくは、11 ページの『ファイル リスト ノード』を参照してください。

拡張子のないファイルまたは末尾にドットの拡張子が付いているファイル (File01 や File01. など) が存在する場合、「拡張子なし」オプションを使用してこれらのファイルを選択してください。

入力エンコード: 出力フィールドに完全一致のテキストが含まれる場合は、次のリストから該当する値を選択してください。

- 自動 (ヨーロッパ)
- UTF-8
- UTF-16
- ISO-8859-1
- ISO-8859-2
- Windows-1250
- US ascii

出力は UTF-8 ドキュメント テキストとして表示されます。

重要: バージョン 14 以降、「ディレクトリーのリスト」オプションは使用できなくなりました。ファイルのリストのみ出力されます。

ファイル リスト ノード:その他のタブ

IBM SPSS Modeler ノードの「データ型」タブは、「注釈」タブ同様、標準タブです。

テキスト マイニングでのファイル リスト ノードの使用

テキスト・データが、Microsoft Word、Microsoft Excel、Microsoft PowerPoint、Adobe PDF、XML、HTML などの形式の、構造のない外部ドキュメント内にある場合、ファイル リスト ノードを使用します。

例として、ファイル リスト ノードをテキスト マイニング・ノードに接続して、外部ドキュメント内にあるテキストを指定します。

1. ファイル リスト ノード（「設定」タブ）。 まず、このノードをストリームに追加して、テキスト・ドキュメントが保存されている場所を指定しました。テキスト マイニングを実行するすべてのドキュメントを含むディレクトリーを選択しました。
2. テキスト マイニング・ノード（「フィールド」タブ）: 次に、テキスト マイニング・ノードを ファイル リスト ノードに追加して接続しました。このノードで、入力形式、リソース・テンプレート、および出力形式を定義しました。ファイル リスト ノードから作成されたフィールド名、テキスト フィールド、およびその他の設定を選択しました。詳しくは、 30 ページの『ストリーム内のテキスト マイニング・ノードの使用』のトピックを参照してください。

テキスト マイニング・ノード使用の詳細は、「 20 ページの『テキスト マイニング・モデル作成ノード』」を参照してください。

Web フィード ノード

Web フィード ノードを使用して、Web フィードのテキスト・データをテキスト マイニング・プロセス向けに準備することができます。このノードは、次の 2 つの形式で Web フィードを受け入れます。

- RSS 形式: RSS は、Web コンテンツ向けの単純な XML ベースの標準化形式です。この形式の URL は、組織化されたニュース ソースやブログなどのリンクした記事のセットがあるページを示します。RSS は標準化された形式であるため、リンクした記事は自動的に特定され、データ・ストリームの個別のレコードとして扱われます。フィルタリング手法をテキストに適用しない限り、フィードの重要なテキスト・データおよびレコードを特定するために、さらなる入力はありません。
- HTML 形式: HTML ページに対する 1 つ以上の URL を「入力」タブで定義できます。「レコード」タブで、レコードの開始タグを定義し、対象の内容を区切るタグを指定して、これらのタグを選択した出力フィールド（説明、タイトル、更新日など）に割り当てます。詳しくは、 15 ページの『Web フィード ノード:「レコード」タブ』のトピックを参照してください。

重要: プロキシ・サーバーを経由して Web の情報を取得しようとしている場合、IBM SPSS Modeler Text Analytics Client および Server の `net.properties` ファイルでプロキシ・サーバーを有効にする必要があります。このファイル内の説明に従ってください。これは、Web フィード ノードを使用して Web にアクセスする場合または Language Weaver ライセンスを取得する場合に適用されます。これらの場合、接続が Java™ を経由するためです。デフォルトでは、ファイルは `C:\Program Files\IBM\SPSS\Modeler\18.1\jre\lib\net.properties` にあります。

このノードの出力は、レコードの説明に使用するフィールドのセットです。「説明」フィールドは、多くのテキスト・コンテンツが含まれているため、最も一般的に使用されます。ただし、レコードの短い説明（「短い説明」フィールド）またはレコードのタイトル（「タイトル」フィールド）など、他のフィールドにも関心がある場合があります。出力フィールドのいずれかを、後続のテキスト マイニング・ノードの入力として選択できます。

注: 「Web フィード」ノードを IBM SPSS Collaboration and Deployment Services - Scoring 構成内のスコアリングに使用することはできません。

このノードは、IBM SPSS Modeler ウィンドウの下部にあるノード・パレットの IBM SPSS Modeler Text Analytics タブにあります。詳しくは、8 ページの『IBM SPSS Modeler Text Analytics ノード』のトピックを参照してください。

Web フィード ノード: 「入力」タブ

「入力」タブを使用して、1 つまたは複数の Web アドレスまたは URL を指定し、テキスト・データをキャプチャーします。テキスト マイニングのコンテキストで、テキスト・データを含むフィードの URL を指定できます。

重要: 非 RSS データを扱う場合、WebQL[®] などの Web スクラッピング・ツールを使用して、コンテンツが異なる入力ノードを使用するツールから出力を収集して参照するよう自動化することをお勧めします。

設定できるパラメーターを次に示します。

URL を入力または貼り付け: 1 つまたは複数の URL を入力または貼り付けることができます。複数の URL を入力する場合、1 行ごとに 1 つの URL だけが入力し、**Enter/Return** キーを使用して、行を区切ります。ファイルへの完全な URL パスを入力します。フィードを示すこれらの URL は次の 2 つの形式のいずれかとなります。

- **RSS 形式:** RSS は、Web コンテンツ向けの単純な XML ベースの標準化形式です。この形式の URL は、組織化されたニュース ソースやブログなどのリンクした記事のセットがあるページを示します。RSS は標準化された形式であるため、リンクした記事は自動的に特定され、データ・ストリームの個別のレコードとして扱われます。フィルタリング手法をテキストに適用しない限り、フィードの重要なテキスト・データおよびレコードを特定するために、さらなる入力はありません。
- **HTML 形式:** HTML ページに対する 1 つ以上の URL を「入力」タブで定義できます。「レコード」タブで、レコードの開始タグを定義し、対象の内容を区切るタグを指定して、これらのタグを選択した出力フィールド (説明、タイトル、更新日など) に割り当てます。非 RSS データを扱う場合、WebQL[®] などの Web スクラッピング・ツールを使用して、コンテンツが異なる入力ノードを使用するツールから出力を収集して参照するよう自動化することをお勧めします。詳しくは、15 ページの『Web フィード ノード: 「レコード」タブ』のトピックを参照してください。

URL ごとに読み込む最新エントリー数: フィード内にある最初のレコードから始まるフィールドに表示された各 URL に読み込む最大レコード数を指定します。テキストの量は、テキストマイニング・ノードまたはテキスト リンク分析ノード下流の抽出の処理速度に影響を与えます。

可能な場合、以前の Web フィードを保存および再利用: このオプションで、Web フィードをスキャンし、処理された結果をキャッシュします。そして、後続のストリームの実行後、指定されたフィードの内容が変わらない場合、またはフィードにアクセスできない場合 (インターネットの機能停止など)、キャッシュされたバージョンを使用して、処理時間を短縮します。これらのフィードで見つかった新しいコンテンツは、次回ノードを実行するときにキャッシュされます。

- **ラベル.** 「可能な場合、以前の Web フィードを保存および再利用」を選択した場合、その結果のラベル名を指定する必要があります。このラベルを使用して、サーバーのキャッシュされたフィードを説明します。ラベルが指定されていない場合、またはラベルが認識されない場合、再利用はできません。これらの Web フィードのキャッシュを、IBM SPSS Deployment Manager に組み込まれた IBM SPSS Text Analytics Administration Console のセッション テーブルで管理できます。詳細は、『Deployment Manager ユーザー・ガイド』を参照してください。

Web フィード ノード: 「レコード」 タブ

「レコード」タブを使用して、新しいレコードの開始点、各レコードに関するその他の関連情報を特定し、非 RSS フィードのテキスト・コンテンツを指定します。非 RSS フィード (HTML) に複数のレコード内にあるテキストが含まれていることがわかっている場合、レコードの開始タグをここで指定する必要があります。指定しない場合、テキストは 1 つのレコードとして扱われます。RSS フィードは標準化され、このタブでタグの指定は必要ありませんが、「プレビュー」タブで内容をプレビューできます。

重要: 非 RSS データを扱う場合、WebQL[®] などの Web スクラッピング・ツールを使用して、コンテンツが異なる入力ノードを使用するツールから出力を収集して参照するよう自動化することをお勧めします。

URL: このドロップダウン・リストには、「入力」タブで入力された URL のリストが表示されます。

HTML 形式および RSS 形式のフィードが表示されます。URL アドレスが長すぎてドロップダウン・リストにすべてを表示できない場合、自動的に省略記号を使用して途中で省略したテキストが表示されます (例: <http://www.spss.com/example/start-of-address...rest-of-address/path.htm>)。

- **HTML** 形式のフィードで、フィードに複数のレコード (エントリー) がある場合、テーブルに表示されたフィールドに対応するデータを含む HTML タグを定義できます。例えば、新しいレコードが開始したことを示す開始タグ、更新日のタグ、または作成者の名前を定義できます。
- **RSS** 形式のフィードの場合、RSS は標準化された形式であるため、タグの入力は要求されません。ただし、必要に応じて「プレビュー」タブでサンプルの結果を表示できます。認識されたすべての RSS フィードの前に、RSS のロゴ・イメージが表示されます。

「ソース」タブ: HTML フィードのソース・コードを表示できます。このコードは編集できません。「検索」フィールドを試用して、このページの特定のタグまたは情報を検索できます。それらは、下のテーブルにコピーして貼り付けることができます。「検索」フィールドでは大文字および小文字の区別はされず、また文字列の一部に一致します。

「プレビュー」タブ: Web フィード ノードでレコードがどのように読み取られるかをプレビューできます。「プレビュー」タブ下のテーブルで HTML タグを定義して、レコードがどのように読み取られるかを変更できるため、このオプションは HTML フィードを使用する場合に特に役立ちます。

非 **RSS** レコードの開始タグ. このオプションは、非 RSS フィードにのみ適用されます。HTML ここで指定します。非 RSS フィードの開始タグを定義しない場合、ページ全体が 1 つのレコードとして扱われ、コンテンツ全体が「説明」フィールドで出力、そしてノードの実行日が「更新日」および「公開日」の両方に使用されます。

フィールド・テーブル: このオプションは、非 RSS フィードにのみ適用されます。このテーブルで、事前定義された出力フィールドの開始タグを入力して、テキスト・コンテンツを特定の出力フィールドに分割することができます。開始タグのみを入力します。HTML を解析し、テーブルの内容を HTML 内のタグ名および属性に一致させることによって、すべての合致が行われます。下部のボタンを使用して、定義したタグをコピーし、他のフィールドにそのタグを再利用します。

表 2. 非 RSS フィード (HTML 形式) に使用できる出力フィールド

出力フィールド名	期待されるタグの内容
タイトル	レコードのタイトルを区切るタグ。(オプション)
短い説明	短い説明またはラベルを区切るタグ。(オプション)
説明	メインのテキストを区切るタグ。空白のままにすると、このフィールドには <body> タグ (レコードが 1 つある場合) 他のすべての内容または現在のレコード内の内容 (レコードの区切り文字が指定されている場合) が入力されます。

表 2. 非 RSS フィード (HTML 形式) に使用できる出力フィールド (続き)

出力フィールド名	期待されるタグの内容
作成者	レコードの作成者を区切るタグ。(オプション)
コントリビュータ	コントリビュータの名前を区切るタグ。(オプション)
公開日	テキストが公開された日付を区切るタグ。空白のままにすると、このフィールドにはノードがデータを読み取った日付が指定されます。
変更日	テキストが更新された日付を区切るタグ。空白のままにすると、このフィールドにはノードがデータを読み取った日付が指定されます。

タグをテーブルに入力すると、このタグを完全一致ではなく一致すべき最小タグとして使用してフィードをスキャンします。つまり、「タイトル」フィールドに <div> と入力した場合、指定した属性を持つタグ (<div class="post three">) など、フィールドの <div> タグに一致し、<div> がルート・タグ (<div>) に等しくなり、属性を含むデリバティブが「タイトル」出力フィールドにその内容を使用します。ルート・タグを入力すると、さらに詳細な属性も含まれます。

表 3. 出力フィールドのテキストの特定に使用される HTML タグの例

たとえば、次のように入力するとします。	一致タグ	その他の一致タグ	一致しないタグ
<div>	<div>	<div class="post">	その他のタグ
<p class="auth">	<p class="auth">	<p color="black" class="auth" id="85643">	<p color="black">

Web フィード ノード: 「コンテンツ フィルター」 タブ

「コンテンツ フィルター」タブを使用して、フィルター手法を RSS フィード・コンテンツに適用します。このタブは、HTML フィードには適用されません。フィードに多くのテキストがヘッダー、フッター、メニュー、広告などの形式で含まれている場合、フィルタリングが必要な場合があります。このタブを使用して、不要な HTML タグ、JavaScript、短い単語または行をコンテンツから除外することができます。

コンテンツのフィルタリング: クリーニング手法を適用したくない場合、「なし」を選択します。適用する場合は、「**RSS** コンテンツ クリーナー」を選択します。

「**RSS** コンテンツ クリーナー オプション」: 「**RSS** コンテンツ クリーナー」を選択すると、特定の基準に基づいて、行を破棄することができます。行は、<p> および <i> などの HTML タグによって区切られます。ただし、、、および などのインライン・タグは使用されません。
 タグは改行として処理されます。

- 短い行を破棄: ここで定義する最小単語数を含まない行が無視されます。
- 短い単語を含む行を破棄: ここで定義する単語の平均長の最小値より長さが短い行が無視されます。
- **1** 文字単語が多い行を破棄: ここで定義する特定の**1** 文字単語の比率より小さい行が無視されます。
- 特定のタグを含む行を破棄: フィールドで指定されたタグのいずれかを含む行のテキストが無視されます。
- 特定のテキストを含む行を破棄: フィールドで指定されたテキストのいずれかを含む行が無視されます。

テキスト マイニングでの Web フィード ノードの使用

Web フィード ノードを使用して、インターネットの Web フィードのテキスト・データをテキスト マイニング・プロセス向けに準備することができます。このノードは、HTML 形式または RSS 形式で Web フィードを受け入れます。これらのフィードは、テキスト マイニングプロセス (後続のテキスト マイニング・ノードまたはテキスト リンク分析ノード) の入力として機能します。

Web フィード ノードを使用する場合、「テキスト」フィールドがテキスト マイニング・ノードまたはテキスト リンク分析ノードの実際のテキストを示すよう指定して、これらのフィードが各記事またはブログのエントリーに直接リンクするようにする必要があります。

重要: プロキシ・サーバーを経由して Web の情報を取得しようとしている場合、IBM SPSS Modeler Text Analytics Client および Server の `net.properties` ファイルでプロキシ・サーバーを有効にする必要があります。このファイル内の説明に従ってください。これは、Web フィード ノードを使用して Web にアクセスする場合または Language Weaver ライセンスを取得する場合に適用されます。これらの場合、接続が Java を経由するためです。デフォルトでは、ファイルは `C:\Program Files\IBM\SPSS\Modeler\18.1\jre\lib\net.properties` にあります。

例:テキストマイニング モデル作成ノードを使用した Web フィード ノード (RSS フィード)

例として、Web フィード ノードをテキスト マイニング・ノードに接続して、RSS フィードのテキスト・データをテキスト マイニング・プロセスに提供するとします。

1. **Web フィード ノード** (「入力」タブ): まず、このノードをストリームに追加して、フィード・コンテンツの場所を指定し、コンテンツの構造を検証しました。最初のタブで、URL を RSS フィードに指定しました。この例は RSS フィード向けであるため、形式は既に定義されており、「レコード」タブで変更を行う必要はありません。オプションのコンテンツ フィルタリング・アルゴリズムを RSS フィードに使用できますが、この場合は適用されませんでした。
2. **テキスト マイニング・ノード** (「フィールド」タブ): 次に、テキスト マイニング・ノードを Web フィード ノードに追加して接続しました。このタブで、Web フィード ノードによってテキスト・フィールド出力を定義しました。この場合、「説明」フィールドを使用する必要がありました。また、「テキスト」フィールドが実際のテキストやその他の設定を示すオプションを選択しました。
3. **テキスト マイニング・ノード** (「モデル」タブ): 次に「モデル」タブで、ビルド モードとリソースを選択しました。この例では、デフォルトのリソース・テンプレートを使用して、このノードから直接コンセプト・モデルを作成するよう選択しました。

テキスト マイニング・ノード使用の詳細は、「20 ページの『テキスト マイニング・モデル作成ノード』」を参照してください。

言語ノード

言語ノードを使用して、ソース データに存在するテキスト フィールドの自然言語を識別できます。

このノードの出力は、検出された言語コードを含む派生フィールドです。

注: 言語ノードを IBM SPSS Collaboration and Deployment Services - Scoring 構成内のスコアリングに使用することはできません。

このノードは、IBM SPSS Modeler ウィンドウの下部にあるノード・パレットの IBM SPSS Modeler Text Analytics タブにあります。詳しくは、8 ページの『IBM SPSS Modeler Text Analytics ノード』のトピックを参照してください。

言語ノード: 「設定」 タブ

このタブでは、選択されたテキスト フィールドの言語詳細を出力する方法を指定します。

テキスト フィールド: 言語を識別するテキスト フィールドを選択します。

派生フィールド名 (**Derive field name**): 検出した言語コードを格納する派生フィールドの名前を入力します。デフォルト値は *Language* です。

言語を識別できない場合のデフォルト値 (**Default value for when language cannot be identified**): 言語を識別できない場合に作成するフィールドの名前を指定します。選択可能な項目は次のとおりです。

- 未定義 (**Undefined**): これを選択すると、派生フィールドにヌル値が入ります。
- サポート対象 (**Supported**): これを選択すると、サポートされる以下の ISO 言語のいずれかを選択できます。
 - 英語 (EN)
 - ドイツ語 (DE)
 - スペイン語 (ES)
 - フランス語 (FR)
 - イタリア語 (IT)
 - オランダ語 (NL)
 - ポルトガル語 (PT)
- ユーザー設定: サポートされているいずれの言語も該当しない場合は、このオプションを使用して、ユーザー指定の値を使用するように指定します。通常は 2 文字の ISO 言語コードにしますが、必要に応じて任意のテキスト文字列を使用できます。

第 3 章 コンセプトおよびカテゴリーのマイニング

テキスト マイニング モデル作成ノードを使用して、次の 2 つのテキスト マイニング モデル ナゲットのうちいずれかを生成します。

- コンセプト モデル ナゲット は構造のあるテキスト・データまたは構造のないテキスト・データの目立ったコンセプトを明らかにし、テキスト・データから抽出します。
- カテゴリー モデル ナゲット はドキュメントおよびレコードをスコアリングし、抽出したコンセプト (およびパターン) で構成されたカテゴリーに割り当てます。

モデル ナゲットから抽出されたコンセプト、パターン、カテゴリーをすべて人口統計など既存の構造化されたデータと組み合わせ、IBM SPSS Modeler のツールの完全パッケージを使用して適用し、より適切で焦点を絞った決定を行うことができます。例えば、顧客が頻繁にログイン問題をオンライン・アカウント管理タスクの完了に対する主な障害として頻繁に一覧化する場合、「ログイン問題」をモデルに組み込むことが必要な場合があります。

また、テキスト マイニング モデル作成ノードは IBM SPSS Modeler で完全に統合されており、PredictiveCallCenter のようなアプリケーションでの構造のないデータをリアルタイムにスコアリングするために、IBM SPSS Modeler Solution Publisher を使用してテキスト マイニング ストリームを展開できます。これらのストリームを展開する機能により、正常でクローズドループのテキスト マイニングの実装を実現します。例えば、組織は、予測モデルを適用してマーケティング メッセージの精度をリアルタイムに向上させることにより、受信者または発信者からのメモ帳のメモを分析することができます。ストリームでテキスト マイニング・モデルの結果を使用すると、予測データ・モデルの精度が向上します。

IBM SPSS Modeler Text Analytics を IBM SPSS Modeler Solution Publisher とともに実行するには、ディレクトリー <install_directory>/ext/bin/spss.TMWServer を \$LD_LIBRARY_PATH 環境変数に追加します。

IBM SPSS Modeler Text Analytics では、抽出されたコンセプトおよびカテゴリーを参照します。探索を目的とする作業およびモデル作成の間、より情報に基づく決定を行うことができるため、コンセプトおよびカテゴリーの意味を理解することは重要です。

コンセプトおよびコンセプト モデル ナゲット

抽出プロセスで、テキスト・データをスキャンして分析し、テキスト内の関心のあるまたは関連する単語 (選挙または平和など) や語句 (大統領選挙、大統領の選挙、または平和条約など) を特定します。これらの単語や句を、まとめて「キーワード」と呼びます。言語リソースを使用して、関連キーワードを抽出し、類似したキーワードをコンセプトと呼ばれる代表語でグループ化します。

このように、コンセプトはテキストおよび使用している言語リソースのセットによって、複数の基本キーワードを示すことができます。例えば、従業員の満足調査を行い、コンセプト「給料」が抽出されたとします。また、「給料」に関連するレコードを参照し、「給料」が常にテキスト内にあるのではなく、特定のレコードにキーワード「賃金」、「報酬」、および「給与」のような類似した単語が含まれているとします。これらのキーワードは、抽出エンジンがキーワードを類似した単語として認識し、処理規則または言語リソースに基づいて類義語であると判断するため、「給与」という名でグループ化されます。この場合、これらのキーワードのいずれかを含むドキュメントまたはレコードは、単語「給料」を含む場合と同じように扱われます。

コンセプトに基づいてグループ化されるキーワードを確認したい場合、インタラクティブ・ワークベンチでコンセプトを検索したり、コンセプト・モデルに示される類義語を確認することができます。詳しくは、34 ページの『コンセプト・モデルの基本キーワード』のトピックを参照してください。

コンセプト モデル ナゲットは、コンセプト (類義語またはグループ化されたキーワードを含む) を含むレコードまたはドキュメントの特定に使用できる一連のコンセプトで構成されています。コンセプト・モデルは次の 2 つの方法で使用できます。1 つ目の方法は、元のソース・テキストで見つかったコンセプトを検証および分析、または関心のあるドキュメントをすばやく特定することです。2 つ目の方法は、このモデルを新しいテキスト・レコードまたはドキュメントに適用し、コール・センターのメモ帳データから主要キーワードをリアルタイムに発見するなど、新しいドキュメント/レコードの同じ主要キーワードをすばやく特定することです。

詳しくは、31 ページの『テキスト マイニング モデル ナゲット:コンセプト・モデル』のトピックを参照してください。

カテゴリーおよびカテゴリー モデル ナゲット

テキスト内の主要なキーワード、情報、属性をキャプチャーする高いレベルのコンセプトまたはトピックを示すカテゴリーを作成できます。カテゴリーは、コンセプト、タイプ、および規則などの一連の記述子で構成されています。また、これらの記述子を共に使用して、レコードまたはドキュメントが指定されたカテゴリーに属するかどうかを特定します。ドキュメントまたはレコードをスキャンして、テキストが記述子に合致するかどうかを確認することができます。合致が見つかった場合は、ドキュメント/レコードはそのカテゴリーに割り当てられます。このプロセスを、カテゴリー化といいます。

製品の自動的手法の頑健なセットを使用して、またはデータに関する詳細な洞察を手動で使用して、あるいはそれらを組み合わせてカテゴリーを自動的に作成することができます。このノードの「モデル」タブを使用してテキスト分析パッケージから事前に作成された一連のカテゴリーを読み込むこともできます。カテゴリーの手動作成またはカテゴリーの調整は、インタラクティブ・ワークベンチでのみ実行できます。詳しくは、24 ページの『テキスト マイニング・ノード:「モデル」タブ』のトピックを参照してください。

カテゴリー モデル ナゲットは、一連のカテゴリーとその記述子で構成されています。モデルを使用して、各ドキュメント/レコードのテキストに基づいて、一連のドキュメントまたはレコードをカテゴライズできます。各ドキュメントまたレコードが読み取られ、記述子の合致が見つかった各カテゴリーに割り当てられます。このように、ドキュメントまたはレコードを複数のカテゴリーに割り当てることができます。カテゴリー モデル ナゲットを使用して、自由記述式アンケートの回答またはブログのエントリーなどの不可欠なキーワードを確認することができます。

詳しくは、40 ページの『テキスト マイニング モデル ナゲット:カテゴリー・モデル』のトピックを参照してください。

テキスト マイニング・モデル作成ノード

テキスト マイニング モデル作成ノードは、言語学的手法および出現頻度に基づく手法を使用して、テキストから主要キーワードを抽出し、これらのコンセプトおよびその他のデータでカテゴリーを作成します。ノードを使用して、テキスト・データの内容を検討、またはコンセプト モデル ナゲットまたはカテゴリー モデル ナゲットのいずれかを作成できます。このモデル作成ノードを実行すると、内部の言語学的抽出エンジンは、自然言語処理手法を使用して、コンセプト、パターンまたはカテゴリー (あるいはそのすべて) を抽出して構成します。

テキストマイニングノードを実行し、「直接生成」オプションを使用して、コンセプトモデルナゲットまたはカテゴリモデルナゲットを自動的に作成できます。また、コンセプトを抽出、カテゴリを作成、および言語リソースを調整するだけでなく、テキストリンク分析を実行してクラスターを検証できる「インタラクティブに作成」モードを使用して、より実践的な探索的アプローチを使用できます。詳しくは、24ページの『テキストマイニング・ノード:「モデル」タブ』のトピックを参照してください。

このノードは、IBM SPSS Modeler ウィンドウの下部にあるノード・パレットの IBM SPSS Modeler Text Analytics タブにあります。詳しくは、8ページの『IBM SPSS Modeler Text Analytics ノード』のトピックを参照してください。

要件: テキストマイニング・モデル作成ノードは、Web フィールド ノード、ファイル リスト ノード、または標準的な入力ノードのいずれかからテキスト・データを受け入れます。このノードは IBM SPSS Modeler Text Analytics と共にインストールされており、IBM SPSS Modeler Text Analytics パレット上で使用できます。

注: このノードは、古いバージョンの製品に付属していた、テキスト抽出ノードに代わるものです。これらのノードまたはモデルナゲットを使用する古いストリームを使用する場合、新しいテキストマイニングノードを使用してストリームを再作成する必要があります。

テキストマイニング・ノード:「フィールド」タブ

「フィールド」タブを使用して、コンセプトを抽出するデータのフィールド設定を指定します。大きいデータセットを扱う場合は、処理時間を短くするために、このノードから上流でサンプル・ノードを使用するようにしてください。詳しくは、30ページの『時間を削減する上流のサンプリング』のトピックを参照してください。

設定できるパラメーターを次に示します。

ID フィールド: テキストレコードの識別子を含むフィールドを選択します。識別子は整数でなければなりません。ID フィールドは、各テキスト・レコードのインデックスとして機能します。テキスト・フィールドがマイニングされるテキストを示す場合、ID フィールドを使用します。

テキスト フィールド: マイニングするテキストが含まれているフィールドを選択します。このフィールドはデータ・ソースによって異なります。

言語フィールド (Language field): 2文字の ISO 言語 ID を含むフィールドを選択します。フィールドを選択しなかった場合、各ドキュメントの言語は、指定されたテンプレートの言語であると想定されます。

ドキュメント タイプ: ドキュメント・タイプは、テキストの構造を指定します。次に示すタイプの1つを選択します。

- **フル テキスト:** このオプションは、多くのドキュメントまたはテキスト・ソースに使用します。テキストのセット全体をスキャンして抽出します。他のオプションとは異なり、このオプションに追加設定はありません。
- **構造のあるテキスト:** このオプションは、参考文献形式、特許、特定および分析できる通常の構造を含むファイルに使用します。このドキュメントタイプを使用して、抽出プロセスのすべてまたは一部をスキップします。キーワードの区切り文字の定義、タイプの割り当て、および出現頻度の最小値の指定ができます。このオプションを選択する場合、「設定」ボタンをクリックして、「ドキュメント設定」ダイアログボックスの「構造のあるテキストの書式設定」領域にテキストの区切り文字を入力します。詳しくは、22ページの『「フィールド」タブのドキュメント設定』のトピックを参照してください。

テキストの単位: 次の実行モードを選択します。

- **ドキュメント・モード:** 通信社からの記事など、短く意味的に同質のドキュメントに使用します。
- **パラグラフ モード:** Web ページおよびタグのないドキュメントに使用します。抽出プロセスでは、内部タグやシンタックスなどの特徴を利用して、ドキュメントを意味的に分割します。このモードを選択すると、パラグラフごとにスコアリングが適用されます。そのため、例えば、apple および orange が同じパラグラフで見つかった場合にのみ、規則「apple & orange」が当てはまります。

注: PDF ドキュメントからテキストを抽出する方法が原因で、これらのドキュメントでは「パラグラフモード」は機能しません。これは、抽出により復帰マーカーが抑制されるためです。

パラグラフ モードの設定: このオプションは、「パラグラフ モード」に「テキストの単位」オプションを設定した場合にのみ選択できます。抽出で使用する文字のしきい値を指定します。実際のサイズは、最も近いピリオドに丸められます。ドキュメント・コレクションのテキストから作成される単語の関連性を典型とするには、抽出サイズが小さすぎないように指定します。

- **最小:** 抽出で使用する文字の最小数を指定します。
- **最大値:** 抽出で使用する文字の最大数を指定します。

データ区分モード データ区分モードを使用して、データ型ノードの設定に基づいて区分するか、別のデータ区分を選択するかを選択します。データ区分によって、データを学習サンプルおよび検定サンプルに分割します。

「フィールド」タグのドキュメント設定

構造のあるテキストの書式設定

構造のあるデータがある、またはテキストの処理方法に規則を強制したいために抽出プロセスの全部または一部をスキップしたい場合、「ドキュメント設定」ダイアログ・ボックスの「構造のあるテキストの書式設定」セクションで「構造のあるテキスト」ドキュメント・タイプのオプションを使用し、テキストを含むフィールドまたはタグを宣言します。抽出されたキーワードは、宣言したフィールドまたはタグ (および下位タグ) に含まれるテキストからのみ派生します。宣言されていないフィールドまたはタグは無視されます。

言語処理が必要ではなく、また言語学的抽出エンジンが明示的な宣言に置き換えることができる場合があります。キーワードのフィールドがセミコロン (;) やコンマ (,) のような区切り文字で区切られている参考文献ファイルでは、2 つの区切り文字の間の文字列を抽出すれば十分です。そのため、フル抽出プロセスを中止し、代わりに特別な処理規則を定義して、キーワードの区切り文字を宣言し、タイプを抽出テキストに割り当てるか、抽出に最小出現頻度を設定します。

構造のあるテキストの要素を宣言する場合、次の規則を使用します。

- 1 行ごとに宣言できるのは、1 つのフィールド、タグ、または要素だけです。それらはデータ内にある必要はありません。
- 宣言では、大文字小文字を区別します。
- `<title id="1234">` のような属性のあるタグを宣言し、すべての変異形、またはこの場合すべての ID を追加したい場合、属性および終わりの山括弧 (`>`) を除いたタグ (例えば `<title)` を追加します。
- フィールド名またはタグ名の後にコロンを追加して、構造のあるテキストを示します。このコロンは、`author:` または `<place>` のように、フィールドまたはタグの直後、そして区切り文字、タイプ、または出現頻度値の前に追加してください。
- 複数のキーワードがフィールドまたはタグに含まれ、また区切り文字を使用して各キーワードを指定することを示すには、`author:;`、または `<section>;` のように、コロンの後に区切り文字を指定します。

- タイプをタグの内容に割り当てるには、`author:,Person` または `<place>;Location` のように、コロンおよび区切り文字の後にタイプ名を指定します。リソース・エディターに表示されるとおりの名前を使用してタイプを宣言します。
- フィールドまたはタグの最小出現頻度を定義するには、`author:,Person1` または `<place>;Location5` のように、行の最後に数字を指定します。n は、定義する出現頻度を示し、フィールドまたはタグ内のキーワードは、抽出するドキュメントまたはレコードのセット全体で少なくとも n 回出現する必要があります。また、区切り文字を定義する必要もあります。
- コロンを含むタグがある場合、コロンの前に円記号を追加し、宣言が無視されないようにします。例えば、`<topic:source>` というフィールドがある場合、`<topic%:source>` のように入力します。

シンタックスを説明するために、次のような、繰り返し出現する参考文献のフィールドがあると仮定しましょう。

```
author:Morel, Kawashima
abstract:This article describes how fields are declared.
publication:Text Mining Documentation
datepub:March 2010
```

この例では、抽出プロセスが作者および要約に焦点を当てるようにしたいが、残りの内容は無視したい場合、次のフィールドのみを宣言します。

```
author:,Person1
abstract:
```

この例の場合、`author:,Person1` フィールドの宣言は、フィールドの内容についての言語処理が中断したことを示しています。代わりに、作者フィールドで、名前を区切るためにコンマを使用して複数の名前を含むことを指定します。そしてこれらの名前は「人名」タイプに割り当て、またドキュメントまたはレコードのセット全体で 1 回以上名前が出現した場合は抽出する必要があることを宣言しています。フィールド `abstract:` が他に宣言せずに表示されているため、抽出時にフィールドがスキャンされ、標準的な言語処理およびタイプ指定が行われます。

XML テキストの書式設定

抽出プロセスを特定の XML タグ内のテキストのみに制限したい場合、「ドキュメント設定」ダイアログ・ボックスの「XML テキストの書式設定」セクションで「XML テキスト」ドキュメント・タイプのオプションを使用し、そのテキストを含むタグを宣言します。抽出されたキーワードは、これらのタグまたは下位タグに含まれるテキストからのみ派生します。

重要: 抽出プロセスをスキップしてキーワードの区切り文字に規則を指定する場合、タイプを抽出したテキストに割り当てるか抽出したキーワードに出現頻度を指定し、次で説明する「構造のあるテキスト」オプションを使用します。

XML テキストの書式設定のタグを宣言する場合、次の規則を使用します。

- 1 行ごとに宣言できるのは、1 つの XML タグだけです。
- タグの要素は、大文字小文字を区別します。
- タグに `<title id="1234">` のような属性があり、すべての変異形、またはこの場合すべての ID を追加したい場合、属性および終わりの山括弧 (`>`) を除いたタグ (例えば `<title`) を追加します。

シンタックスを説明するために、次のような XML ドキュメントがあると仮定しましょう。

```
<section>Rules of the Road
  <title id="01234">Traffic Signals</title>
  <p>Road signs are helpful.</p>
</section>
<p>Learning the rules is important.</p>
```

この例の場合、次のタグを宣言します。

```
<section>
<title
```

この例では、タグ `<section>` を宣言しているため、このタグのテキストと入れ子になっているタグ、交通信号および道路標識は役立ちますが抽出プロセスでスキャンされます。ただし、タグ `<p>` が明示的に宣言されず、宣言されたタグの入れ子にもなっていないため、規則について学ぶことは重要です。は無視されません。

テキスト マイニング・ノード: 「モデル」 タブ

「モデル」タブを使用して、ノード出力の作成方法と一般的なモデル設定を指定します。

設定できるパラメーターを次に示します。

モデル名: 対象または ID フィールド (そのようなフィールドが指定されていない場合はモデル タイプ) に基づいてモデル名を自動的に生成するか、カスタム名を指定することができます。

データ区分データを使用。 データ区分フィールドが定義されている場合、このオプションでは学習用データ区分からのデータのみがモデル構築に使用されるようにします。

ビルド モード このテキスト マイニング・ノードでストリームが実行される場合にモデル ナゲットが作成される方法を指定します。また、コンセプトを抽出、カテゴリを作成、および言語リソースを調整するだけでなく、テキスト リンク分析を実行してクラスターを検証できる「インタラクティブに作成」 モードを使用して、より実践的な探索的アプローチを使用できます。

- **インタラクティブに作成:** ストリームを実行する場合、このオプションでインタラクティブ ワークベンチが起動します。インタラクティブ ワークベンチでは、コンセプトやパターンの抽出、抽出結果の探索および調整、カテゴリの作成および微調整、言語リソース (テンプレート、類義語、タイプ、ライブラリなど) の微調整、カテゴリ モデル ナゲットの作成を行うことができます。詳しくは、25 ページの『インタラクティブに作成』のトピックを参照してください。
- **直接生成** ストリームの実行時、モデルが自動的に作成され、「モデル」パレットに追加されることを指示します。インタラクティブ・ワークベンチとは異なり、ノードで定義された設定のほか、実行時に必要な追加の操作はありません。このオプションを選択すると、作成するモデルのタイプを定義できるモデル特有のオプションが表示されます。詳しくは、26 ページの『直接生成』のトピックを参照してください。

大規模なモデルを **AS** に格納: IBM SPSS Analytic Server に接続している場合にこのオプションを選択すると、モデルがサーバーにリモート保存されます。

注: サーバーで構築して格納したモデルは、いずれもそのサーバーのみでスコアリングできます。このようなモデルを含むインタラクティブ ワークベンチ セッションを再開するには、セッションの作成に使用した、元のサーバーに接続する必要があります。

リソースのコピー元: テキスト マイニング時、抽出は、「エキスパート」タブの設定だけでなく、言語リソースに基づいて行われます。これらのリソースは、抽出時のテキストの処理方法の基本として機能し、コンセプト、タイプ、そしてときには TLA パターンを取得します。リソース・テンプレートまたはテキス

ト分析パッケージのいずれかから、リソースをテキスト マイニングモデル作成ノードにコピーできます。いずれかを選択して「読み込み」をクリックし、リソースのコピー元となるパッケージまたはテンプレートを定義します。読み込んでいるときに、リソースのコピーがノードに保存されます。そのため、更新されたテンプレートまたは TAP を使用したい場合、ここで、インタラクティブ・ワークベンチ・セッションで再読み込みを行う必要があります。リソースがコピーされ、読み込まれた日時が、ノードに表示されます。詳しくは、27 ページの『テンプレートおよび TAP からのリソースのコピー』のトピックを参照してください。

テキストの言語。マイニングされるテキストの言語を識別します。ノードでコピーされたリソースが、表示される言語オプションを制御します。リソースを調整した言語を選択してください。

インタラクティブに作成

テキスト マイニング・モデル作成ノードの「モデル」タブで、モデル ナゲットのビルド モードを選択できます。「インタラクティブに作成」を選択すると、ストリーム実行時にインタラクティブ・インターフェースが開きます。このインタラクティブ・ワークベンチで、次の作業を行えます。

- コンセプトおよびタイプなどの抽出結果を抽出および検証し、テキスト・データで目立つキーワードを探索します。
- さまざまな手法を使用してコンセプト、タイプ、TLA パターンおよび規則からカテゴリを作成および展開して、ドキュメントおよびレコードをこれらのカテゴリにスコアリングできるようにします。
- 言語リソース (リソース・テンプレート、ライブラリー、辞書、類義語など) を調整し、コンセプトが抽出、検証、調整されるインタラクティブ プロセスによって結果が改善されるようにします。
- テキスト リンク分析 (TLA) を実行し、検出された TLA パターンを使用して、よりよいカテゴリモデル ナゲットを作成します。テキスト リンク分析ノードには、同じ探索オプションやモデル作成機能はありません。
- クラスターを生成して、検証ペインで新しい関係性を探索し、コンセプト、タイプ、パターン、およびカテゴリの間関係性を検証します。
- IBM SPSS Modeler の「モデル」パレットに調整済みカテゴリ モデル ナゲットを生成し、それらを他のストリームで使用します。

注: IBM SPSS Collaboration and Deployment Services ジョブを作成する場合は、インタラクティブ モデルは作成できません。

前回更新したノードのセッション作業 (カテゴリ、TLA、リソースなど) を使用: インタラクティブ・ワークベンチ・セッションで作業している場合、セッション・データ (抽出パラメーター、リソース、カテゴリ定義など) でノードを更新できます。 「セッション作業を使用」 オプションを選択すると、保存されたセッション・データを使用してインタラクティブ・ワークベンチを再起動できます。初めてこのノードを使用する場合、セッション・データが保存されていないため、このオプションは無効になります。このオプションを使用できるセッション・データでのノードの更新方法については、78 ページの『モデル作成ノードの更新および保存』を参照してください。

このオプションを使用してセッションを起動すると、前回インタラクティブ・ワークベンチ・セッションからノードを更新した時の抽出設定、カテゴリ、リソースおよびその他の作業が、次回セッションを起動したときに使用できます。保存したセッションデータはこのオプションで使用されるため、下のテンプレートからコピーしたリソースなど、特定の内容やその他のタブが無効となり、無視されます。このオプションを使用せずにセッションを起動すると、定義されたおりのノードの内容のみが使用されます。つまり、ワークベンチで実行した前回の作業は使用できなくなります。

注:抽出結果を「セッション作業を使用」 オプションを使用してキャッシュした後、ストリームの入力ノードを変更する際、抽出結果を更新する場合にインタラクティブ・ワークベンチ・セッションが起動したら新しい抽出を実行する必要があります。

抽出をスキップしてキャッシュデータおよび結果を再利用: インタラクティブ・ワークベンチ・セッションで、キャッシュされた抽出結果およびデータを再利用できます。このオプションは、セッションが起動したときに実行されるまったく新しい抽出を待機するのではなく、時間を節約して抽出結果を再利用したい場合に特に役立ちます。このオプションを使用するには、インタラクティブ・ワークベンチ・セッションからこのノードを事前に更新し、オプション「セッション作業を保存し、再利用するために抽出結果とともにテキストデータをキャッシュに格納する」を選択する必要があります。このオプションを使用できるセッション・データでのノードの更新方法については、78 ページの『モデル作成ノードの更新および保存』を参照してください。

セッションの開始: インタラクティブ・ワークベンチ・セッションの起動時に最初に表示したいビューおよび実行したい操作を示します。開始時のビューに関係なく、セッション内の任意のビューに一度切り替えることができます。

- 抽出結果を使用してカテゴリーを作成: カテゴリーとコンセプト・ビューでインタラクティブ・ワークベンチを起動し、必要に応じて抽出を実行します。このビューでは、カテゴリーを作成し、カテゴリー・モデルを生成できます。また、別のビューに切り替えることもできます。詳しくは、67 ページの『第 7 章 インタラクティブ・ワークベンチ・モード』のトピックを参照してください。
- テキスト リンク分析 (TLA) 結果を探索: まず、意見またはテキスト リンク分析ビューの他のリンクなど、テキスト内のコンセプト間の関係性を抽出および特定します。このオプションを使用して結果を抽出するには、TLA パターン規則を含むテンプレートまたはテキスト分析パッケージを選択する必要があります。より大きいデータセットを扱う場合、TLA 抽出に時間がかかる場合があります。この場合、上流でサンプル・ノードの使用を検討する必要があります。詳しくは、149 ページの『第 11 章 テキスト リンク分析の検証』のトピックを参照してください。
- 共通語クラスターを分析: このオプションはクラスター・ビューで起動し、古い抽出結果を更新します。このビューで、共通語クラスター分析を実行し、一連のクラスターを作成できます。共通語クラスターリングは、まず指定されたレコードまたはドキュメントの共起に基づいて 2 つのコンセプト間のリンク値の強度を評価し、最後に強くリンクしたコンセプトをクラスターにグループ化するプロセスです。詳しくは、67 ページの『第 7 章 インタラクティブ・ワークベンチ・モード』のトピックを参照してください。

直接生成

テキスト マイニング・モデル作成ノードの「モデル」タブで、モデル ナゲットのビルド モードを選択できます。「直接生成」を選択すると、ノードでオプションを設定し、ストリームを実行できます。出力はコンセプト モデル ナゲットで、「モデル」パレットに直接投入されます。インタラクティブ・ワークベンチとは異なり、ノードのオプションで定義された出現頻度設定のほか、実行時に必要な追加の操作はありません。

モデルが含む最大コンセプト数: 自動的にモデルを作成する場合 (非インタラクティブ) にのみ適用され、コンセプト・モデルを作成することを示します。また、このモデルには、指定した数以下のコンセプトが含まれることも示します。

- 最も高い頻度に基づいてコンセプトをチェックする。上位のコンセプト:チェックされるコンセプトの数です。最も頻度の高いコンセプトからチェックします。ここで、頻度は、ドキュメント/レコードのセット全体の中でコンセプト (およびすべての基本キーワード) が出現する回数を示します。レコード内にコンセプトが複数回出現する場合があるため、この数値がレコード数を上回る場合があります。

- 多くのレコードに出現するコンセプトはチェックを外す。レコードの割合。レコード数の割合が、指定した数を上回るコンセプトのチェックを解除します。このオプションは、テキストまたはすべてのレコードで頻繁に出現するが、分析においては重要でないコンセプトを除外する場合に役立ちます。

スコアリングの速度の最適化。このオプションはデフォルトで選択され、作成されたモデルがコンパクトで高速でスコアリングするようにします。このオプションの選択を解除すると、より低速でスコアリングを行う大規模なモデルが作成されます。ただし、大規模なモデルの場合、生成されたコンセプト・モデルで最初に表示されるスコアは、モデル ナゲットで同じテキストをスコアリングした場合に取得されるスコアと同じになります。

テンプレートおよび TAP からのリソースのコピー

テキスト マイニング時、抽出は、「エキスパート」タブの設定だけでなく、言語リソースに基づいて行われます。これらのリソースは、抽出時のテキストの処理方法の基本として機能し、コンセプト、タイプ、そしてときには TLA パターンを取得します。リソースをリソース・テンプレートからこのノードにコピーすることができます。また、テキスト マイニング・ノードを使用している場合は、テキスト分析パッケージ (TAP) を選択することもできます。

デフォルトでは、ノードを領域に追加すると、製品のライセンスされた言語の基本テンプレートからノードに、リソースがコピーされます。複数の言語のライセンスがある場合、最初に選択された言語を使用して、自動的に読み込むテンプレートを決定します。

読み込んでいるときに、選択したリソースのコピーがノードに保存されます。テンプレートまたは TAP の内容のみがコピーされますが、テンプレートまたは TAP 自体はノードにリンクしません。つまり、このテンプレートまたは TAP を後で更新すると、これらの更新が自動的にノードで使用できることはありません。ノードに読み込まれたリソースは、テンプレートまたは TAP のコピーが再読み込みされないかぎり、またはテキスト マイニング・ノードを更新して「セッション作業を使用」オプションを選択しないかぎり、かならず使用されます。「セッション作業を使用」の詳細は、このトピックを参照してください。

テンプレートまたは TAP を選択する場合、テキスト・データと言語が同じものを選択してください。ライセンスが付与された言語でのみテンプレートまたは TAP を使用できます。テキスト リンク分析を実行したい場合、TLA パターンを含むテンプレートを選択する必要があります。テンプレートに TLA パターンが含まれている場合、「リソース・テンプレートを読み込む」ダイアログ・ボックスの「TLA」列にアイコンが表示されます。

注: TAP をテキスト リンク分析ノードに読み込むことはできません。

リソース・テンプレート

リソース・テンプレートは、特定のドメインまたは使用向けに調整された、事前定義済みライブラリーおよび詳細な言語リソースおよび非言語リソースです。テキスト マイニング・モデル作成ノードでは、ノードをストリームに追加するときには基本テンプレートのリソースのコピーが既にノードに読み込まれています。ただし、テンプレートを変更するか、「リソース テンプレート」または「テキスト分析パッケージ」を選択して「読み込み」をクリックし、テキスト分析パッケージを読み込むことができます。テンプレートの場合、「リソース テンプレートの読み込み」ダイアログ・ボックスでテンプレートを選択できます。

注: 必要なテンプレートがリストに表示されないにもかかわらず、コンピューターにコピーがエクスポートされている場合、今すぐインポートできます。このダイアログ・ボックスからエクスポートして、他のユーザーと共有することもできます。詳しくは、175 ページの『テンプレートのインポートおよびエクスポート』のトピックを参照してください。

テキスト分析パッケージ (TAP)

テキスト分析パッケージ (TAP) は、1 つまたは複数の事前定義されたカテゴリのセットとまとめられた、ライブラリーと高度な言語リソースおよび非言語リソースの事前定義されたセットです。IBM SPSS Modeler Text Analytics では、英語テキスト向けの事前に作成された TAP がいくつか用意され、それぞれが特定のドメイン向けに調整されています。これらの TAP を編集できませんが、それらを使用してカテゴリ・モデル作成を開始できます。インタラクティブ・セッションで独自の TAP を作成することもできます。詳しくは、137 ページの『テキスト分析パッケージの読み込み』のトピックを参照してください。

注: TAP をテキスト リンク分析ノードに読み込むことはできません。

「セッション作業を使用」 オプションの使用 (「モデル」 タブ)

「モデル」 タブでリソースがノードにコピーされますが、後からインタラクティブ・セッションでリソースを変更し、これら最近の変更でテキスト マイニング・モデル作成ノードの更新が必要な場合があります。この場合、テキスト マイニング・モデル作成ノードの「モデル」 タブの「セッション作業を使用」 オプションを選択します。

「セッション作業を使用」 を選択すると、ノードの「読み込み」 ボタンが無効となり、インタラクティブ・ワークベンチのこれらのリソースが、以前ここで読み込まれたリソースの代わりに使用されることを示します。

「セッション作業を使用」 オプションで選択したリソースに変更を行うには、リソース・エディタービューを使用して、インタラクティブ・ワークベンチ・セッション内で直接リソースを編集または切り替えることができます。詳しくは、173 ページの『読み込み後のノード・リソースの更新』のトピックを参照してください。

テキスト マイニング・ノード:「エキスパート」 タブ

「エキスパート」 タブには、テキストの抽出方法および処理方法に影響を与える高度なパラメーターがあります。このダイアログ・ボックスのパラメーターは、抽出プロセスの基本的な操作、そしていくつかの高度な操作を制御します。ただし、使用できるオプションの部分のみを示します。また、抽出結果に影響を与える言語リソースやオプションも数多くあり、「モデル」 タブで選択するリソース・テンプレートによって制御します。詳しくは、24 ページの『テキスト マイニング・ノード:「モデル」 タブ』のトピックを参照してください。

注: 「モデル」 タブで保存されたインタラクティブ ワークベンチ情報に基づいて、「インタラクティブに作成」モードを選択した場合、このタブ全体が無効になります。この場合、抽出設定は、最近保存されたワークベンチ セッションから取得されます。

抽出時には、以下のパラメーターを設定できます。

グローバル頻度が次の値以上のコンセプトに抽出を制限: 抽出するために、単語または句が出現する必要がある最低限の回数を指定します。値に 5 を指定すると、抽出するこれらの単語または句が、レコードまたはドキュメントのセット全体で少なくとも 5 回出現するよう、制限します。

この制約を変更すると、抽出結果、つまり作成されるカテゴリに大きな違いが生じる場合があります。あるレストランのデータを処理し、このオプションの制約に1より大きい値を設定しないものとします。この場合、抽出結果がピザ (1)、薄いピザ (2)、ほうれん草のピザ (2)、および好きなピザ (2) となります。ただし、抽出のグローバル出現頻度を 5 以上に設定して抽出すると、これらのコンセプトのうち 3 つが取得されなくなります。代わりに、ピザが最も簡単な形で、この単語は考えられる候補として既に存在するため、ピザ (7) が取得されます。また、残りのテキストにピザという単語を含む他の句があるかどうかによ

て、7より大きい出現頻度がある場合があります。また、ほうれん草のピザがカテゴリーの記述子である場合、すべてのレコードをキャプチャーする代わりに、記述子としてピザの追加が必要な場合があります。このため、カテゴリーが既に作成されている場合は、注意してこの制約を変更してください。

これは抽出のみの機能であることに注意してください。つまり、テンプレートに用語が含まれる場合 (通常そのようになります) でテンプレートの用語がテキスト内で見つかった場合、その用語は頻度に関わらずインデックス付けされます。

例えば、コア・ライブラリーの <Location> タイプに「ロサンゼルス」が含まれている基本リソース・テンプレートを使用するとします。この場合、ドキュメント内での「ロサンゼルス」の出現回数が 1 回だけでも、ロサンゼルスがコンセプト・リストに含まれることとなります。これを回避するには、「グローバル頻度が次の値以上のコンセプトに抽出を制限」フィールドに入力された値以上の出現回数を持つコンセプトだけを表示するように、フィルターを設定する必要があります。

句読点エラーを調整: 抽出時に句読点エラー (不適切な使用方法など) を含むテキストを一時的に正規化し、コンセプトの抽出可能性を向上させます。自由記述式アンケートの回答、電子メール、CRM データなど、テキストが短く品質が悪い場合、またはテキストに略語が多く含まれている場合に特に役立ちます。

文字数が次の最小値以上のときにスペルを調整する: Fuzzy Grouping の手法を適用し、共通してミススペルのある単語またはスペルの近い単語を 1 つのコンセプトにグループ化できるようにします。Fuzzy Grouping アルゴリズムでは、最初の母音を除くすべての母音を一時的に抜き取った後抽出した単語から 2 つ/3 つの子音を抜き取り、それらを比較して、それらが同じで modeling と modelling が同じグループに分けられるかどうかを確認します。ただし、各キーワードが <Unknown> タイプを除いて、別のタイプに割り当てられた場合、Fuzzy Grouping 手法は適用されません。

Fuzzy Grouping を使用する前に必要な、語幹文字数の制限を定義することもできます。キーワード内の語幹文字数は、すべての文字を合計し、活用語尾、複合語キーワードの場合は区切り文字および前置詞を形成する文字を差し引いて計算します。例えば、キーワード exercises の語幹文字数は「exercise」という形式で 8 文字と数えられます。語末の s は活用語尾 (複数形) であるためです。同様に、apple sauce の語幹文字は 10 文字 (「apple sauce」)、そして manufacturing of cars の語幹文字は 16 文字 (「manufacturing car」) となります。この算出方法は、Fuzzy Grouping を適用するべきかどうかを確認するためにのみ使用されますが、単語がどのように一致するかについては影響を与えません。

注: 特定の単語が後で不適切にグループ化されていることが分かった場合、「拡張リソース」タブの **Fuzzy Grouping**: 例外 セクションで 明示的に宣言することによって、単語のペアをこの手法から除外できます。詳しくは、205 ページの『Fuzzy Grouping』のトピックを参照してください。

ユニタームを抽出 単語が複合語の一部でない限り、または名詞、またはスピーチ内の認識できない品詞である場合、このオプションは単一の単語 (ユニターム) を抽出します。

固有表現を抽出 電話番号、セキュリティ番号、時間、日付、通貨、数字、パーセント、電子メールアドレス、HTTP アドレスなどの固有表現を抽出します。「拡張リソース」タブの「固有表現: 設定」セクションで、特定の種類の固有表現を追加したり除外したりできます。不要な固有表現を無効にすることにより、抽出エンジンは処理時間を節約できます。詳しくは、210 ページの『構成』のトピックを参照してください。

大文字アルゴリズム キーワードの最初の文字が大文字である場合、組み込み辞書にない単純キーワードおよび複合キーワードを抽出します。このオプションには、最も適切な名詞を抽出するのに優れた方法があります。

可能な場合は、個人名の一部または全部をグループ化 テキスト内で別々の形式で同時に出現する名前をグループ化します。名前はテキストの始めでは完全な形式で、後は短い形式でのみ参照されるため、この機能が役立ちます。このオプションでは、タイプが <Unknown> のユニタームが、タイプ <Person> の複合キーワードの最後の単語に一致するようにします。例えば、*doe* があり、最初タイプが <Unknown> である場合、抽出エンジンは、<Person> タイプの複合キーワードに最後の単語として *doe* が含まれているかどうか (例: *john doe*) を確認します。ほとんどがユニタームとして抽出されることがないため、人の名前に適用されることはありません。

機能語による倒置を次の値を最大値として考慮する 倒置手法を適用する場合に指定されている場合がある非機能的単語の最大数を指定します。この倒置手法では、活用語尾に関係なく、含まれる非機能的単語 (*of* や *the* など) によってお互いに異なる類似した句をグループ化します。例えば、この値を最大 2 単語に設定し、*company officials* および *officials of the company* が抽出されたとします。この場合、両方の抽出キーワードは、*of the* が無視されると同じであるとみなされるため、最終コンセプト・リストに共にグループ化されます。

マルチタームをグループ化するときに派生関係を使用: ビッグデータを処理するときこのオプションを選択すると、派生規則を使用してマルチタームがグループ化されます。

注: テキスト リンク分析結果の抽出を有効にするには、「テキスト リンク分析 (TLA) 結果を探索」オプションでセッションを開始する必要があります。また、TLA 定義を含むリソースを選択する必要があります。「抽出設定」ダイアログで、インタラクティブ・ワークベンチ・セッション中に後から TLA 結果を抽出できます。詳しくは、82 ページの『データの抽出』のトピックを参照してください。

時間を削減する上流のサンプリング

大きいデータがある場合、特にインタラクティブ・ワークベンチ・セッションを使用している場合、処理時間は数分から数時間かかる場合があります。データのサイズが大きいほど、抽出およびカテゴリー化処理の時間がより長くなります。より効率的に作業するために、テキスト・マイニング・ノードの上流の IBM SPSS Modeler のサンプル・ノードを追加します。このサンプル・ノードを使用して、ドキュメントまたはレコードの小さい部分集合で無作為サンプルを取得し、最初のいくつかの通過を実行します。

小さいサンプルは、リソースの編集方法を決定し、すべてではなくても多くのカテゴリーを作成するのに適しています。小さいデータセットで実行し、結果が適切であれば、カテゴリー作成の手法をデータのセット全体に適用できます。作成したカテゴリーに適合しないドキュメントまたはレコードを検出し、必要に応じて調整を行うことができます。

注: サンプル・ノードは、標準的な IBM SPSS Modeler ノードです。

ストリーム内のテキスト マイニング・ノードの使用

テキスト マイニング モデル作成ノードを使用して、データにアクセスし、ストリーム内のコンセプトを抽出します。データベース ノードなどのように、データにアクセスするにはどのソース ノードも使用できます。可変長ファイル ノード、Web フィールド ノード、または固定長ファイル ノード。外部ドキュメント内のテキストの場合、ファイル リスト ノードを使用できます。

例 1: コンセプト モデル ナゲットを直接作成するファイル リスト ノードおよびテキスト マイニング・ノード

次の例では、テキスト マイニング モデル作成ノードと共にファイル リスト ノードを使用して、コンセプト モデル ナゲットを生成する方法が示されています。ファイル リスト ノード使用の詳細は、「11 ページの『ファイル リスト ノード』」を参照してください。

1. ファイル リスト ノード（「設定」タブ）。まず、このノードをストリームに追加して、テキスト・ドキュメントが保存されている場所を指定しました。テキスト マイニングを実行するすべてのドキュメントを含むディレクトリーを選択しました。
2. テキスト マイニング・ノード（「フィールド」タブ）: 次に、テキスト マイニング・ノードを ファイル リスト ノードに追加して接続しました。このノードで、入力形式、リソース・テンプレート、および出力形式を定義しました。ファイル リスト ノードから作成されたフィールド名を選択し、テキスト フィールドやその他の設定を選択しました。詳しくは、30 ページの『ストリーム内のテキスト マイニング・ノードの使用』のトピックを参照してください。
3. テキスト マイニング・ノード（「モデル」タブ）:次に「モデル」タブで、ビルド モードとリソースを選択し、このノードから直接コンセプト モデル ナゲットを生成しました。異なるリソース・テンプレートを選択するか、基本リソースを保持します。

例 2:インタラクティブにカテゴリーモデルを作成する Excel ファイルノードおよびテキストマイニングノード

この例では、テキスト マイニング・ノードがどのようにインタラクティブ・ワークベンチ・セッションを起動できるかを示しています。インタラクティブ・ワークベンチの詳細は、67 ページの『第 7 章 インタラクティブ・ワークベンチ・モード』を参照してください。

1. **Excel** ソース・ノード（「データ」タブ）。まず、このノードをストリームに追加して、テキストが保存されている場所を指定しました。
2. テキスト マイニング・ノード（「フィールド」タブ）: 次に、テキスト マイニング・ノードを追加して接続しました。最初のタブで入力形式を定義しました。入力ノードからフィールド名を選択しました。
3. テキスト マイニング・ノード（「モデル」タブ）: 次に、「モデル」タブで、インタラクティブにカテゴリー モデル ナゲットを作成し、抽出結果を使用して自動的にカテゴリーを作成するよう選択しました。この例では、リソースのコピーおよび一連のカテゴリーをテキスト分析パッケージから読み込みました。
4. インタラクティブ・ワークベンチ・セッション:次に、ストリームを実行し、インタラクティブ・ワークベンチ・インターフェースが開きました。抽出が実行された後、データの探索およびカテゴリーの改善を開始しました。

テキスト マイニング モデル ナゲット:コンセプト・モデル

テキスト マイニング・コンセプト モデル ナゲットは、「モデル」タブで「モデルを直接生成」のオプションを選択してテキスト マイニング・モデル・モードが正常に実行されると作成されます。テキスト マイニング・コンセプト モデル ナゲットは、コール・センターのメモ帳データなど、その他のテキスト・データの主要キーワードをリアルタイムで発見するために使用されます。

コンセプト モデル ナゲット自体は、タイプに割り当てられているコンセプトのリストを組み合わせます。そのモデルのいずれかまたはすべてのコンセプトを選択し、その他のデータに対してスコアリングができます。テキスト マイニング モデル ナゲットを含むストリームを実行すると、モデル作成の前に、テキスト マイニングモデル作成ノードの「モデル」タブで選択されたビルド モードに従って、新しいフィールドがデータに追加されます。詳しくは、32 ページの『コンセプト・モデル:「モデル」タブ』のトピックを参照してください。

モデル ナゲットが翻訳されたドキュメントを使用して生成された場合、翻訳された言語でスコアリングが実行されます。同様に、モデル ナゲットが英語で生成された場合、ドキュメントが英語に翻訳されるため、モデル ナゲットで翻訳言語を指定できます。

テキストマイニングモデルナゲットは生成時、モデルナゲットパレット (IBM SPSS Modeler ウィンドウの右上の「モデル」タブ) 内にあります。

結果の表示

モデルナゲットに関する情報を表示するには、モデルナゲットパレットでノードを右クリックし、コンテキスト・メニューから「参照」を選択します (ストリーム中のノードの場合は「編集」)。

モデルのストリームへの追加

モデルナゲットをストリームに追加するには、モデルナゲットパレット内でアイコンをクリックし、ノードを配置するストリーム領域をクリックします。アイコンを右クリックし、コンテキスト・メニューから「ストリームに追加」をクリックします。次に、ストリームをノードに接続すれば、データを渡して予測を生成する準備が整います。

注意: スコアリングナゲットを使用して、カテゴリモデルおよび使用するテンプレートの両方が含まれるモデル作成ノードを再生成したい場合は、スコアリングナゲットを生成する前に、モデル作成ノードの代わりに、TAP を作成してインタラクティブセッションでそれを使用することをお勧めします。

コンセプト・モデル: 「モデル」タブ

コンセプト・モデルの「モデル」タブには、抽出されたコンセプトのセットが表示されます。コンセプトは、各コンセプトに 1 行ずつのテーブル形式で表示されます。このタブでは、スコアリングに使用されるコンセプトを選択します。

注: 代わりにカテゴリモデルナゲットを生成すると、このタブには異なる情報が表示されます。詳しくは、41 ページの『カテゴリモデルナゲット: 「モデル」タブ』のトピックを参照してください。

一番左の列のチェック・ボックスで示されているように、デフォルトではすべてのコンセプトがスコアリングに選択されています。チェック・ボックスがオンの場合、コンセプトはスコアリングに使用されます。チェック・ボックスがオフの場合、コンセプトはスコアリングから除外されます。複数の行を選択して、選択部分のいずれかのチェック・ボックスをクリックすると、複数の行をオンにできます。

各コンセプトの詳細については、次の各列に表示された追加情報を参照してください。

コンセプト: 抽出された代表語句です。コンセプトが、このコンセプトに関連する基本キーワードのほか、コンセプト名を示す場合があります。コンセプトの一部である基本キーワードを確認するには、このタブ内の用語ペインを表示してコンセプトを選択し、ダイアログ・ボックスの下部にある該当するキーワードを確認します。詳しくは、34 ページの『コンセプト・モデルの基本キーワード』のトピックを参照してください。

グローバル: ここで、グローバル (出現頻度) は、ドキュメント/レコードのセット全体の中でコンセプト (およびすべての基本キーワード) が出現する回数を示します。

- **棒グラフ:** 棒グラフで表示されたこのコンセプトがテキスト・データに出現したグローバル出現頻度。棒の色は、タイプを視覚的に区別するためにコンセプトに割り当てられたタイプの色です。
- **%:** このコンセプトがテキスト・データに出現したグローバル出現頻度 (パーセント表示)。
- **N:** テキスト・データにおけるこのコンセプトの出現数。

ドキュメント: ここで、ドキュメントはドキュメント数、つまりコンセプト (およびすべての基本キーワード) が出現するドキュメントまたはレコード数を示します。

- 棒グラフ: 棒グラフとして表示されたこのコンセプトのドキュメント数。棒の色は、タイプを視覚的に区別するためにコンセプトに割り当てられたタイプの色です。
- %、パーセントで表示されたこのコンセプトのドキュメント数。
- N: このコンセプトを含むドキュメントまたはレコードの数。

タイプ: コンセプトが割り当てられるタイプ。各コンセプトについて、「グローバル」列および「ドキュメント」列は色付きで表示され、コンセプトに割り当てられたタイプを示します。タイプは、コンセプトの意味上のグループです。詳しくは、189 ページの『キーワード辞書』のトピックを参照してください。

コンセプトの作業

テーブル内のセルを右クリックすると、次のようなコンテキスト・メニューが表示されます。

- すべて選択: テーブル内のすべての行が選択されます。
- コピー: 選択したコンセプトがクリップボードにコピーされます。
- フィールドもコピー: 選択したコンセプトが列の見出しと共にクリップボードにコピーされます。
- 選択項目をチェック: スコアリングするこれらのコンセプトを含むテーブルで選択した行のすべてのチェック・ボックスをオンにします。
- 選択項目をチェック解除: テーブルで選択した行のすべてのチェック・ボックスをオフにします。
- すべてチェック: テーブルのすべてのチェック・ボックスをオンにします。これにより、すべてのコンセプトが最終的な出力に使用されます。
- すべてチェック解除: テーブルのすべてのチェック・ボックスをオフにします。コンセプトのチェックを解除すると、最終出力では使用されません。
- コンセプトを含む: 「コンセプトを含む」ダイアログ・ボックスが表示されます。詳しくは、『スコアリングのコンセプト追加のオプション』のトピックを参照してください。

スコアリングのコンセプト追加のオプション

それらのコンセプトをすばやく選択または選択解除するには、「コンセプトを含む」のツールバー・ボタンをクリックします。



図 1. 「コンセプトを含む」ツールバー・ボタン

このツールバー・ボタンをクリックすると「コンセプトを含む」ダイアログ・ボックスが開き、規則に基づいてコンセプトを選択できます。「モデル」タブでチェック マークの付いたすべてのコンセプトは、スコアリングに追加されます。このサブダイアログで規則を適用し、スコアリングに使用するコンセプトを変更します。

以下のオプションの中から選択することができます。

最も高い頻度に基づいてコンセプトをチェックする。上位のコンセプト:チェックされるコンセプトの数です。最もグローバル頻度の高いコンセプトからチェックします。ここで、頻度は、ドキュメント/レコードのセット全体の中でコンセプト (およびすべての基本キーワード) が出現する回数を示します。レコード内にコンセプトが複数回出現する場合があるため、この数値がレコード数を上回る場合があります。

ドキュメント数に基づいてコンセプトをチェックする。最小数。チェックするコンセプトに必要な最低限のドキュメント数です。ここで、ドキュメント数は、コンセプト (およびすべての基本キーワード) が出現するドキュメントの数を示します。

タイプに割り当てられたコンセプトをチェックする: ドロップダウン・リストからタイプを選択して、このタイプに割り当てられるすべてのコンセプトをチェックします。抽出時、コンセプトはタイプに自動的に割り当てられます。タイプは、コンセプトの意味上のグループです。タイプには、上位レベルのコンセプト、肯定語および否定語および識別子、コンテキスト識別子、人名、地名、組織名などが含まれます。詳しくは、189 ページの『キーワード辞書』のトピックを参照してください。

多くのレコードに出現するコンセプトはチェックを外す。レコードの割合。レコード数の割合が、指定した数を上回るコンセプトのチェックを解除します。このオプションは、テキストまたはすべてのレコードで頻繁に出現するが、分析においては重要でないコンセプトを除外する場合に役立ちます。

タイプに割り当てられたコンセプトはチェックを解除する: ドロップダウン・リストから選択したタイプと合致するコンセプトのチェックを解除します。

コンセプト・モデルの基本キーワード

テーブルで選択したコンセプトに定義されている基本キーワードが表示されます。ツールバーで基本キーワードの切り替えボタンをクリックすると、ダイアログの分割したパネルに基本キーワード表が表示されます。

これらの基本キーワードには、モデル ナゲットの生成に使用されるテキストにある抽出された複数形/単数形、置換キーワード、Fuzzy Grouping のキーワードなどのほか、テキスト内にあったかどうかに関係なく、言語リソースで定義された類義語が含まれます。



図 2. 「基本キーワードを表示」ツールバー・ボタン

注:基本キーワードのリストは編集できません。このリストは、すべて言語リソースで定義されている置換、類義語辞書、Fuzzy Grouping によって生成されます。コンセプトに基づいてキーワードをどのようにグループ化するか、またはそれらをどのように処理するかを変更するには、リソースで直接変更し (インタラクティブ・ワークベンチの リソース・エディター または テンプレート・エディター で編集でき、ロードに再読み込み)、ストリームを再実行して、更新された結果を持つ新しいモデル ナゲットを取得します。

基本キーワードまたはコンセプトを含むテーブル内のセルを右クリックすると、次のようなコンテキスト・メニューが表示されます。

- コピー: 選択したセルをクリップボードにコピーします。
- フィールドもコピー: 選択したセルが列の見出しと共にクリップボードにコピーされます。
- すべて選択: テーブル内のすべてのセルが選択されます。

コンセプト モデル: 「設定」 タブ

必要に応じて、「設定」タブを使用して、新しい入力データのテキスト・フィールド値を定義します。また、ここで出力のデータ・モデル (スコアリング・モード) も定義します。

注: このタブは、モデル ナゲットが領域内にある場合にのみ表示されます。「モデル」パレットでこのダイアログ・ボックスを使用している場合、このタブは存在しません。

スコアリング・モード:レコードとしてのコンセプト

このスコアリングモードで、新しいレコードが各 concept/document ペアに作成されます。通常、入力に比べて出力に多くのレコードがあります。

入力フィールドに加えて、次の新しいフィールドがデータに追加されます。

表 4. 「レコードとしてのコンセプト」の出力フィールド:

フィールド	説明
Concept	テキスト・データ・フィールドの抽出したコンセプト名が指定されます。
Type	地名または人名などの完全なタイプ名としてコンセプトのタイプを格納します。タイプは、コンセプトの意味上のグループです。詳しくは、189 ページの『キーワード辞書』のトピックを参照してください。
Count	テキスト本文 (レコード/ドキュメント) の該当するコンセプト (および基本キーワード) の出現数が表示されます。

このオプションを選択すると、「句読点エラーを調整」を除くすべてのオプションが無効になります。

スコアリング・モード:フィールドとしてのコンセプト

各入力フィールドのコンセプト・モデルで、新しいレコードが指定されたドキュメントで見つかった各コンセプトに作成されます。そのため、入力内と同じ数の出力レコードがあります。ただし、各レコード (行) では、「モデル」タブで選択された (チェック マークあり) コンセプトまたはカテゴリーに 1 つずつ新しいフィールド (列) があります。各コンセプト・フィールドの値は、このタブのフィールド値として「フラグ」または「度数」のどちらを選択するかによって異なります。

注: 非常に大規模なデータ・セットを使用している場合 (Db2 データベースの使用など)、「フィールドとしてのコンセプト」を使用すると、データ量が原因で処理上の問題が発生する可能性があります。この場合は、代わりに「レコードとしてのコンセプト」を使用することをお勧めします。

フィールド値: 各コンセプトの新しいフィールドに度数またはフラグ値を指定するかを選択します。

- フラグ。はい/いいえ、真/偽、T/F、または 1 と 2 など、出力に 2 つの値を持つフラグを取得します。ストレージ・タイプは、選択した値を反映するよう自動的に設定されます。例えば、フラグに数値を入力すると、自動的に整数値として処理されます。フラグ型では、文字列、整数、実数、または日付/時間のストレージ・タイプを利用することができます。「真 (True)」および「偽 (False)」のフラグ値を入力します。
- 度数: 指定したレコードにコンセプトが出現した回数を取得します。

フィールド名拡張子: フィールド名の拡張子を指定します。コンセプト名に加えてこの拡張子を使用して、フィールド名が生成されます。

- 追加方法: 拡張子をフィールド名に追加する場所を指定します。「接頭辞」を選択すると、文字列の頭に拡張子が追加されます。「接尾辞」を選択すると、文字列の終わりに拡張子が追加されます。

句読点エラーを調整: 抽出時に句読点エラー (不適切な使用方法など) を含むテキストを一時的に正規化し、コンセプトの抽出可能性を向上させます。自由記述式アンケートの回答、電子メール、CRM データなど、テキストが短く品質が悪い場合、またはテキストに略語が多く含まれている場合に特に役立ちます。

コンセプト モデル:「フィールド」タブ

「フィールド」タブは、必要に応じて、新しい入力データのテキスト フィールド値を定義します。

注: このタブは、モデル ナゲットがストリーム内に配置されているときにのみ表示されます。「モデル」パレットでこの出力を使用している場合、このタブは存在しません。

テキスト フィールド: マイニングするテキストが含まれているフィールドを選択します。このフィールドはデータ・ソースによって異なります。

ドキュメント タイプ: ドキュメント・タイプは、テキストの構造を指定します。次に示すタイプの 1 つを選択します。

- フル テキスト: このオプションは、多くのドキュメントまたはテキスト・ソースに使用します。テキストのセット全体をスキャンして抽出します。他のオプションとは異なり、このオプションに追加設定はありません。
- 構造のあるテキスト: このオプションは、参考文献形式、特許、特定および分析できる通常の構造を含むファイルに使用します。このドキュメントタイプを使用して、抽出プロセスのすべてまたは一部をスキップします。キーワードの区切り文字の定義、タイプの割り当て、および出現頻度の最小値の指定ができます。このオプションを選択する場合、「設定」ボタンをクリックして、「ドキュメント設定」ダイアログ ボックスの「構造のあるテキストの書式設定」領域にテキストの区切り文字を入力します。詳しくは、22 ページの『「フィールド」タブのドキュメント設定』のトピックを参照してください。

入力エンコード: テキスト・フィールドが「ドキュメントへのパス名」を示すよう指定した場合にのみ使用できます。デフォルトのテキスト・エンコードを指定します。指定された、または認識されたエンコードから ISO-8859-1 への変換が行われます。そのため、別のエンコードが指定されている場合であっても、抽出エンジンは処理前にテキストを ISO-8859-1 に変換します。ISO-8859-1 エンコード定義に一致しない文字は、スペースに変換されます。

テキストの言語: マイニングされるテキストの言語を識別します。これは、抽出中に選択したメインの言語になります。現在使用できないサポート言語のライセンス購入については、営業担当者に連絡してください。

コンセプト・モデル:「要約」タブ

「要約」タブには、モデルそのもの (分析フォルダー)、モデルで使用するフィールド (フィールド フォルダー)、モデル作成時に使用する設定 (ビルド設定 フォルダー)、およびモデルの学習 (学習の要約 フォルダー) についての情報を表示します。

モデル作成ノードを初めて参照した場合、「要約」タブのフォルダーは閉じられています。関心のある結果を表示するには、フォルダーの左側にある展開コントロールを使用するか、または「すべて展開」ボタンをクリックしてすべての結果を表示してください。見終わった結果を隠すには、展開コントロールを使用して特定のフォルダーを閉じるか、または「すべて閉じる」ボタンをクリックしてすべてのフォルダーを非表示にします。

ストリームでのコンセプト モデル ナゲットの使用

テキスト マイニング・モデル作成ノードを使用すると、コンセプト モデル ナゲットまたはカテゴリ モデル ナゲット (インタラクティブ・ワークベンチ・セッションを使用) のいずれかを生成できます。次の例では、単純なストリームでコンセプト・モデルの使用方法について示しています。

例:コンセプト モデル ナゲットを含む **Statistics** ファイル ノード

次の例は、テキストマイニング コンセプト モデル ナゲットの使用方法を表示しています。



図 3. ストリームの例:テキスト マイニング コンセプト モデル ナゲットを含む *Statistics* ファイル ノード

1. **Statistics** ファイル・ノード (「データ」タブ)。 まず、このノードをストリームに追加して、テキスト・ドキュメントが保存されている場所を指定しました。

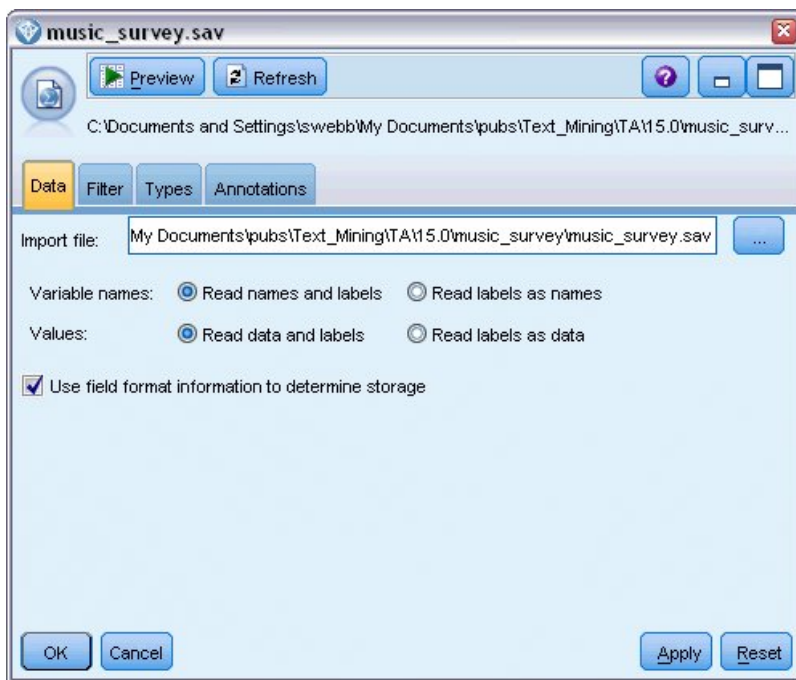


図 4. *Statistics* ファイル・ノード・ダイアログ・ボックス:「データ」タブ

2. テキスト マイニング コンセプト モデル ナゲット (「モデル」タブ): 次に、コンセプト モデル ナゲットを *Statistics* ファイル ノードに追加して接続しました。データのスコアリングに使用したいコンセプトを選択しました。

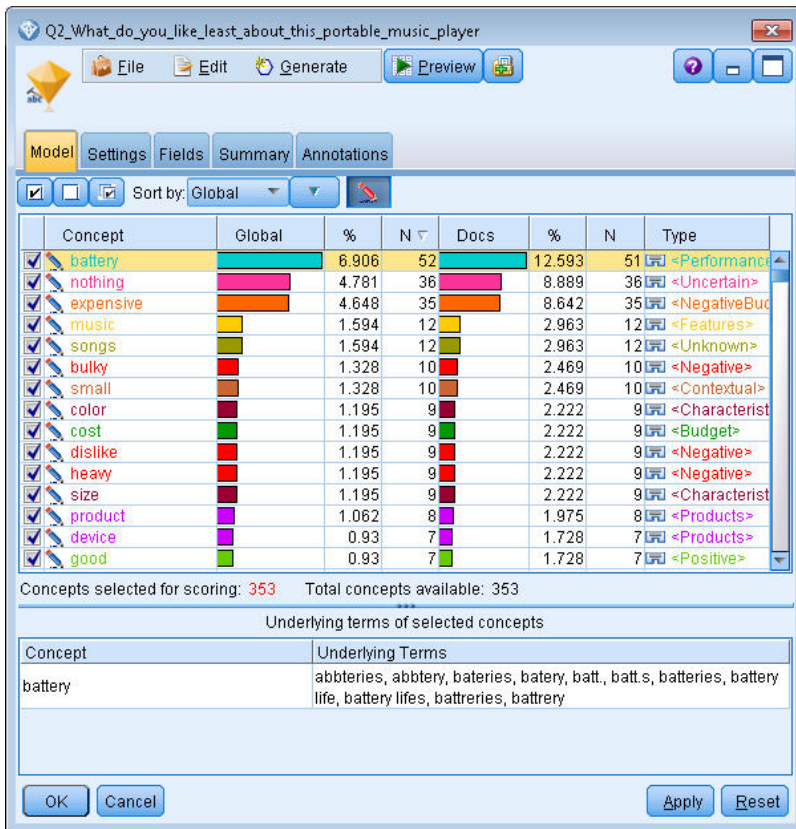


図 5. 「テキスト マイニング モデル ナゲット」 ダイアログ・ボックス: 「モデル」 タブ

3. テキスト マイニング コンセプト モデル ナゲット (「設定」 タブ): 次に出力形式を定義し、 「フィールドとしてのコンセプト」 を選択しました。1 つの新しいフィールドが、 「モデル」 タブで選択した各コンセプトの出力に作成されます。各フィールド名は、コンセプト名と、接頭辞 「Concept_」 で成り立っています。

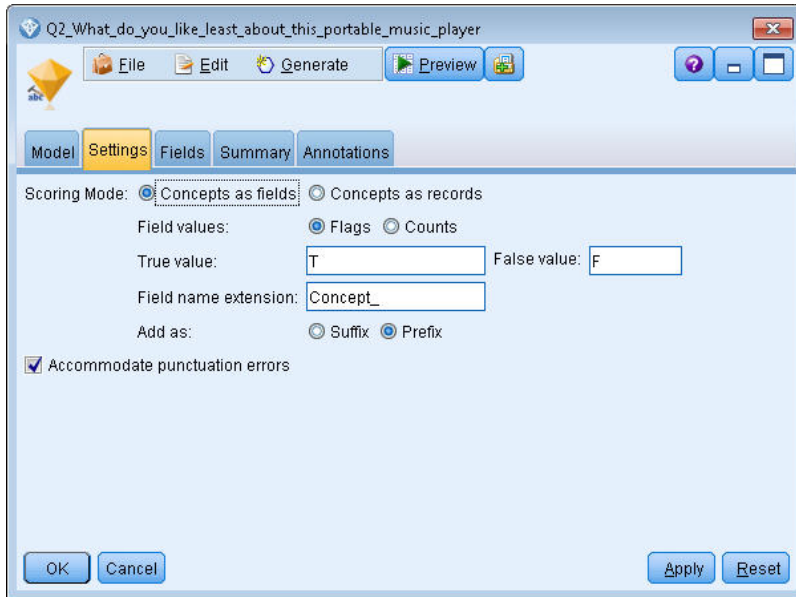


図 6. 「テキスト マイニング コンセプト モデル ナゲット」 ダイアログ・ボックス:「設定」タブ

4. テキスト マイニング コンセプト モデル ナゲット (「フィールド」タブ): 次に、テキスト・フィールド「Q2_What_do_you_like_least_about_this_portable_music_player」を選択しました。次に、テキスト・フィールド Q2_What_do_you_like_least_about_this_portable_music_player を選択します。これは Statistics ファイル ノードに由来しています。また、オプション「テキスト・フィールドの表示: 実際のテキスト」を選択しました。

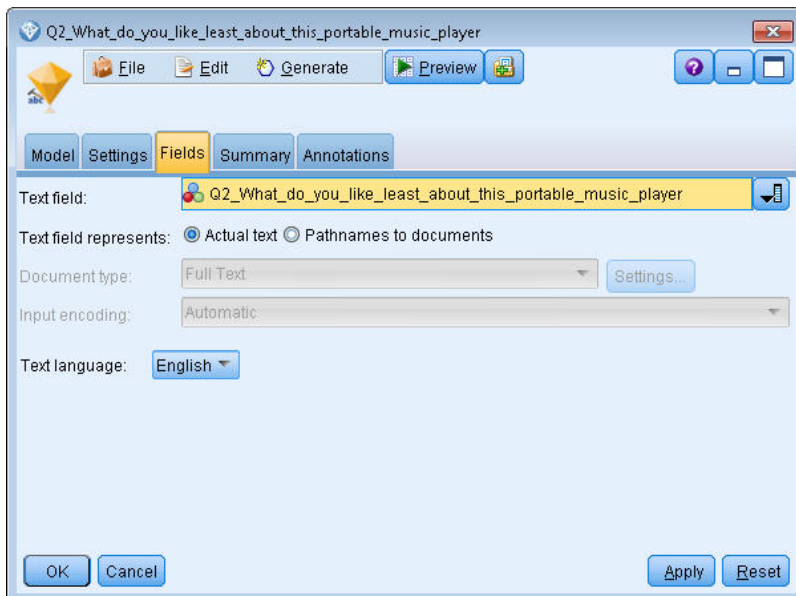


図 7. 「テキスト マイニング コンセプト モデル ナゲット」 ダイアログ・ボックス:「フィールド」タブ

5. テーブル・ノード: 次に、 テーブル・ノードを接続して結果を表示し、ストリームを実行しました。テーブル出力が画面上に表示されます。

Respondent_ID	Q1_WV...	Q2_What_do_you_like_least_about_this_portable_music_player	Concept_reliable	Concept_downloading...	Concept_white color	Concept_limited
1	1	little, it... expensive	F	F	F	F
2	2	The ba... The screen is hard to see when outside.	F	F	F	F
3	3	cost a... difficult software	F	F	F	F
4	4	Having... Nothing, I love it!	F	F	F	F
5	5	The sh... Battery life seems shorter than advertised.	F	F	F	F
6	6	Batter... Ubiquitousness; everyone has one.	F	F	F	F
7	7	I like it... I wish the 40GB model was still available. I have a 20GB model and need more memory.	F	F	F	F
8	8	portabi... it doesn't have a light.	F	F	F	F
9	9	Small, ... Nothing, I love it.	F	F	F	F
10	10	Able t... It is in the shop due to a hardware failure.	F	F	F	F
11	11	It's por... smudges on the display	F	F	F	F
12	12	Living i... Battery life	F	F	F	F
13	13	mobility Technical difficulties setting it up initially and managing the library of songs on my PC.	F	F	F	F
14	14	I like th... it is a little heavy, and the battery life isn't long enough.	F	F	F	F
15	15	It hold... Battery life.	F	F	F	F
16	16	It's fun... nothing	F	F	F	F
17	17	its cool battery	F	F	F	F
18	18	lots of ... it was very expensive	F	F	F	F
19	19	Others... I find the controls hard to use.	F	F	F	F
20	20	lightw... so small afraid I'll lose it easily	F	F	F	F

図 8. コンセプト・フラグを表示するためにスクロールしたテーブル出力

テキスト マイニング モデル ナゲット:カテゴリー・モデル

テキスト マイニング カテゴリー モデル ナゲットは、インタラクティブ・ワークベンチからカテゴリー・モデルが生成されると作成されます。このモデル作成ナゲットには、一連のカテゴリーが含まれ、その定義はコンセプト、タイプ、TLA パターンおよびカテゴリー規則で構成されています。ナゲットを使用して、アンケートの回答、ブログ エントリー、その他の Web フィード、およびその他テキスト・データをカテゴリー化します。

モデル作成ノードでインタラクティブ・ワークベンチ・セッションを起動すると、カテゴリー・モデルを生成する前に抽出結果を探索し、リソースを調整、カテゴリーを調整できます。テキスト マイニング モデル ナゲットを含むストリームを実行すると、モデル作成の前に、テキスト マイニング モデル作成ノードの「モデル」タブで選択されたビルド モードに従って、新しいフィールドがデータに追加されます。詳しくは、41 ページの『カテゴリー モデル ナゲット:「モデル」タブ』のトピックを参照してください。

モデル ナゲットが翻訳されたドキュメントを使用して生成された場合、翻訳された言語でスコアリングが実行されます。同様に、モデル ナゲットが英語で生成された場合、ドキュメントが英語に翻訳されるため、モデル ナゲットで翻訳言語を指定できます。

テキスト マイニング モデル ナゲットは生成時、モデル ナゲット パレット (IBM SPSS Modeler ウィンドウの右上の「モデル」タブ) 内にあります。

結果の表示

モデル ナゲットに関する情報を表示するには、モデル ナゲット パレットでノードを右クリックし、コンテキスト・メニューから「参照」を選択します (ストリーム中のノードの場合は「編集」)。

モデルのストリームへの追加

モデル ナゲットをストリームに追加するには、モデル ナゲット パレット内でアイコンをクリックし、ノードを配置するストリーム領域をクリックします。アイコンを右クリックし、コンテキスト・メニューから「ストリームに追加」をクリックします。次に、ストリームをノードに接続すれば、データを渡して予測を生成する準備が整います。

注意: スコアリング ナゲットを使用して、カテゴリー モデルおよび使用するテンプレートの両方が含まれるモデル作成ノードを再生成したい場合は、スコアリング ナゲットを生成する前に、モデル作成ノードの代わりに、TAP を作成してインタラクティブ セッションでそれを使用することをお勧めします。

カテゴリー モデル ナゲット: 「モデル」 タブ

カテゴリー・モデルの場合、「モデル」タブの左側にカテゴリー・モデルのカテゴリーのリスト、右側に選択したカテゴリーの記述子のリストが表示されます。各カテゴリーは、多くの記述子で構成されています。選択した各カテゴリーについて、関連する記述子がテーブルに表示されます。記述子には、コンセプト、カテゴリー規則、タイプ、および TLA パターンが含まれます。記述子が示す内容の例のほか、各記述子のタイプも表示されます。

このタブでは、スコアリングに使用されるカテゴリーを選択します。カテゴリー・モデルについて、ドキュメントおよびレコードがカテゴリーにスコアリングされます。ドキュメントまたはレコードにテキストまたは基本キーワードの 1 つまたは複数の記述子がある場合、そのドキュメントまたはレコードは記述子が含まれるカテゴリーに割り当てられます。これらの基本キーワードには、モデル ナゲットの生成に使用されるテキストにある抽出された複数形/単数形のキーワード、置換キーワード、Fuzzy Grouping のキーワードなどのほか、テキスト内にあったかどうかに関係なく、言語リソースで定義された類義語が含まれます。




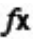
注:代わりにコンセプト モデル ナゲットを生成すると、このタブには異なる結果が表示されます。詳しくは、32 ページの『コンセプト・モデル: 「モデル」 タブ』のトピックを参照してください。

カテゴリー・ツリー

各カテゴリーの詳細については、該当するカテゴリーを選択し、そのカテゴリーの記述子に表示される情報を参照してください。記述子については、次の情報が表示されます。

- 記述子名:このフィールドには記述子名のほか記述子の種類を示すアイコンが指定されます。

表 5. 記述子アイコン

	コンセプト		TLA パターン
	データ型		カテゴリー規則

- タイプ: このフィールドには記述子のタイプ名が指定されます。タイプは、組織名、商品名、肯定的な意見など、類似したコンセプトの集合 (意味的なグループ) です。条件規則はタイプに割り当てられません。
- 詳細: その記述子に含まれる内容のリストが指定されます。合致数によっては、ダイアログ・ボックスのサイズ制限のため、各記述子のリスト全体を表示できない場合があります。

カテゴリーの選択およびコピー

左のパネルのチェック・ボックスで示されているように、デフォルトではすべてのカテゴリーがスコアリングに選択されています。チェック・ボックスがオンの場合、カテゴリーはスコアリングに使用されます。チ

チェック・ボックスがオフの場合、カテゴリーはスコアリングから除外されます。複数の行を選択して、選択部分のいずれかのチェック・ボックスをクリックすると、複数の行をオンにできます。また、カテゴリーまたはサブカテゴリーが選択されているが、サブカテゴリーの 1 つが選択されていない場合、チェックボックスの背景が青で表示され、選択されたカテゴリーの下位の選択が一部であることを示します。

テーブル内のカテゴリーを右クリックすると、次のようなコンテキスト・メニューが表示されます。

- 選択項目をチェック: テーブルで選択した行のすべてのチェック・ボックスをオンにします。
- 選択項目をチェック解除: テーブルで選択した行のすべてのチェック・ボックスをオフにします。
- すべてチェック: テーブルのすべてのチェック・ボックスをオンにします。これにより、すべてのカテゴリーが最終的な出力に使用されます。ツールバーの対応するチェックボックス・アイコンを使用することもできます。
- すべてチェック解除: テーブルのすべてのチェック・ボックスをオフにします。カテゴリーのチェックを解除すると、カテゴリーは最終的な出力で使用されなくなります。ツールバーの対応する空白のチェックボックス・アイコンを使用することもできます。

記述子テーブル内のセルを右クリックすると、次のようなコンテキスト・メニューが表示されます。

- コピー: 選択したコンセプトがクリップボードにコピーされます。
- フィールドもコピー: 選択した記述子が列の見出しと共にクリップボードにコピーされます。
- すべて選択: テーブル内のすべての行が選択されます。

カテゴリー モデル ナゲット: 「設定」 タブ

必要に応じて、「設定」タブを使用して、新しい入力データのテキスト・フィールド値を定義します。また、ここで出力のデータ・モデル (スコアリング・モード) も定義します。

注: このタブは、モデル ナゲットがストリーム内の領域にある場合にのみ、ノードのダイアログ・ボックスに表示されます。「モデル」パレットでこのナゲットを使用している場合、このタブは存在しません。

スコアリング・モード: フィールドとしてのカテゴリー

このオプションの場合、入力内と同じ数の出力レコードがあります。しかし、これによって各レコードは、モデル・タブで選択された全てのカテゴリー (チェック・マークを使用) 毎に新しいフィールドを含みます。各フィールドに、はい/いいえ、True/False、T/F または 1 および 2 などの **True** および **False** のフラグ値を入力します。例えば、フラグに数値を入力すると、自動的に整数値として処理されます。フラグ型では、文字列、整数、実数、または日付/時間のストレージ・タイプを利用することができます。

注: 非常に大規模なデータ・セットを使用している場合 (Db2 データベースの使用など)、「フィールドとしてのカテゴリー」を使用すると、データ量が原因で処理上の問題が発生する可能性があります。この場合は、代わりに「レコードとしてのカテゴリー」を使用することをお勧めします。

フィールド名拡張子: フィールド名の拡張接頭辞/接尾辞を指定したり、カテゴリー・コードを使用したりできます。カテゴリー名に加えてこの拡張子を使用して、フィールド名が生成されます。

- 追加方法: 拡張子をフィールド名に追加する場所を指定します。「接頭辞」を選択すると、文字列の頭に拡張子が追加されます。「接尾辞」を選択すると、文字列の終わりに拡張子が追加されます。

サブカテゴリーが選択されていない場合: スコアリングに選択されていないサブカテゴリーに含まれる記述子の処理方法を指定できます。2 つのオプションがあります。

- オプション 「記述子をスコアリングから完全に除外する」を選択すると、チェック記号のない (選択されていない) サブカテゴリーは無視され、スコアリングに使用されません。

- オプション 「記述子を上位カテゴリ内の記述子と合計する」 を選択すると、チェック記号のない (選択されていない) サブカテゴリの記述子は上位カテゴリ (このサブカテゴリの上位にあるカテゴリ) の記述子として使用されます。複数レベルのサブカテゴリが選択されない場合、記述子は使用できる最初の上位カテゴリにロール・アップされます。

句読点エラーを調整: 抽出時に句読点エラー (不適切な使用方法など) を含むテキストを一時的に正規化し、コンセプトの抽出可能性を向上させます。自由記述式アンケートの回答、電子メール、CRM データなど、テキストが短く品質が悪い場合、またはテキストに略語が多く含まれている場合に特に役立ちます。

スコアリング・モード:レコードとしてのカテゴリ

このオプションでは、新しいレコードが各 category、document ペアに作成されます。通常、入力に比べて出力に多くのレコードがあります。入力フィールドのほか、モデルの種類によって、新しいフィールドもデータに追加されます。

表 6. 「レコードとしてのカテゴリ」の出力フィールド:

新しい出力フィールド	説明
カテゴリ	テキスト・ドキュメントが割り当てられるカテゴリ名が指定されます。カテゴリが別のカテゴリのサブカテゴリである場合、ダイアログで選択した値によってカテゴリ名への完全パスが制御されます。

階層カテゴリの値: サブカテゴリの名前を出力内でどのように表示するかを制御します。

- 完全カテゴリ・パス。 カテゴリ名と、該当する場合、カテゴリ名とサブカテゴリ名をスラッシュを使用して区切り、上位カテゴリの完全パスを出力します。
- 省略したカテゴリ パス: カテゴリ名のみを出力します。ただし、省略記号を使用して、該当するカテゴリの上位カテゴリ数を示します。
- 下位レベルのカテゴリ: 完全パスまたは上位カテゴリを表示せず、カテゴリ名のみを出力します。

サブカテゴリが選択されていない場合: スコアリングに選択されていないサブカテゴリに含まれる記述子の処理方法を指定できます。2つのオプションがあります。

- オプション 「記述子をスコアリングから完全に除外する」 を選択すると、チェック記号のない (選択されていない) サブカテゴリは無視され、スコアリングに使用されません。
- オプション 「記述子を上位カテゴリ内の記述子と合計する」 を選択すると、チェック記号のない (選択されていない) サブカテゴリの記述子は上位カテゴリ (このサブカテゴリの上位にあるカテゴリ) の記述子として使用されます。複数レベルのサブカテゴリが選択されない場合、記述子は使用できる最初の上位カテゴリにロール・アップされます。

句読点エラーを調整: 抽出時に句読点エラー (不適切な使用方法など) を含むテキストを一時的に正規化し、コンセプトの抽出可能性を向上させます。自由記述式アンケートの回答、電子メール、CRM データなど、テキストが短く品質が悪い場合、またはテキストに略語が多く含まれている場合に特に役立ちます。

カテゴリ モデル ナゲット:その他のタブ

カテゴリ モデル ナゲットの「フィールド」タブと「設定」タブは、コンセプト モデル ナゲットと同じです。

- 「フィールド」タブ。詳しくは、35 ページの『コンセプト モデル:「フィールド」タブ』のトピックを参照してください。

- 「要約」タブ。詳しくは、36 ページの『コンセプト・モデル:「要約」タブ』のトピックを参照してください。

ストリームでのカテゴリ モデル ナゲットの使用

テキスト マイニング カテゴリ モデル ナゲットは、インタラクティブ・ワークベンチ・セッションから生成されます。このモデル ナゲットはストリームで使用できます。

例:カテゴリ モデル ナゲットを含む **Statistics** ファイル ノード

次の例では、テキスト マイニング モデル ナゲットの使用方法について示しています。

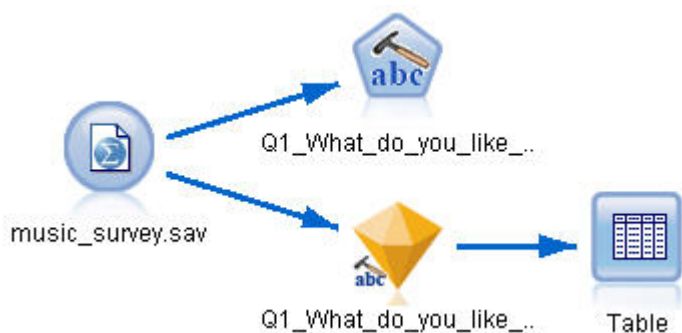


図 9. ストリームの例:テキスト マイニング カテゴリ モデル ナゲットを含む *Statistics* ファイル ノード

1. **Statistics** ファイル・ノード (「データ」タブ)。まず、このノードをストリームに追加して、テキスト・ドキュメントが保存されている場所を指定しました。



図 10. *Statistics* ファイル・ノード・ダイアログ・ボックス:「データ」タブ

2. テキスト マイニング カテゴリ モデル ナゲット (「モデル」タブ): 次に、カテゴリ モデル ナゲットを *Statistics* ファイル ノードに追加して接続しました。データのスコアリングに使用したいカテ

ゴリーを選択しました。

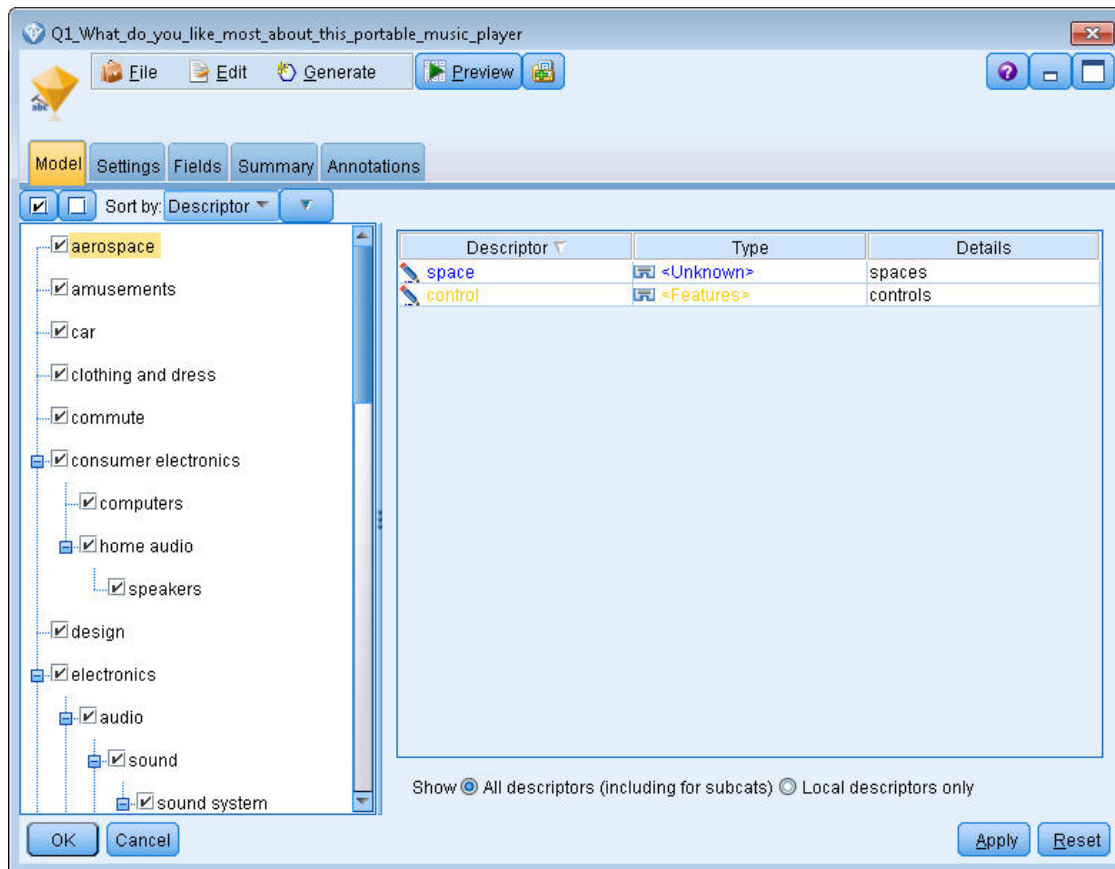


図 11. 「テキスト マイニング モデル ナゲット」 ダイアログ・ボックス: 「モデル」 タブ

3. テキスト マイニング モデル ナゲット (「設定」 タブ): 次に出力形式 「フィールドとしてのカテゴリー」 を定義しました。

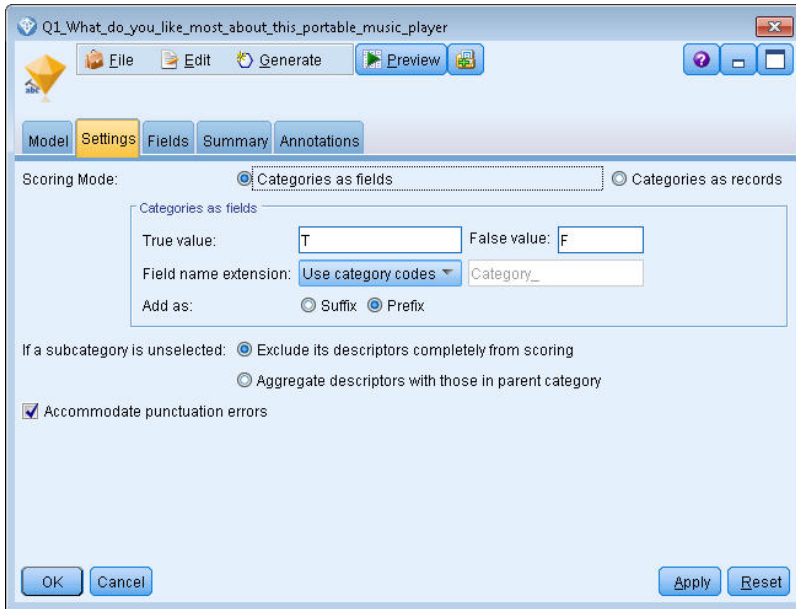


図 12. 「カテゴリ モデル ナゲット」 ダイアログ・ボックス: 「設定」 タブ

4. テキスト マイニング カテゴリ モデル ナゲット (「フィールド」 タブ): 次に Statistics ファイル ノードから作成されたフィールド名であるテキスト・フィールド変数を選択し、テキスト・フィールド が示す内容のオプション 「実際のテキスト」 やその他の設定を選択しました。

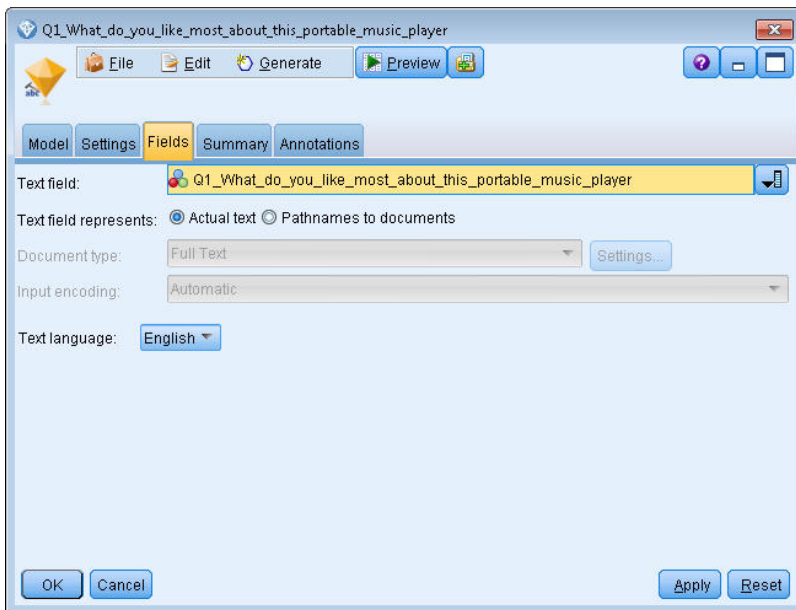


図 13. 「テキスト マイニング モデル ナゲット」 ダイアログ・ボックス: 「フィールド」 タブ

5. テーブル・ノード: 次に、 テーブル・ノードを接続して結果を表示し、ストリームを実行しました。

Table (3 fields, 844 records)

	ID	Q1_What_do_you_like_most_about_this_portable_music_player	Category
1	1	little, light	light
2	2	The battery power is great.	light
3	2	The battery power is great.	electronics/battery
4	2	The battery power is great.	electronics
5	3	cost and size	size
6	6	Battery life. Portability. Accessories. Style.	light
7	6	Battery life. Portability. Accessories. Style.	electronics/battery
8	6	Battery life. Portability. Accessories. Style.	electronics
9	7	I like its ability to store all of my music. I also like the ability to create playlists.	playlists
10	7	I like its ability to store all of my music. I also like the ability to create playlists.	light
11	7	I like its ability to store all of my music. I also like the ability to create playlists.	music
12	8	portability, capacity, sound quality, durability	light
13	8	portability, capacity, sound quality, durability	electronics/audio/sound
14	8	portability, capacity, sound quality, durability	electronics/audio

図 14. テーブル出力

第 4 章 テキスト・リンクのマイニング

テキスト リンク分析ノード

テキスト リンク分析 (TLA) ノードは、パターンマッチ・テクノロジーをテキスト マイニングのコンセプト抽出に追加し、既知のパターンに基づいてテキスト・データのコンセプト間の関連性を特定します。これらの関連性は、顧客が製品についてどのように感じているか、どの企業と組んでビジネスを行うか、または遺伝子または医薬品の間の関連性について説明が可能です。

例えば、競合他社の製品名を抽出しても、重要でない場合があります。このノードを使用して、データ内に人々がこの製品に関してどのように感じているかという意見がある場合、それについて学習することができます。関連性および関係性は、既知のパターンをテキスト・データに合致させることによって、特定および抽出します。

IBM SPSS Modeler Text Analytics に付属する特定のリソース・テンプレート内の TLA パターンを使用または独自のパターンを作成/編集できます。パターン規則は、マクロ、単語リスト、およびブール型質問を形成する単語の空所、または入力テキストと比較される条件規則で構成されています。TLA パターン規則がテキストに一致する場合、このテキストを TLA 結果として抽出し、出力データとして再構築できます。詳しくは、215 ページの『第 18 章 テキスト リンク規則について』のトピックを参照してください。

テキスト リンク分析ノードを使用して、より直接的にテキストから TLA パターン結果を特定および抽出し、パターンの結果をストリーム内のデータセットに追加できます。ただし、テキスト リンク分析を実行できるのは、テキスト リンク分析ノードだけではありません。テキスト マイニングモデル作成ノードのインタラクティブ・ワークベンチ・セッションを使用することもできます。

インタラクティブ・ワークベンチで、TLA パターン結果を検証してその結果をカテゴリ記述子として使用し、ドリルダウンおよびグラフを使用して結果についてより詳細に学習することができます。詳しくは、149 ページの『第 11 章 テキスト リンク分析の検証』のトピックを参照してください。実際に、テキスト マイニング・ノードを使用して TLA 結果を抽出すると、後で TLA ノードで直接使用するためにテンプレートを検証し、データ向けに調整することができます。

出力は、最大 6 つのスロット、または 6 つの部分で表示されます。詳しくは、51 ページの『TLA ノード出力』のトピックを参照してください。

このノードは、IBM SPSS Modeler ウィンドウの下部にあるノード・パレットの IBM SPSS Modeler Text Analytics タブにあります。詳しくは、8 ページの『IBM SPSS Modeler Text Analytics ノード』のトピックを参照してください。

要件: テキスト リンク分析ノードは、標準的な入力ノード (データベース・ノード、フラット・ファイル・ノードなど) のいずれかを使用してフィールドに読み込まれた、またはファイル リスト ノードまたは Web フィールド ノードによって生成された外部ドキュメントへのフィールド・リスト・パスに読み込まれたテキスト・データを受け入れます。

利点: テキスト リンク分析 ノードは、基本的なコンセプト抽出以上の機能により、コンセプト間の関連性、およびデータ内にあると考えられる関連する意見や識別子に関する情報を提供します。

「テキスト リンク分析ノード」: 「フィールド」タブ

「フィールド」タブを使用して、コンセプトを抽出するデータのフィールド設定を指定します。設定できるパラメーターを次に示します。

ID フィールド: テキスト・レコードの識別子を含むフィールドを選択します。識別子は整数でなければなりません。ID フィールドは、各テキスト・レコードのインデックスとして機能します。テキスト・フィールドがマイニングされるテキストを示す場合、ID フィールドを使用します。

テキスト フィールド: マイニングするテキストが含まれているフィールドを選択します。このフィールドはデータ・ソースによって異なります。

言語フィールド (Language field): 2 文字の ISO 言語 ID を含むフィールドを選択します。フィールドを選択しなかった場合、各ドキュメントの言語は、指定されたテンプレートの言語であると想定されます。

ドキュメント タイプ: ドキュメント・タイプは、テキストの構造を指定します。次に示すタイプの 1 つを選択します。

- **フル テキスト:** このオプションは、多くのドキュメントまたはテキスト・ソースに使用します。テキストのセット全体をスキャンして抽出します。他のオプションとは異なり、このオプションに追加設定はありません。
- **構造のあるテキスト:** このオプションは、参考文献形式、特許、特定および分析できる通常の構造を含むファイルに使用します。このドキュメントタイプを使用して、抽出プロセスのすべてまたは一部をスキップします。キーワードの区切り文字の定義、タイプの割り当て、および出現頻度の最小値の指定ができます。このオプションを選択する場合、「設定」ボタンをクリックして、「ドキュメント設定」ダイアログ ボックスの「構造のあるテキストの書式設定」領域にテキストの区切り文字を入力します。詳しくは、22 ページの『「フィールド」タブのドキュメント設定』のトピックを参照してください。

テキストの単位: 次の実行モードを選択します。

- **ドキュメント・モード:** 通信社からの記事など、短く意味的に同質のドキュメントに使用します。

- **パラグラフ モード:** Web ページおよびタグのないドキュメントに使用します。抽出プロセスでは、内部タグやシンタックスなどの特徴を利用して、ドキュメントを意味的に分割します。このモードを選択すると、パラグラフごとにスコアリングが適用されます。そのため、例えば、apple および orange が同じパラグラフで見つかった場合にのみ、規則「apple & orange」が当てはまります。

注: PDF ドキュメントからテキストを抽出する方法が原因で、これらのドキュメントでは「パラグラフモード」は機能しません。これは、抽出により復帰マーカーが抑制されるためです。

パラグラフ モードの設定: このオプションは、「パラグラフ モード」に「テキストの単位」オプションを設定した場合にのみ選択できます。抽出で使用する文字のしきい値を指定します。実際のサイズは、最も近いピリオドに丸められます。ドキュメント・コレクションのテキストから作成される単語の関連性を典型とするには、抽出サイズが小さすぎないように指定します。

- **最小:** 抽出で使用する文字の最小数を指定します。
- **最大値:** 抽出で使用する文字の最大数を指定します。

リソースのコピー元: テキスト マイニング時、抽出は、「エキスパート」タブの設定だけでなく、言語リソースに基づいて行われます。これらのリソースは、抽出時のテキストの処理方法の基本として機能し、コンセプト、タイプ、および TLA パターンを取得します。リソース・テンプレートからリソースをテキスト マイニングモデル作成ノードにコピーできます。

リソース・テンプレートは、特定のドメインまたは使用向けに調整された、事前定義済みライブラリーおよび詳細な言語リソースおよび非言語リソースです。これらのリソースは、抽出時のデータの処理方法についての基本として機能します。「読み込み」をクリックし、リソースをコピーするテンプレートを選択します。

テンプレートを選択したときにテンプレートが読み込まれ、ストリームが実行されているときには読み込まれません。読み込んでいるときに、リソースのコピーがノードに保存されます。そのため、更新されたテンプレートを使用したい場合、ここで再読み込みを行う必要があります。詳しくは、27 ページの『テンプレートおよび TAP からのリソースのコピー』のトピックを参照してください。

テキストの言語。 マイニングされるテキストの言語を識別します。ノードでコピーされたリソースが、表示される言語オプションを制御します。リソースを調整した言語を選択してください。

テキスト リンク分析ノード:「エキスパート」タブ

このノードでは、テキスト リンク分析 (TLA) パターン結果の抽出が自動的に有効化されています。「エキスパート」タブには、テキストの抽出方法および処理方法に影響を与える追加パラメーターがあります。このダイアログ・ボックスのパラメーターは、抽出プロセスの基本的な操作、そしていくつかの高度な操作を制御します。また、抽出結果にも影響を与える言語リソースやオプションも数多くあり、選択するリソース・テンプレートによって制御します。

グローバル頻度が次の値以上のコンセプトに抽出を制限: 抽出するために、単語または句が出現する必要がある最低限の回数を指定します。値に 5 を指定すると、抽出するこれらの単語または句が、レコードまたはドキュメントのセット全体で少なくとも 5 回出現するよう、制限します。

この制約を変更すると、抽出結果、つまり作成されるカテゴリーに大きな違いが生じる場合があります。あるレストランのデータを処理し、このオプションの制約に1より大きい値を設定しないものとします。この場合、抽出結果がピザ (1)、薄いピザ (2)、ほうれん草のピザ (2)、および好きなピザ (2) となります。ただし、抽出のグローバル出現頻度を 5 以上に設定して抽出すると、これらのコンセプトのうち 3 つが取得されなくなります。代わりに、ピザが最も簡単な形で、この単語は考えられる候補として既に存在するため、ピザ (7) が取得されます。また、残りのテキストにピザという単語を含む他の句があるかどうかによ

て、7より大きい出現頻度がある場合があります。また、ほうれん草のピザがカテゴリーの記述子である場合、すべてのレコードをキャプチャーする代わりに、記述子としてピザの追加が必要な場合があります。このため、カテゴリーが既に作成されている場合は、注意してこの制約を変更してください。

これは抽出のみの機能であることに注意してください。つまり、テンプレートに用語が含まれる場合 (通常そのようになります) でテンプレートの用語がテキスト内で見つかった場合、その用語は頻度に関わらずインデックス付けされます。

例えば、コア・ライブラリーの <Location> タイプに「ロサンゼルス」が含まれている基本リソース・テンプレートを使用するとします。この場合、ドキュメント内での「ロサンゼルス」の出現回数が 1 回だけでも、ロサンゼルスがコンセプト・リストに含まれることとなります。これを回避するには、「グローバル頻度が次の値以上のコンセプトに抽出を制限」フィールドに入力された値以上の出現回数を持つコンセプトだけを表示するように、フィルターを設定する必要があります。

句読点エラーを調整: 抽出時に句読点エラー (不適切な使用方法など) を含むテキストを一時的に正規化し、コンセプトの抽出可能性を向上させます。自由記述式アンケートの回答、電子メール、CRM データなど、テキストが短く品質が悪い場合、またはテキストに略語が多く含まれている場合に特に役立ちます。

文字数が次の最小値以上のときにスペルを調整する: Fuzzy Grouping の手法を適用し、共通してミススペルのある単語またはスペルの近い単語を 1 つのコンセプトにグループ化できるようにします。Fuzzy Grouping アルゴリズムでは、最初の母音を除くすべての母音を一時的に抜き取った後抽出した単語から 2 つ/3 つの子音を抜き取り、それらを比較して、それらが同じで modeling と modelling が同じグループに分けられるかどうかを確認します。ただし、各キーワードが <Unknown> タイプを除いて、別のタイプに割り当てられた場合、Fuzzy Grouping 手法は適用されません。

Fuzzy Grouping を使用する前に必要な、語幹文字数の制限を定義することもできます。キーワード内の語幹文字数は、すべての文字を合計し、活用語尾、複合語キーワードの場合は区切り文字および前置詞を形成する文字を差し引いて計算します。例えば、キーワード exercises の語幹文字数は「exercise」という形式で 8 文字と数えられます。語末の s は活用語尾 (複数形) であるためです。同様に、apple sauce の語幹文字は 10 文字 (「apple sauce」)、そして manufacturing of cars の語幹文字は 16 文字 (「manufacturing car」) となります。この算出方法は、Fuzzy Grouping を適用するべきかどうかを確認するためにのみ使用されますが、単語がどのように一致するかについては影響を与えません。

注: 特定の単語が後で不適切にグループ化されていることが分かった場合、「拡張リソース」タブの **Fuzzy Grouping**: 例外 セクションで 明示的に宣言することによって、単語のペアをこの手法から除外できます。詳しくは、205 ページの『Fuzzy Grouping』のトピックを参照してください。

ユニタームを抽出 単語が複合語の一部でない限り、または名詞、またはスピーチ内の認識できない品詞である場合、このオプションは単一の単語 (ユニターム) を抽出します。

固有表現を抽出 電話番号、セキュリティ番号、時間、日付、通貨、数字、パーセント、電子メールアドレス、HTTP アドレスなどの固有表現を抽出します。「拡張リソース」タブの「固有表現: 設定」セクションで、特定の種類の固有表現を追加したり除外したりできます。不要な固有表現を無効にすることにより、抽出エンジンは処理時間を節約できます。詳しくは、210 ページの『構成』のトピックを参照してください。

大文字アルゴリズム キーワードの最初の文字が大文字である場合、組み込み辞書にない単純キーワードおよび複合キーワードを抽出します。このオプションには、最も適切な名詞を抽出するのに優れた方法があります。

可能な場合は、個人名の一部または全部をグループ化 テキスト内で別々の形式で同時に出現する名前をグループ化します。名前はテキストの始めでは完全な形式で、後は短い形式でのみ参照されるため、この機能が役立ちます。このオプションでは、タイプが <Unknown> のユニタームが、タイプ <Person> の複合キーワードの最後の単語に一致するようにします。例えば、*doe* があり、最初タイプが <Unknown> である場合、抽出エンジンは、<Person> タイプの複合キーワードに最後の単語として *doe* が含まれているかどうか (例: *john doe*) を確認します。ほとんどがユニタームとして抽出されることがないため、人の名前に適用されることはありません。

機能語による倒置を次の値を最大値として考慮する 倒置手法を適用する場合に指定されている場合がある非機能的単語の最大数を指定します。この倒置手法では、活用語尾に関係なく、含まれる非機能的単語 (of や the など) によってお互いに異なる類似した句をグループ化します。例えば、この値を最大 2 単語に設定し、*company officials* および *officials of the company* が抽出されたとします。この場合、両方の抽出キーワードは、of the が無視されると同じであるとみなされるため、最終コンセプト・リストに共にグループ化されます。

マルチタームをグループ化するとき派生関係を使用: ビッグデータを処理するときこのオプションを選択すると、派生規則を使用してマルチタームがグループ化されます。

TLA ノード出力

テキスト リンク分析ノードを実行した後、データが再構築されます。テキスト マイニングでデータを再構築する方法を理解することは重要です。データ・マイニングに異なる構造が必要な場合、「フィールド操作」パレットのノードを使用して、これを実行できます。例えば、各行にテキスト・レコードが表示されているデータを処理している場合、ソース・テキスト・データで見つかったパターンごとに 1 行ずつ作成されます。出力の各行について、次の 15 個のフィールドがあります。

- 6 つのフィールド (コンセプト 1、コンセプト 2、およびコンセプト 6 のようなコンセプト #) は、パターン・マッチで見つかったコンセプトを示します。
- 6 つのフィールド (タイプ 1、タイプ 2、およびタイプ 6 のようなタイプ #) は、各コンセプトのタイプを示します。
- 条件規則名は、テキストを合致させ、出力を生成するのに使用するテキスト リンク規則の名前を示します。
- ノードで指定した ID フィールドの名前を使用し、入力データと同じようにレコード ID またはドキュメント ID を示すフィールド。
- 「マッチ テキスト」は、元のレコードまたはドキュメント内にある、TLA パターンに合致したテキスト・データの部分を示します。

注: リリース 5.0 より前のテキスト リンク分析ノードを含む既存のストリームは、ノードを更新しない限り、完全には実行できない可能性があります。IBM SPSS Modelerの最新バージョンでの特定の機能改善には、古いノードを新しいバージョンに置き換える必要があります。これにより、さらに展開可能で強力となります。

特定の言語の自動翻訳も実行できます。この機能を使用して、話せないまたは読めない言語のドキュメントをマイニングできます。翻訳機能を使用したい場合は、SDL Software as a Service (SaaS)へのアクセスが必要です。詳しくは、翻訳設定のトピックを参照してください。

TLA 結果のキャッシュ

キャッシュすると、テキスト リンク分析結果がストリーム内に置かれます。ストリームの実行ごとにテキスト リンク分析結果の抽出が繰り返されないようにするには、テキスト リンク分析ノードを選択して、メ

ニューから「編集」>「ノード」>「キャッシュ」>「有効化」を選択します。次回ストリームが実行される場合、出力がノードにキャッシュされます。ノードのアイコンには、小さい「ドキュメント」のグラフィックが表示され、キャッシュがいっぱいになると白から緑に変わります。キャッシュはセッションの間保存されます。ストリームを閉じて再び開いた後など、キャッシュを別の日にも保持するには、ノードを選択して、メニューから「編集」>「ノード」>「キャッシュ」>「キャッシュの保存」を選択します。次にストリームを開く場合、翻訳を再度行わずに保存されたキャッシュを再度読み込むことができます。

また、ノードを右クリックして、コンテキスト・メニューから「キャッシュ」を選択して、ノードのキャッシュを保存したり有効にしたりできます。

ストリーム内のテキスト リンク分析ノードの使用

テキスト リンク分析ノードを使用して、データにアクセスし、ストリーム内のコンセプトを抽出します。入力ノードを使用して、データにアクセスできます。

例:Statistics ファイル・ノードとテキスト リンク分析ノード

次の例には、テキスト リンク分析ノードの使用方法が示されています。



図 15. 例:Statistics ファイル・ノードとテキスト リンク分析ノード

1. **Statistics** ファイル・ノード (「データ」タブ)。 まず、このノードをストリームに追加して、テキストが保存されている場所を指定しました。

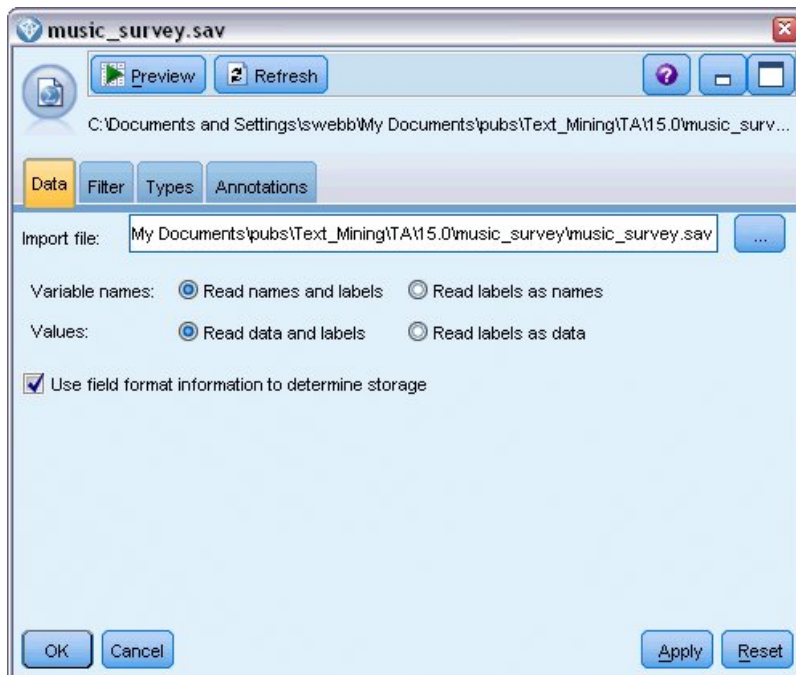


図 16. Statistics ファイル・ノード・ダイアログ・ボックス:「データ」タブ

2. テキスト リンク分析ノード (「フィールド」タブ): このノードをストリームに接続して、コンセプトを抽出し、下流でモデル作成または表示しました。ID フィールドおよびデータを含むテキスト・フィールド名、そしてその他の設定を指定しました。

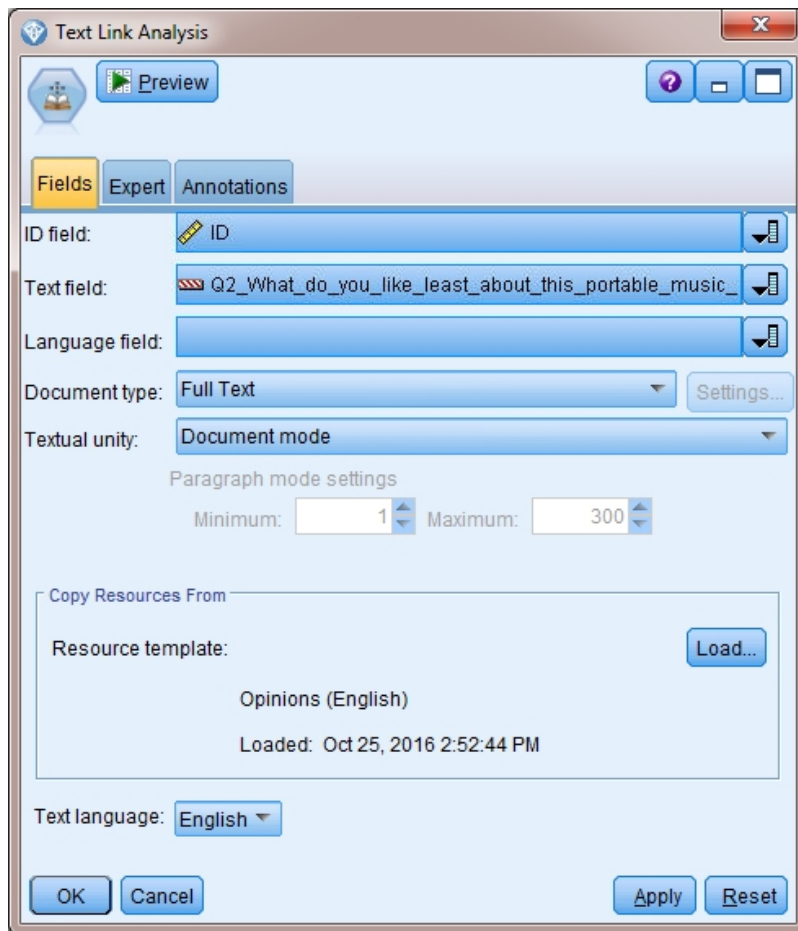


図 17. 「テキスト リンク分析ノード」ダイアログ・ボックス:「フィールド」タブ

3. テーブル・ノード: 最後にテーブル・ノードを接続して、テキスト・ドキュメントから抽出されたコンセプトを表示しました。表示されるテーブル出力で、このストリームがテキスト リンク分析ノードで実行された後、データ内の TLA パターン結果を確認できます。いくつかの結果で、合致したコンセプト/タイプが 1 つだけであることを示します。他の結果はより複雑で、いくつかのタイプおよびコンセプトが含まれています。また、テキスト リンク分析ノードを使用してデータを実行し、コンセプトを抽出した結果、データのいくつかの部分が変化しています。例の元のデータには、8 つのフィールドと 405 件のレコードが含まれていました。テキスト リンク分析ノードを実行した後、フィールド数は 15 で、レコード数は 640 件となります。TLA パターン結果ごとに 1 行ずつ割り当てられます。例えば、ID 7 は、3 つの TLA パターン結果が抽出されているため、3 行となります。この出力を元のデータに結合したい場合、結合ノードを使用できます。

	Concept1	Type1	Concept2	Type2	Conc...	Type3	Con...	Type4	Conc...	Type5	Con...	Type6	Rule Number	ID	Matched Text
1	expensive	NegativeBudget	Null	Null	Null	Null	Null	Null	Null	Null	Null	Null	0/0350_opinion	1	<!*expensive!*
2	screen	Unknown	difficult	Nega...	Null	Null	Null	Null	Null	Null	Null	Null	0/0145_topic + opinion	2	The <!*screen!*
3	software	Unknown	difficult	Nega...	Null	Null	Null	Null	Null	Null	Null	Null	0/0211_opinion + topic	3	<!*difficult!*
4	nothing	Uncertain	Null	Null	Null	Null	Null	Null	Null	Null	Null	Null	0/0153_topic/opinion	4	<!*Nothing!*
5	like	Positive	Null	Null	Null	Null	Null	Null	Null	Null	Null	Null	0/0350_opinion	4	Nothing , <!*I love it!*
6	battery life	Unknown	too long	Nega...	Null	Null	Null	Null	Null	Null	Null	Null	0/0145_topic + opinion	5	<!*Battery life!*
7	ubiquitousness	Unknown	Null	Null	Null	Null	Null	Null	Null	Null	Null	Null	0/0500_topic	6	<!*Ubiquitousness!*
8	40gb model	Unknown	available	Positi...	Null	Null	Null	Null	Null	Null	Null	Null	0/0145_topic + opinion	7	I wish the <!*40GB model!*
9	20gb model	Unknown	Null	Null	Null	Null	Null	Null	Null	Null	Null	Null	0/0102_topic + Negative + topic	7	I have a <!*20GB model!*
10	memory	Unknown	need more	Nega...	Null	Null	Null	Null	Null	Null	Null	Null	0/0102_topic + Negative + topic	7	I have a <!*20GB model!*

図 18. テーブル出力ノード

第 5 章 外部ソース・テキストの参照

ファイル ビューアー ノード

ドキュメントのコレクションをマイニングしている場合、ファイルの完全パス名をテキスト マイニング モデル作成ノードに直接指定できます。ただし、テーブル・ノードに出力する場合、テーブル・ノード内のテキストではなく、ドキュメントの完全パス名のみ表示されます。ファイル・ビューアー・ノードをテーブル・ノードのアナログとして使用でき、すべてのドキュメントを 1 つのファイルに結合することなくドキュメント内の実際のテキストを使用できるようになります。

ストリームではアクセスできないため、ファイル・ビューアー・ノードを使用して、コンセプトが抽出されたソース・テキストまたは翻訳されていないテキストにアクセスを提供することによって、テキスト抽出の結果をより深く理解することができます。このノードは、ファイル リスト ノードの後のストリームに追加され、すべてのファイルへのリンクのリストを取得します。

このノードの結果として、コンセプトを抽出するために読み込み、使用されたすべてのドキュメント要素を示すウィンドウが表示されます。このウィンドウで、ツールバー・アイコンをクリックして、ドキュメント名をハイパーリンクで表示する外部ブラウザでレポートを起動することができます。リンクをクリックして、コレクションの該当するドキュメントを開くことができます。詳しくは、56 ページの『ファイル ビューアー ノードの使用』のトピックを参照してください。

このノードは、IBM SPSS Modeler ウィンドウの下部にあるノード・パレットの IBM SPSS Modeler Text Analytics タブにあります。詳しくは、8 ページの『IBM SPSS Modeler Text Analytics ノード』のトピックを参照してください。

注: クライアントサーバー・モードで作業し、ファイル・ビューアー・ノードがストリームの一部である場合、ドキュメント・コレクションはサーバーの Web サーバー・ディレクトリーに保存する必要があります。テキスト マイニング出力ノードは、Web サーバー・ディレクトリー保存されているドキュメントのリストを作成するため、Web サーバーのセキュリティ設定により、これらのドキュメントに対する権限を管理します。

ファイル ビューアー ノード設定

ファイル・ビューアー・ノード設定には、以下のオプションを指定できます。

ドキュメント・フィールド: 表示するドキュメントの名前全体およびパスを含むデータからフィールドを選択します。

生成された **HTML** ページのタイトル: ドキュメントのリストを表示するページの最上部に表示されるタイトルを作成します。

ファイル ビューアー ノードの使用

次の例には、ファイル ビューアー ノードの使用方法が示されています。

例:ファイル リスト ノードおよびファイル・ビューアー・ノード



図 19. ファイル・ビューアー・ノードの使用を示すストリーム

1. ファイル リスト ノード (「設定」タブ)。最初にこのノードを追加して、ドキュメントの場所を指定しました。

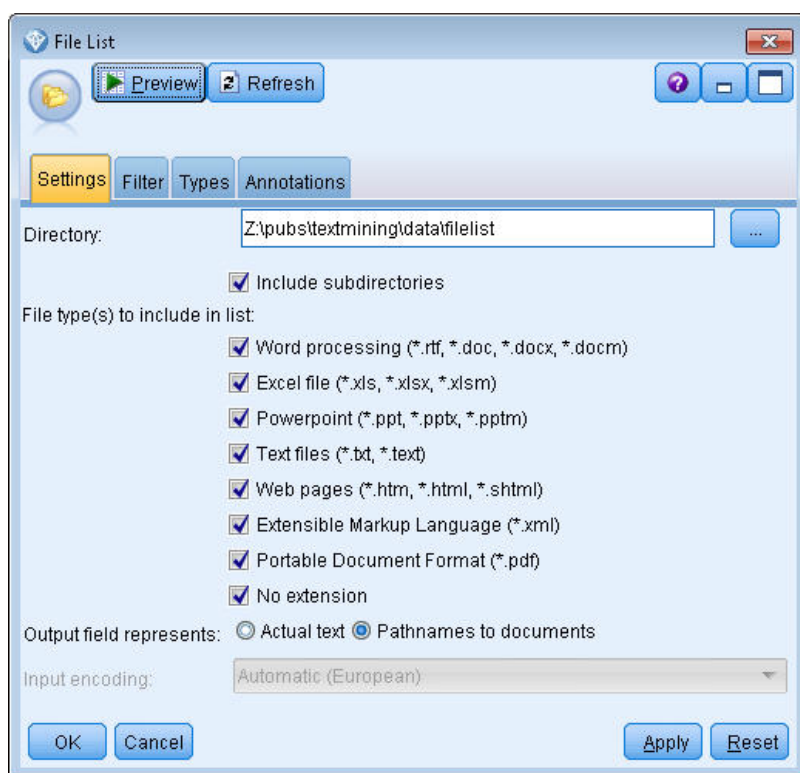


図 20. 「ファイル リスト ノード」 ダイアログ・ボックス: 「設定」タブ

2. ファイル・ビューアー・ノード (「設定」タブ): 次に、ファイル・ビューアー・ノードを接続して、ドキュメントの HTML リストを作成しました。

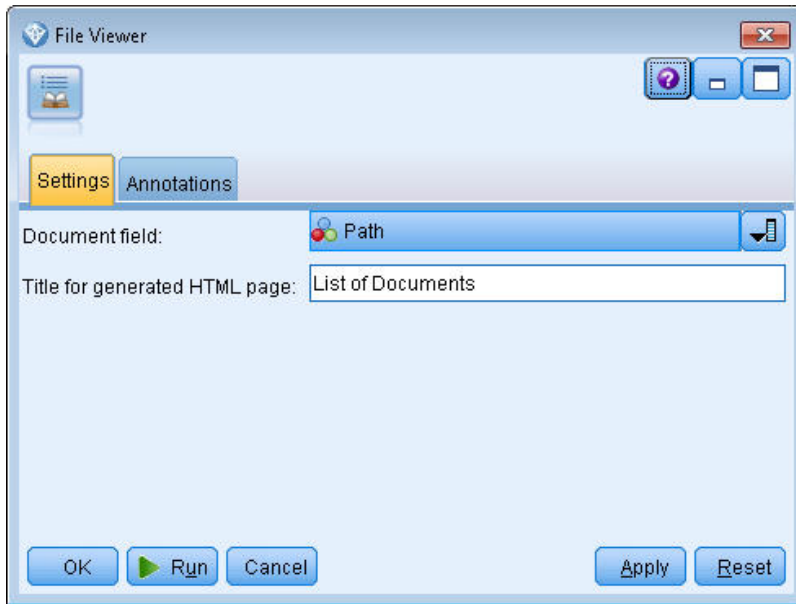


図 21. 「ファイル・ビューアー・ノード」ダイアログ・ボックス:「設定」タブ

3. 「ファイル・ビューアー出力」ダイアログ: 次に、新しいウィンドウでドキュメントのリストを出力するストリームを実行しました。

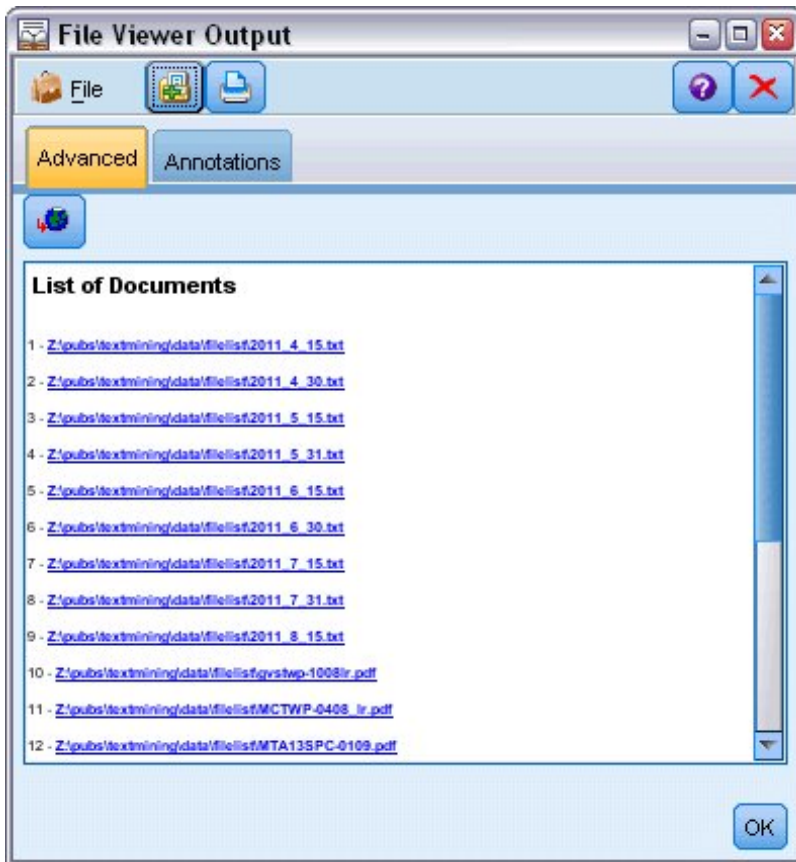


図 22. ファイル・ビューアー出力

4. ドキュメントを参照するために、赤い矢印の付いた地球儀の表示されたツールバーをクリックしました。ブラウザーのドキュメントのハイパーリンクのリストが開きました。

第 6 章 スクリプト用のノードのプロパティ

IBM SPSS Modeler には、コマンド・ラインからストリームを実行できるスクリプト言語があります。ここでは、IBM SPSS Modeler Text Analytics に付属する各ノードに固有のノードのプロパティについて説明します。IBM SPSS Modeler に付属するノードの標準セットの詳細については、『Scripting and Automation Guide』を参照してください。

ファイル リスト ノード:filelistnode

スクリプトには、次の表のプロパティを使用できます。このノードを `filelistnode` といいます。

表 7. ファイル リスト ノードのスクリプトのプロパティ

スクリプトのプロパティ	データ型
<code>path</code>	文字列
<code>recurse</code>	フラグ
<code>word_processing</code>	フラグ
<code>excel_file</code>	フラグ
<code>powerpoint_file</code>	フラグ
<code>text_file</code>	フラグ
<code>web_page</code>	フラグ
<code>xml_file</code>	フラグ
<code>pdf_file</code>	フラグ
<code>no_extension</code>	フラグ

注: 'Create list' パラメーターは使用できなくなり、そのオプションを含むスクリプトは自動的に 'Files' 出力に変換されます。

Web フィード ノード:webfeednode

スクリプトには、次の表のプロパティを使用できます。このノードを `webfeednode` といいます。

表 8. Web フィード ノードのスクリプトのプロパティ

スクリプトのプロパティ	データ型	プロパティの説明
<code>urls</code>	<i>string1 string2 ...stringn</i>	各 URL は、リスト構造で指定されます。「\n」で区切った URL リスト
<code>recent_entries</code>	フラグ	
<code>limit_entries</code>	整数	URL ごとに読み込む最新エントリー数
<code>use_previous</code>	フラグ	Web フィード・キャッシュを保存して再利用
<code>use_previous_label</code>	文字列	保存した Web キャッシュの名前
<code>start_record</code>	文字列	非 RSS の開始タグ

表 8. Web フィード ノードのスキプトのプロパティ (続き)

スキプトのプロパティ	データ型	プロパティの説明
url <i>n</i> .title	文字列	リスト内の各 URL について、ここでも定義する必要があります。最初の URL は url1.title となり、番号は URL リストの位置に対応します。コンテンツのタイトルを含む開始タグです。
url <i>n</i> .short_description	文字列	url の場合と同じ <i>n</i> .title.
url <i>n</i> .description	文字列	url の場合と同じ <i>n</i> .title.
url <i>n</i> .authors	文字列	url の場合と同じ <i>n</i> .title.
url <i>n</i> .contributors	文字列	url の場合と同じ <i>n</i> .title.
url <i>n</i> .published_date	文字列	url の場合と同じ <i>n</i> .title.
url <i>n</i> .modified_date	文字列	url の場合と同じ <i>n</i> .title.
html_alg	None HTMLCleaner	コンテンツのフィルタリング方法。
discard_lines	フラグ	短い行を破棄。 min_words とともに使用
min_words	整数	最小単語数。
discard_words	フラグ	短い行を破棄。 min_avg_len とともに使用
min_avg_len	整数	
discard_scw	フラグ	1 文字単語が多い行を破棄。 max_scw とともに使用
max_scw	整数	1 行につき最大 0 から 100 パーセントの割合の 1 文字単語
discard_tags	フラグ	特定のタグを含む行を破棄。
tags	文字列	特殊文字は、バックスラッシュ (\) でエスケープする必要があります。
discard_spec_words	フラグ	特定の文字列を含む行を破棄。
words	文字列	特殊文字は、バックスラッシュ (\) でエスケープする必要があります。

言語ノード: languageidentifier

スキプトには、次の表のプロパティを使用できます。このノードを languageidentifier といいます。

表 9. 言語ノードのスキプトのプロパティ

スキプトのプロパティ	データ型	プロパティの説明
テキスト	フィールド	
language_field_name	文字列	出力として生成するフィールド名。
unidentified_language_value	Undefined Supported Custom	言語を識別できない場合に使用するデフォルト値。

表 9. 言語ノードのスキプトのプロパティ (続き)

スキプトのプロパティ	データ型	プロパティの説明
unidentified_language_supported	en de es fr it ja nl pt	ISO コード。unidentified_language_value が Supported である場合にのみ使用できます。
unidentified_language_custom	文字列	unidentified_language_value が Custom である場合にのみ使用できます。

テキスト マイニング・ノード:TextMiningWorkbench

次のパラメーターを使用して、スキプトを介してノードを定義または更新することができます。このノードを TextMiningWorkbench といいます。

重要: スキプトを介して異なるリソース・テンプレートを指定できません。テンプレートが必要な場合は、ノードのダイアログ・ボックスで選択する必要があります。

表 10. テキスト マイニング・モデル作成ノードのスキプトのプロパティ

スキプトのプロパティ	データ型	プロパティの説明
テキスト	フィールド	
method	ReadText ReadPath	
docType	整数	正の数 (0,1,2) を指定します。ここでは 0 = Full Text、1 = Structured Text、2 = XML になります
encoding	Automatic (自動) "UTF-8" "UTF-16" "ISO-8859-1" "US-ASCII" "CP850" "EUC-JP" "SHIFT-JIS" "ISO2022-JP"	"UTF-8" のような特殊文字と値の組み合わせは、算術演算子と区別するために引用符 (") で囲む必要があります。
unity	整数	正の数 (0,1) を指定します。ここでは、0 = Paragraph および 1 = Document になります
para_min	整数	
para_max	整数	
mtag	文字列	すべての mtag 設定が含まれます (XML ファイルの「設定」ダイアログ・ボックス)
mclef	文字列	すべての mclef 設定が含まれます (構造テキスト・ファイルの「設定」ダイアログ・ボックス)
partition	フィールド	
custom_field	フラグ	分割フィールドを指定するかどうかを示します。

表 10. テキストマイニング・モデル作成ノードのスクリプトのプロパティ (続き)

スクリプトのプロパティ	データ型	プロパティの説明
use_model_name	フラグ	
model_name	文字列	
use_partitioned_data	フラグ	データ区分フィールドが定義されている場合、学習データだけがモデルの構築に使用されます。
model_output_type	Interactive Model	Interactive はカテゴリモデルになります。Model はコンセプトモデルになります。
use_interactive_info	フラグ	ワークベンチ・セッションでインタラクティブに作成する場合のみ。
reuse_extraction_results	フラグ	ワークベンチ・セッションでインタラクティブに作成する場合のみ。
interactive_view	Categories TLA Clusters	ワークベンチ・セッションでインタラクティブに作成する場合のみ。
extract_top	整数	このパラメーターは model_type = Concept のときに使用されます。
use_check_top	フラグ	
check_top	整数	
use_uncheck_top	フラグ	
uncheck_top	整数	
language	de en es fr it ja nl pt	
frequency_limit	整数	バージョン 14.0 では廃止。
concept_count_limit	整数	抽出を少なくともこの値以上のグローバル度数のコンセプトに制限。
fix_punctuation	フラグ	
fix_spelling	フラグ	
spelling_limit	整数	
extract_uniterm	フラグ	
extract_nonlinguistic	フラグ	
upper_case	フラグ	
group_names	フラグ	
permutation	整数	語順が異なる類義語で無視する非機能語の数 (デフォルトは 3)。

テキスト マイニング モデル ナゲット:TMWBModelApplier

スクリプトには、次の表のプロパティを使用できます。このノードを TMWBModelApplier といいます。

表 11. テキスト マイニング モデル ナゲットのプロパティ

スクリプトのプロパティ	データ型	プロパティの説明
scoring_mode	フィールド Records	
field_values	Flags カウント	このオプションは、カテゴリーのモデル ナゲットでは使用できません。Flags は、TRUE または FALSE にセットしてください
true_value	文字列	Flags で、true の値を定義。
false_value	文字列	Flags で、false の値を定義。
extension_concept	文字列	フィールド名の拡張子を指定します。コンセプト名に加えてこの拡張子を使用して、フィールド名が生成されます。add_as 値を使用して、この拡張子を使用する場所を指定します。
extension_category	文字列	フィールド名拡張子: フィールド名の拡張接頭辞/接尾辞を指定したり、カテゴリー・コードを使用したりできます。カテゴリー名に加えてこの拡張子を使用して、フィールド名が生成されます。add_as 値を使用して、この拡張子を使用する場所を指定します。
add_as	Suffix Prefix	
fix_punctuation	フラグ	
excluded_subcategories_descriptors	RollUpToParent 無視	<p>カテゴリー・モデル のみ。サブカテゴリーが選択されていない場合、スコアリングに選択されていないサブカテゴリーに含まれる記述子の処理方法を指定できます。2 つのオプションがあります。</p> <ul style="list-style-type: none"> • Ignore. オプション 「記述子をスコアリングから完全に除外する」 を選択すると、チェック記号のない (選択されていない) サブカテゴリーは無視され、スコアリングに使用されません。 • RollUpToParent. オプション 「記述子を上位カテゴリー内の記述子と合計する」 を選択すると、チェック記号のない (選択されていない) サブカテゴリーの記述子は上位カテゴリー (このサブカテゴリーの上位にあるカテゴリー) の記述子として使用されます。複数レベルのサブカテゴリーが選択されない場合、記述子は使用できる最初の上位カテゴリーにロール・アップされます。
check_model	フラグ	バージョン14 では廃止。
テキスト	フィールド	
method	ReadText ReadPath	

表 11. テキストマイニングモデルナゲットのプロパティ (続き)

スクリプトのプロパティ	データ型	プロパティの説明
docType	整数	正の数 (0,1,2) を指定します。ここでは 0 = Full Text、1 = Structured Text、2 = XML になります
encoding	Automatic (自動) "UTF-8" "UTF-16" "ISO-8859-1" "US-ASCII" "CP850" "EUC-JP" "SHIFT-JIS" "ISO2022-JP"	"UTF-8" のような特殊文字と値の組み合わせは、算術演算子と区別するために引用符 (") で囲む必要があります。
language	de en es fr it ja nl pt	

テキストリンク分析ノード: textlinkanalysis

次の表のパラメーターを使用して、スクリプトを介してノードを定義または更新することができます。このノードを textlinkanalysis といいます。

重要: スクリプトを介してリソース・テンプレートを指定できません。テンプレートを選択するには、ノードのダイアログ・ボックスで選択する必要があります。

表 12. テキストリンク分析 (TLA) ノードのスクリプトのプロパティ

スクリプトのプロパティ	データ型	プロパティの説明
id_field	フィールド	
テキスト	フィールド	
method	ReadText ReadPath	
docType	整数	正の数 (0,1,2) を指定します。ここでは 0 = Full Text、1 = Structured Text、2 = XML になります
encoding	Automatic (自動) "UTF-8" "UTF-16" "ISO-8859-1" "US-ASCII" "CP850" "EUC-JP" "SHIFT-JIS" "ISO2022-JP"	"UTF-8" のような特殊文字と値の組み合わせは、算術演算子と区別するために引用符 (") で囲む必要があります。

表 12. テキスト リンク分析 (TLA) ノードのスキプトのプロパティ (続き)

スキプトのプロパティ	データ型	プロパティの説明
unity	整数	正の数 (0,1) を指定します。ここでは、0 = Paragraph および 1 = Document になります
para_min	整数	
para_max	整数	
mtag	文字列	すべての mtag 設定が含まれます (XML ファイルの「設定」ダイアログ・ボックス)
mclef	文字列	すべての mclef 設定が含まれます (構造テキスト・ファイルの「設定」ダイアログ・ボックス)
language	de en es fr it ja nl pt	
concept_count_limit	整数	抽出を少なくともこの値以上のグローバル度数のコンセプトに制限。
fix_punctuation	フラグ	
fix_spelling	フラグ	
spelling_limit	整数	
extract_uniterm	フラグ	
extract_nonlinguistic	フラグ	
upper_case	フラグ	
group_names	フラグ	
permutation	整数	語順が異なる類義語で無視する非機能語の数 (デフォルトは 3)。

第 7 章 インタラクティブ・ワークベンチ・モード

テキストマイニングモデル作成ノードから、ストリーム実行時にインタラクティブ・ワークベンチ・セッションを起動することができます。このワークベンチでは、テキスト・データから主要なコンセプトを抽出、カテゴリーを作成、そしてテキストリンク分析パターンおよびクラスターを検討、カテゴリー・モデルを生成できます。この章では、次のような主要な要素とともに、ワークベンチ・インターフェースについて高レベルな視点から説明します。

- 抽出結果。抽出が実行された後、テキストデータから特定および抽出されるキー・ワードおよび句で、「コンセプト」とも呼ばれます。これらのコンセプトは、「タイプ」にグループ化されます。これらのコンセプトおよびタイプを使用して、カテゴリーを作成するほか、データを検討できます。これらは、コンセプトおよびカテゴリービューで管理されます。
- カテゴリー: 抽出結果、パターン、ルールなどの記述子を定義として使用し、カテゴリー定義の一部を含むかどうかに基づいてドキュメントおよびレコードが割り当てられるカテゴリーのセットを手動でまたは自動的に作成できます。これらは、コンセプトおよびカテゴリービューで管理されます。
- クラスター。クラスターは、間にリンクがあるコンセプトのグループ化で、コンセプト間の関係を示します。コンセプトは、その他の要素間で、2つのコンセプトがそれぞれ現れる頻度と比較して2つのコンセプトが同時に現れる頻度を使用する複雑なアルゴリズムを使用してグループ化します。これらは、クラスタービューで管理されます。クラスターを構成するコンセプトをカテゴリーに追加することもできます。
- テキストリンク分析パターン: 言語リソースにテキストリンク分析 (TLA) パターンの規則がある場合、または既にいくつかの TLA 規則があるリソース・テンプレートを使用している場合、テキスト・データからパターンを抽出できます。これらのパターンを使用して、データのコンセプト間に興味深い関連を見つけることができます。また、これらのパターンをカテゴリーの記述子として使用することもできます。これらは、「テキストリンク分析」ビューで管理されます。
- 言語リソース: 抽出プロセスは、テキストの抽出方法および処理方法を支配する一連のパラメーターおよび言語定義によって異なります。これらは、リソース・エディタービューのテンプレートおよびライブラリーのフォームで管理されます。

インタラクティブ ワークベンチで発生する可能性がある問題

- 複数のインタラクティブ ワークベンチ セッションを使用すると、動作が遅くなる可能性があります。SPSS Modeler Text Analytics および SPSS Modeler は、インタラクティブ ワークベンチ セッションを起動するときに共通の Java ランタイム エンジンと共有します。SPSS Modeler セッション中に起動するインタラクティブ ワークベンチ セッションの数によっては、同じセッションを開いた後に閉じる場合であっても、システムメモリーのためにアプリケーションの動作が遅くなる場合があります。大量のデータを処理する場合や、推奨 RAM 設定 (4GB) に満たないマシンを使用する場合は、この影響が特に顕著になります。マシンの応答が遅いと感じる場合は、すべての作業を保存して SPSS Modeler をシャットダウンし、アプリケーションを再起動することをお勧めします。推奨メモリー未満のマシンで SPSS Modeler Text Analytics を実行すると、特に大規模なデータ・セットを処理する場合や長時間にわたって作業する場合に、Java のメモリーが不足してシャットダウンしてしまふことがあります。大量のデータを処理する場合は、推奨メモリー設定以上にアップグレードする (または SPSS Modeler Text Analytics Server を使用する) ことを強くお勧めします。

- アプリケーションを再起動せずに複数の SPSS Modeler Text Analytics インタラクティブ ワークベンチ セッションを実行した後に、SPSS Modeler Client のメモリーが不足する場合があります。状況表示行でメモリー使用量をモニタし、少なくなったら、SPSS Modeler Client を閉じてから再び開いてください。

カテゴリーとコンセプト・ビュー

アプリケーション・インターフェースは、いくつかのビューで構成されています。カテゴリーとコンセプト・ビューは、抽出結果を検討および調整するほか、カテゴリーを作成および検討するウィンドウです。カテゴリーは、スコアリング・プロセスでドキュメントおよびレコードが割り当てられる、密接に関連するキーワードおよびパターンのグループを参照します。コンセプトは、カテゴリーの構築ブロック（記述子）として使用できる最も基本的なレベルの抽出結果を参照します。

Category	Descriptors	Docs
exercise	1	
feature		5
hardware		3
headphones		2
home		3
internet		2
listening		3
look		2
memory device		12
music		27
Neg: General Dissatisfaction		24
Neg: Pricing and Billing		9
Neg: Product Dissatisfaction		43
Neg: Service Dissatisfaction		42
occupation		2

Concept	In	Global	Docs	Type
small		58 (5%)	58 (14%)	<Contextual>
music		54 (4%)	51 (13%)	<Features>
easy to use		45 (4%)	44 (11%)	<Positive>
like		55 (5%)	43 (11%)	<Positive>
portable		44 (4%)	43 (11%)	<Positive>
size		36 (3%)	36 (9%)	<Characteristics>
sound		34 (3%)	33 (8%)	<Features>
excellent		39 (3%)	32 (8%)	<Positive>
good		31 (3%)	30 (7%)	<Positive>
listening		30 (2%)	29 (7%)	<Unknown>
songs		29 (2%)	26 (6%)	<Unknown>
large		20 (2%)	20 (5%)	<Contextual>
product		19 (2%)	18 (4%)	<Products>
battery		16 (1%)	16 (4%)	<Performance>
design		15 (1%)	15 (4%)	<Characteristics>
cds		13 (1%)	13 (3%)	<Products>
lightweight		12 (1%)	12 (3%)	<Contextual>

Q1: What do you like most about this portable music player? (33)	Categories
1 like that Product A has a lot of storage. Also, the interface is very easy to use.	memory device/memory
2 Everything! Product A rules! I can't wait to get a [redacted] one!	memory device/recording/video
3 I can store a lot of music on it.	memory device/memory
4 Convenience of storing all my music in one device	memory device/memory
5 Large storage capacity	memory
6 Small size. It has 512Mb of add-on memory, so it is quick to load and play music. It can also encode directly from external devices from the radio or a CD player.	consumer electronics memory device/memory music radio size
7 storage capacity	memory
8 Small but lots of space (60 GB). [redacted] is a bit of a toy but cool.	memory device/recording/video space

図 23. 「カテゴリーとコンセプト」ビュー

カテゴリーとコンセプト・ビューは 4 つのパネルで構成され、「表示」メニューから名前を選択して隠したり表示したりできます。詳しくは、95 ページの『第 9 章 テキストデータのカテゴリー化』のトピックを参照してください。

カテゴリー・ペイン

左上にあるこの領域は、構築したカテゴリーを管理できる表が表示されます。テキスト・データからコンセプトとタイプを抽出した後、セマンティック・ネットワークや内包関係のコンセプトなどの方法を使用して、または手動で作成してカテゴリーを構築できます。カテゴリー名をダブルクリックすると、「カテゴリー

一定義」ダイアログボックスが開き、コンセプト、タイプ、規則など、定義を構成するすべての記述子が表示されます。詳しくは、95 ページの『第 9 章 テキストデータの 카테고리化』のトピックを参照してください。すべての言語ですべての自動的手法が使用できるわけではありません。

パネルで 1 行選択すると、データ・ペインおよび視覚化ペインに該当するドキュメント/レコードまたは記述子に関する情報が表示されます。

抽出結果ペイン

左下のこの領域には、抽出結果が表示されます。抽出を実行すると、抽出エンジンがテキスト・データを読み込み、関連コンセプトを特定し、それぞれにタイプを割り当てます。コンセプトは、テキスト・データから抽出した単語や句です。タイプは、キーワード辞書の形式で保存されたコンセプトのセマンティックグループです。抽出が完了すると、コンセプトとタイプが抽出結果ペインにカラー・コード化されて表示されます。詳しくは、81 ページの『抽出結果: コンセプトとタイプ』のトピックを参照してください。

コンセプト名の上にマウスポインタを置くと、コンセプトの基本キーワードのセットが表示されます。これにより、コンセプト名とそのコンセプトにグループ化された数行のキーワードを示すツールヒントが表示されます。これらの基本キーワードには、抽出された複数形/単数形のキーワード、置換キーワード、Fuzzy Grouping のキーワードなどのほか、テキスト内にあったかどうかに関係なく、言語リソースで定義された類義語が含まれます。コンセプト名を右クリックし、コンテキスト・メニューのオプションを選択して、これらのキーワードをコピーしたり、基本キーワードの完全セットを表示したりできます。

テキストマイニングは、抽出結果をテキスト・データのコンテキストに従ってレビューし、新しい結果を作成するよう調整、そして再評価するインタラクティブプロセスです。言語リソースを修正することによって、抽出結果を修正できます。この調整は、抽出結果ペインまたはデータ・ペインから部分的に直接、またはリソース・エディター・ビューから直接実行できます。詳しくは、74 ページの『リソース・エディター・ビュー』のトピックを参照してください。

注: 表示中のペインに結果を表示しきれない場合は、ペインの下部にあるコントロールを使用して前後の結果に移動したり、移動先のページ番号を入力したりすることができます。

視覚化パネル

右上にあるこの領域には、ドキュメント/レコードのカテゴリ化の共通性について、さまざまな観点が表示されます。各グラフやチャートは類似の情報を提供しますが、異なる方法または異なる詳細レベルで表示します。これらの図表やグラフを使用して、カテゴリ化の結果を分析したり、カテゴリまたはレポートの調整を行うことができます。例えば、グラフを使用して、あまりに類似している (75% を超えるレコードを共有しているなど) またはあまりに異なるカテゴリを見つけることができます。グラフまたは図表の内容は、その他のパネルでの選択内容に対応しています。詳しくは、155 ページの『カテゴリ・グラフおよび図表』のトピックを参照してください。

データ パネル

データ・ペインは、右下に表示されます。このウィンドウには、ビューの別の領域での選択内容に対応するドキュメントまたはレコードを示すテーブルが表示されます。選択内容に応じて、対応するテキストのみがデータ・ペインに表示されます。選択した後、「表示」ボタンをクリックすると、データ・ペインに対応するテキストが入力されます。

別のパネルで選択した場合、該当するドキュメントまたはレコードにはコンセプトが色付きで強調表示され、テキスト内のコンセプトを特定しやすくなります。カラーコード化された項目上でマウス・ポインタを停止させて、項目が抽出されたコンセプトの名前と、項目が割り当てられたタイプを示すヒントを表示するこ

ともできます。詳しくは、104 ページの『データ・ペイン』のトピックを参照してください。

コンセプトとカテゴリー・ビューでの検索

特定のセクションで、情報の迅速な検索が必要な場合があります。検索ツールバーを使用して、検索する文字列を入力し、大文字や小文字の区別および検索方向など、その他の検索基準を定義することができます。そして、検索するパネルを選択できます。

検索機能を使用するには

1. カテゴリーとコンセプト・ビューのメニューで、「編集」>「検索」を選択します。カテゴリー・ペインおよび視覚化ペインの上に検索ツールバーが表示されます。
2. テキスト・ボックスに検索したい文字列を入力します。ツールバー・ボタンを使用して、大文字と小文字の区別、部分一致、検索の方向を制御します。
3. ツールバーで、検索するパネル名をクリックします。一致が見つかった場合は、テキストがウィンドウで強調表示されます。
4. 次の一致を検索するには、パネルの名前をもう一度クリックします。

クラスター・ビュー

クラスター・ビューでは、テキスト・データのクラスター結果を構築および検討できます。クラスターは、コンセプトが出現する頻度およびいっしょに出現する頻度に基づいてアルゴリズムをクラスター化することによって生成されるコンセプトのグループです。カテゴリーの目的は、含まれるテキストが各カテゴリーの記述子 (コンセプト、条件規則、パターン) にどのように合致するかに基づいてドキュメントまたはレコードをグループ化することですが、クラスターの目的は共起するコンセプトをグループ化することです。

クラスター内のコンセプトが、その他のコンセプトと共に出現する低い頻度のコンセプトと共により頻繁に出現すればするほど、クラスターが興味深いコンセプトの関係性をより適切に特定します。2 つのコンセプトが同じドキュメントまたはレコードで現れると (または類義語かキーワードのいずれかが現れると) それらは共起します。詳しくは、143 ページの『第 10 章 クラスターの分析』のトピックを参照してください。

クラスターを構築し、見つけるのに時間がかかるコンセプト間の関係を発見することができる一連の図表およびグラフでクラスターを検討できます。クラスター全体をカテゴリーに追加できませんが、「クラスター定義」ダイアログ・ボックスを使用してクラスターのコンセプトをカテゴリーに追加できます。詳しくは、147 ページの『クラスター定義』のトピックを参照してください。

クラスターリングの設定に変更を行うと、結果に影響を与える場合があります。詳しくは、144 ページの『クラスターの作成』のトピックを参照してください。

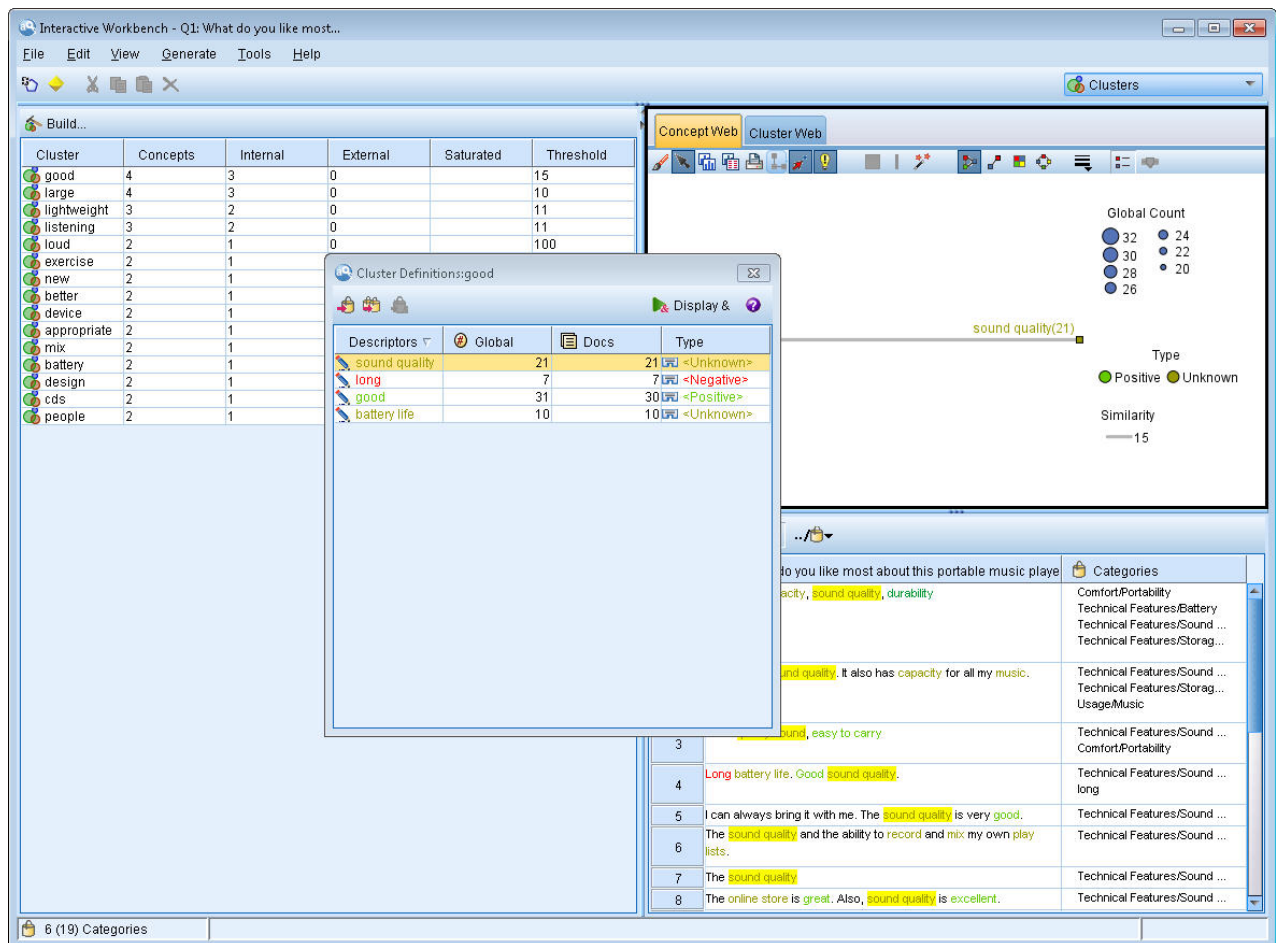


図 24. クラスタ・ビュー

クラスタ・ビューは 3 つのパネルで構成され、「表示」メニューから名前を選択して隠したり表示したりできます。通常、表示されるのはクラスタ・ペインと視覚化ペインだけです。

クラスタ・ペイン

左側にあるこの領域には、テキスト・データで見つかったクラスタが表示されます。「作成」ボタンをクリックして、クラスタリングの結果を作成できます。クラスタは、クラスタリング・アルゴリズムによって形成され、頻繁に共起するコンセプトを特定しようとしています。

新しい抽出が行われると、クラスタ結果は消去され、クラスタを再構築して最新の結果を取得する必要があります。クラスタを構築している場合、作成する最大クラスタ数、含むことができる最大コンセプト数、使用できる外部コンセプトとの最大リンク数など、いくつかの設定を変更できます。詳しくは、146 ページの『クラスタの検証』のトピックを参照してください。

視覚化パネル

右上の角にあるこのパネルはクラスタ化の2つの観点を提供します：それは、コンセプト Web グラフとクラスタ Web グラフです。表示されない場合、「表示」メニュー（「表示」>「視覚化」）からこの領域にアクセスできます。クラスタのパネルでの選択内容によって、クラスタ間またはクラスタ内の該当する交互作用を表示できます。次のような形式で結果を表示します。

- コンセプト **Web**: クラスター外のリンクしたコンセプトのほか、選択したクラスター内のすべてのコンセプトを表示する Web グラフ。
- クラスター **Web**: その他のクラスター間のリンクのほか、選択したクラスターからのリンクを表示する Web グラフ。

注: クラスター Web グラフを表示するには、外部リンクを持つクラスターを先に構築する必要があります。外部リンクは、別々のクラスターにあるコンセプトのペア (あるクラスターのコンセプトと外部の別のクラスターのコンセプトとの間) 間のリンクです。詳しくは、157 ページの『クラスター・グラフ』のトピックを参照してください。

データ パネル

データ・ペインは、右下にあり、デフォルトでは表示されません。これらのクラスターの範囲は複数のドキュメント/レコードにわたり、データ結果は興味深いものではなくなるため、「クラスター」結果からデータ・ペインの結果を表示できません。ただし、「クラスター定義」ダイアログ・ボックス内の選択に対応するデータを表示できます。そのダイアログ・ボックスの選択内容に応じて、データ・ペインに対応するテキストのみが表示されます。選択を行うと、「表示」& ボタンをクリックして、データ・ペインに、すべてのコンセプトを同時に含むドキュメントまたはレコードを表示します。

該当するドキュメントまたはレコードには、コンセプトを色付きで強調表示し、テキスト内のコンセプトを特定しやすくします。カラーコード化された項目上でマウス・ポインタを停止させて、項目が抽出されたコンセプトと、項目が割り当てられたタイプを表示することもできます。データ・ペインには、複数の列が表示されますが、テキスト・フィールド列は常に表示されます。抽出時に使用されたテキスト・フィールドの名前、またはテキスト・データがさまざまなファイルにある場合はドキュメント名が表示されます。その他の列も使用できます。詳しくは、104 ページの『データ・ペイン』のトピックを参照してください。

テキスト リンク分析ビュー

テキスト リンク分析ビューでは、テキスト・データで見つかったテキスト リンク分析パターンを作成および検証できます。テキスト リンク分析 (TLA) はパターンマッチ手法で、TLA 規則を定義し、それらをテキスト内の実際の抽出されたコンセプトおよび関連性と比較することができます。

パターンは、コンセプト間の関連性または特定のサブジェクトに関する意見を探索する場合に最も役立ちます。例えば、調査データで製品に関する意見、医療調査アンケートから遺伝子的関連性、または情報データから人名と地名との関連性を抽出したい場合などです。

TLA パターンを抽出すると、データ・ペインまたは視覚化ペインで検証し、カテゴリとコンセプト・ビューでそれらの結果をカテゴリに追加することができます。いくつかの TLA 規則が、TLA 結果を抽出するために使用するリソース・テンプレートまたはライブラリーで定義されています。詳しくは、215 ページの『第 18 章 テキスト リンク規則について』のトピックを参照してください。

TLA パターン規則の抽出を選択している場合、結果はこのビューに表示されます。抽出を選択していない場合は、「抽出」ボタンを使用して、パターンの抽出を可能にするオプションを選択する必要があります。

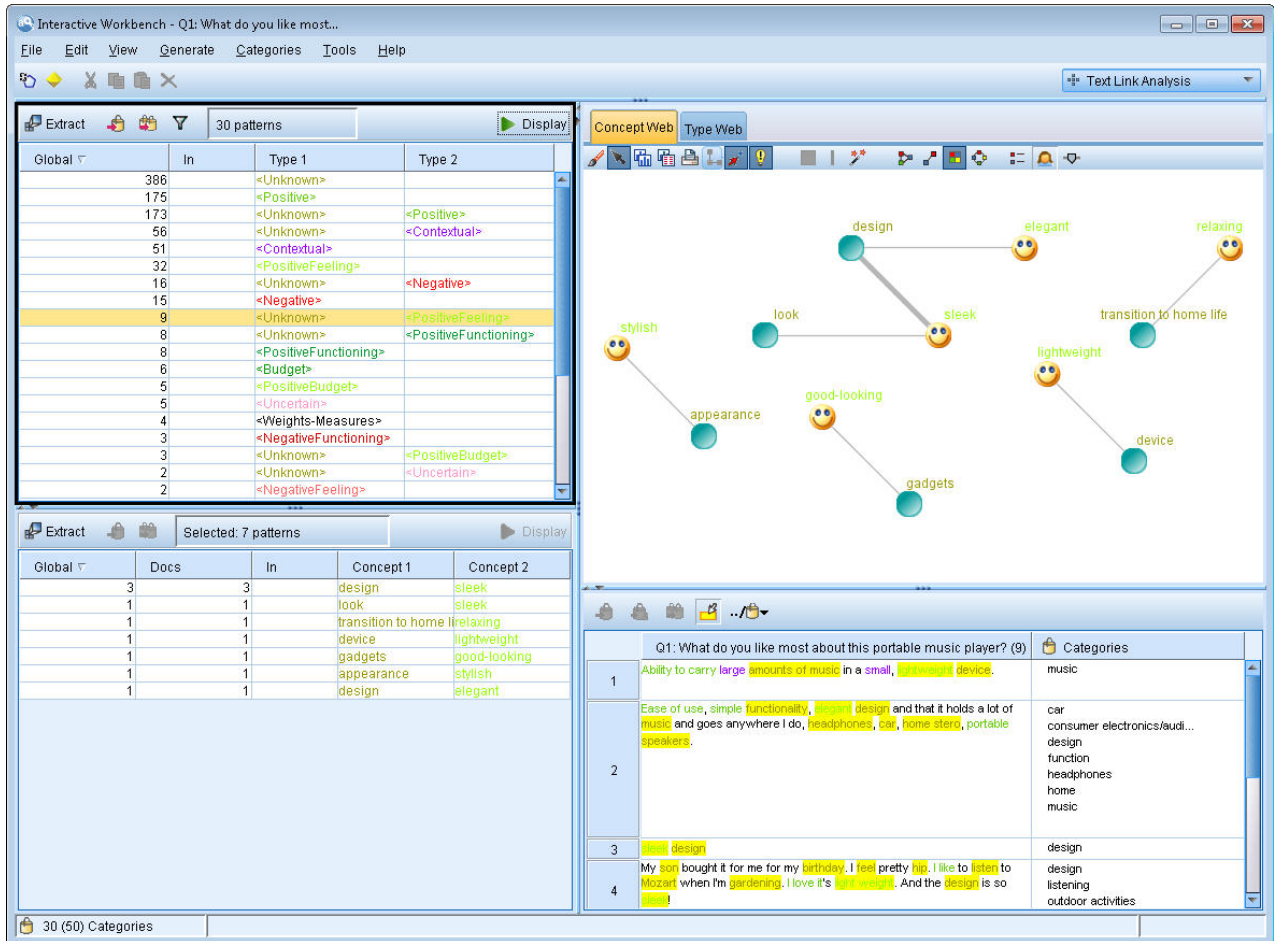


図 25. テキスト リンク分析ビュー

テキスト リンク分析ビューは 4 つのパネルで構成され、「表示」メニューから名前を選択して隠したり表示したりできます。詳しくは、149 ページの『第 11 章 テキスト リンク分析の検証』のトピックを参照してください。

タイプ・パターン・ペインおよびコンセプト・パターン・ペイン

左側のタイプ・パターン・ペインおよびコンセプト・パターン・ペインは、TLA パターン結果を検証および選択できる 2 つの関連したパネルです。パターンは、6 つのタイプまたは 6 つのコンセプトで構成されています。言語リソースに定義されているため、TLA パターン規則は、パターン結果の複雑さを示します。詳しくは、215 ページの『第 18 章 テキスト リンク規則について』のトピックを参照してください。

パターン結果はまずタイプ レベルでグループ化され、コンセプト・パターンに分割されます。このため、2 つの異なる結果パネルがあります：それが、タイプ・パターン（左上）とコンセプト・パターン（左下）です。

- **タイプ・パターン**：「タイプ・パターン」は、TLA パターン規則を満たす 2 つ以上の関連タイプで構成されている抽出パターンが表示されます。タイプ・パターンは、<組織名> + <地名> + <肯定的> と表され、特定の場所の組織について、肯定的なフィードバックを提供します。

- **コンセプト・パターン:**コンセプト・パターン・ペインには、上の「タイプ・パターン」で現在選択されているすべてのタイプ・パターンのコンセプト・レベルで抽出パターンが表示されます。コンセプト・パターンは、ホテル + パリ + すばらしい などの構造に従います。

カテゴリとコンセプト・ビューの抽出結果と同様、ここで結果を確認できます。これらのパターンを構成するタイプおよびコンセプトに調整を行う場合、カテゴリとコンセプト・ビューの抽出結果ペインまたはリソース・エディターで変更を行うか、パターンを再抽出します。

視覚化パネル

テキスト リンク分析ビューの右上のこのパネルには、選択したパターンの Web グラフがタイプ・パターンまたはコンセプト・パターンのいずれかとして表示されます。表示されない場合、「表示」メニュー（「表示」 > 「視覚化」）からこの領域にアクセスできます。その他のパネルでの選択内容によって、ドキュメント/レコードおよびパターンの間の該当する相互作用を表示できます。

次のような形式で結果を表示します。

- **コンセプト グラフ:** このグラフには、選択したパターンのすべてのコンセプトを示します。コンセプトグラフの線の幅およびノードのサイズ (タイプ・アイコンが表示されていない場合) には、選択したテーブルのグローバル出現値を示します。
- **タイプ グラフ:** このグラフには、選択したパターンのすべてのタイプを示します。グラフの線の幅およびノードのサイズ (タイプ・アイコンが表示されていない場合) には、選択したテーブルのグローバル出現値を示します。ノードは、タイプ カラーまたはアイコンによって示されます。

詳しくは、158 ページの『テキスト リンク分析のグラフ』のトピックを参照してください。

データ パネル

データ・ペインは、右下に表示されます。このウィンドウには、ビューの別の領域での選択内容に対応するドキュメントまたはレコードを示すテーブルが表示されます。選択内容に応じて、対応するテキストのみがデータ・ペインに表示されます。選択した後、「表示」 ボタンをクリックすると、データ・ペインに対応するテキストが入力されます。

別のパネルで選択した場合、該当するドキュメントまたはレコードにはコンセプトが色付きで強調表示され、テキスト内のコンセプトを特定しやすくします。カラーコード化された項目上でマウス・ポインタを停止させて、項目が抽出されたコンセプトの名前と、項目が割り当てられたタイプを示すヒントを表示することもできます。詳しくは、104 ページの『データ・ペイン』のトピックを参照してください。

リソース・エディター・ビュー

IBM SPSS Modeler Text Analytics は、頑健な抽出エンジンを使用して、主要なコンセプトをテキストデータから迅速かつ正確にキャプチャーします。このエンジンは、大容量の非構造化テキスト・データをどのように分析および解釈するかを示す言語リソースによって大きく異なります。

リソース・エディター ビューでは、コンセプトの抽出、タイプに基づいたグループ化、テキスト・データでのパターンの検出などに使用する言語リソースを表示および調整できます。IBM SPSS Modeler Text Analytics では、事前設定されたリソース・テンプレートをいくつか用意しています。また、一部の言語では、テキスト分析パッケージのリソースを使用することもできます。詳しくは、136 ページの『テキスト分析パッケージの使用』のトピックを参照してください。

これらのリソースは、常にデータのコンテキストに完全に対応しているとは限らないため、リソース・エディターで特定のコンテキストまたはドメインの独自のリソースを作成、編集および管理できます。詳しくは、179 ページの『第 15 章 ライブラリーの使用』のトピックを参照してください。

言語リソースの調整プロセスを簡略化するために、一般的な辞書タスクを、抽出結果ペインおよびデータ・ペインのコンテキスト・メニューを使用して、カテゴリとコンセプトビューから直接実行できます。詳しくは、89 ページの『抽出結果の調整』のトピックを参照してください。

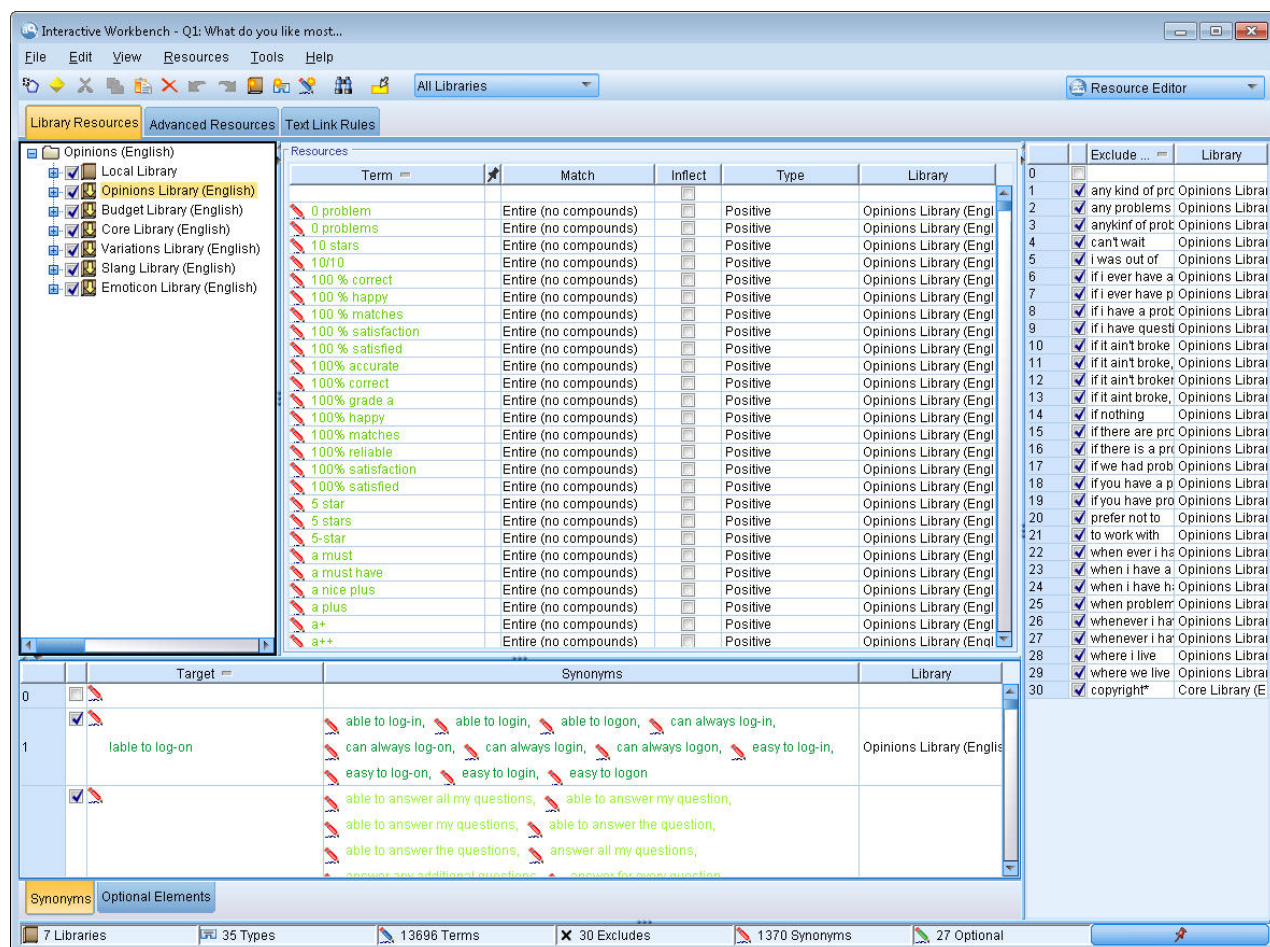


図 26. リソース・エディター・ビュー

リソース・エディターで実行する操作は、言語リソースの管理および調整を中心に展開しています。これらのリソースは、テンプレートおよびライブラリーの形で保存されています。リソース・エディタービューは4つに編成されています：ライブラリー・ツリー・ペイン、キーワード辞書ペイン、代替辞書ペイン、それと除外辞書ペイン。

注: 詳しくは、168 ページの『エディターのインターフェース』のトピックを参照してください。

オプションの設定

「オプション」ダイアログ・ボックスで IBM SPSS Modeler Text Analytics の一般的なオプションを設定できます。このダイアログボックスには、以下のようなタブがあります。

- セッション。 このタブには、一般オプションおよび区切りがあります。
- ディスプレイ。 インターフェイスで使用される色についてのオプションがあります。
- サウンド。 サウンド・キューについてのオプションがあります。

オプションを編集するには：

1. メニューの「ツール」>「オプション」を選択します。「オプション」ダイアログボックスが開きます。
2. 変更する情報を含むタブを選択します。
3. オプションのいずれかを変更します。
4. 「OK」をクリックして、変更を保存します。

オプション：「セッション」タブ

このタブで、基本的な設定をいくつか定義できます。

データウィンドウ枠とカテゴリー・グラフの表示：カテゴリーとコンセプト・ビューのデータ・ペインおよび視覚化ペインにデータがどのように表示されるかを指定します。

- データ・パネルおよびカテゴリー **Web** の制限を表示：カテゴリーとコンセプト・ビューのデータ・ペインまたはグラフおよび図表の入力に表示または使用するドキュメントの最大数を設定します。
- 表示時にドキュメント/レコードのカテゴリーを表示：「表示」をクリックするとドキュメントまたはレコードがスコアリングされ、それらが属するカテゴリーがデータ・ペインの「カテゴリー」列およびカテゴリー・グラフに表示されます。特にデータセットが大きい場合は、データおよびグラフをより早く表示するよう、このオプションを無効にする必要があります。

データパネルからカテゴリーに追加：ドキュメントおよびレコードがデータ・ペインから追加される場合、カテゴリーに追加する内容を設定します。

- 「「カテゴリーとコンセプト」ビューでコピー このビューのデータ・ペインからドキュメントまたはレコードを追加すると、「コンセプトのみ」または「コンセプトとパターン」でコピーされます。
- テキスト リンク分析ビューでコピー：このビューのデータ・ペインからドキュメントまたはレコードを追加すると、「パターンのみ」または「コンセプトとパターン」でコピーされます。

リソース・エディターの区切り文字：コンセプト、類義語、オプション要素などの要素をリソース・エディター・ビューで入力する場合に区切り文字として使用する文字を選択します。

オプション：「表示」タブ

このタブでは、全体的な外観に関するオプションや、要素を区別するための色を編集できます。

注：製品の表示を以前のリリースのクラシック表示に切り替えるには、IBM SPSS Modeler メイン・ウィンドウにある「ツール」メニューの「ユーザー オプション」ダイアログを開きます。

ユーザー定義の色：画面上に表示される要素の色を編集します。表内のそれぞれの要素に関して、色を変更できます。ユーザー定義の色を指定するには、変更したい要素の右側にある色をクリックし、色のドロップダウンリストから色を変更します。

- 未抽出のテキスト：データ・ペインで抽出されていないが表示されるテキスト・データ。

- **強調背景:** パネルの要素またはデータ・ペインのテキストを選択する場合のテキスト選択の背景色。
- **抽出が必要な背景:** 「抽出結果」、「パターン」、クラスター・ペインの背景色はライブラリーに変更が行われたことを示し、抽出が必要です。
- **カテゴリー フィードバック背景:** 操作後に出現するカテゴリー背景色。
- **デフォルト タイプ:** データ・ペインおよび抽出結果ペインに出現するタイプおよびコンセプトのデフォルト色。この色は、リソース・エディターで作成するカスタム タイプに適用されます。リソース・エディター でこれらのキーワード辞書のプロパティを編集し、カスタム・キーワード辞書のこのデフォルト色を上書きします。詳しくは、191 ページの『キーワード辞書の作成』のトピックを参照してください。
- **テーブルの縞 1:** 各セットの行を区別するための、「強制コンセプトを編集」ダイアログボックスのテーブルで、交互に使用される 2 つの色のうちの最初の色。
- **テーブルの縞 2:** 各セットの行を区別するための、「強制コンセプトを編集」ダイアログボックスのテーブルで、交互に使用される 2 つの色のうちの 2 番目の色。

注: 「デフォルトに戻す」 ボタンをクリックすると、このダイアログ・ボックスのすべてのオプションは、この製品を最初にインストールしたときに設定されていた値に戻されます。

オプション: 「サウンド」 タブ

このタブでは、サウンドに関するオプションを編集できます。サウンドを鳴らすイベントの部分で、あるイベントが起こった際にこれを知らせるサウンドを指定できます。使用できるサウンドはたくさんあります。サウンドを参照して選択する場合は、「...」 ボタンを使用します。IBM SPSS Modeler Text Analytics のサウンドを作成するために使われる .wav ファイルは、インストール・ディレクトリー中の *media* サブディレクトリーにあります。サウンドを鳴らしたくない場合、「すべてのサウンドをミュート」 を選択します。デフォルトでは、音が鳴らないようになっています。

注: 「デフォルトに戻す」 ボタンをクリックすると、このダイアログ・ボックスのすべてのオプションは、この製品を最初にインストールしたときに設定されていた値に戻されます。

Microsoft Internet Explorer ヘルプの設定

Microsoft Internet Explorer の設定

このアプリケーションのほとんどのヘルプ機能では、Microsoft Internet Explorer に基づいたテクノロジーが使用されています。Internet Explorer のバージョンによっては (Microsoft Windows XP、Service Pack 2 と共に提供されるバージョンも含む)、ローカル・コンピューター上の「Internet Explorer」ウィンドウ内で「アクティブなコンテンツ」と見なされる対象が、デフォルトにより封鎖されます。このデフォルトの設定により、ヘルプ機能内である種のコンテンツが表示されなくなります。すべてのヘルプ・コンテンツを表示するために、Internet Explorer のデフォルトの動作を変更できます。

1. Internet Explorer のメニューから次の項目を選択します。

「ツール」 > 「インターネット オプション...」

2. 「詳細設定」 タブをクリックします。
3. 「セキュリティ」 セクションまで下方へスクロールします。
4. 「マイ コンピューターのファイルでのアクティブ・コンテンツの実行を許可する」を選択します。

モデル ナゲットおよびモデル作成ノードの生成

インタラクティブ・セッションの場合、実行した作業を使用して、次のいずれかを作成する必要があります。

- **テキスト マイニング モデル作成ノード:** インタラクティブ・ワークベンチ・セッションから生成したモデル作成ノードは、設定とオプションが、オープン・インタラクティブ・セッションに保存されている設定およびオプションを反映するテキスト マイニング・ノードです。元のテキスト マイニング・ノードがない場合、または新しいバージョンを作成したい場合、役立ちます。詳しくは、19 ページの『第 3 章 コンセプトおよびカテゴリーのマイニング』のトピックを参照してください。
- **カテゴリー モデル ナゲット:** インタラクティブ・ワークベンチ・セッションから生成されたモデル ナゲットは、カテゴリー モデル ナゲットです。カテゴリー モデル ナゲットを生成するには、カテゴリーとコンセプト・ビューに 1 つ以上のカテゴリーが必要です。詳しくは、40 ページの『テキスト マイニング モデル ナゲット:カテゴリー・モデル』のトピックを参照してください。

テキスト マイニング モデル作成ノードを生成するには

1. メニューの「生成」>「モデル作成ノードの生成」を選択します。テキスト マイニング モデル作成ノードは、ワークベンチ・セッションで現在すべての設定を使用する作業キャンバスに追加されます。ノードには、テキスト・フィールドの名前が付きます。

カテゴリー モデル ナゲットを生成するには

1. メニューの「生成」>「モデルの生成」を選択します。モデル ナゲットは、モデル・パレットに直接生成され、デフォルト名が付きます。

モデル作成ノードの更新および保存

インタラクティブ・セッションで作業する場合、時々モデル作成ノードを変更して、変更を保存することをお勧めします。また、インタラクティブ・ワークベンチ・セッションで作業を終了し、作業を保存する場合も、モデル作成ノードを更新する必要があります。モデル作成ノードを更新する場合、ワークベンチ・セッションの内容が、インタラクティブ・ワークベンチ・セッションに由来するテキスト マイニング・ノードに保存します。更新しても、出力ウィンドウは閉じません。

重要: 更新するとストリームは保存されません。ストリームを保存するには、モデル作成ノードを更新した後、IBM SPSS Modeler のメイン・ウィンドウで保存します。

モデル作成ノードを更新するには

1. メニューの「ファイル」>「モデル作成ノードを更新」を選択します。オプションおよびカテゴリーとともに、作成設定および抽出設定でモデル作成ノードを更新します。

セッションの終了

セッションの作業を終了する場合、次の 3 つの終了でセッションを離れることができます。

- **保存:** まず、別のセッションで再利用するためにライブラリーを公開し、次のセッションのために元のモデル作成ノードに作業を保存します。詳しくは、184 ページの『ライブラリーの共有』のトピックを参照してください。保存した後、セッションが終了し、IBM SPSS Modeler ウィンドウの出力マネージャーからセッションが削除されます。
- **終了:** 保存していない作業を破棄し、セッション・ウィンドウを閉じて、IBM SPSS Modeler ウィンドウの出力マネージャーからセッションが削除されます。メモリを確保するために、重要な作業を保存して、セッションを終了することをお勧めします。

- 閉じる: 作業は保存されず、また破棄もされません。セッション・ウィンドウが閉じますが、セッションは稼働し続けます。IBM SPSS Modeler ウィンドウの出力マネージャーでこのセッションを選択すると、セッション・ウィンドウを開くことができます。

ワークベンチのセッションを終了するには

1. メニューから、「ファイル」>「閉じる」を選択します。

キーボード・アクセシビリティ

インタラクティブ・ワークベンチのインターフェイルには、製品の機能によりアクセスしやすくするキーボード・ショートカットが用意されています。基本的には、Alt キーと他の適切なキーを同時に押してメニュー項目を選択したり (例: Alt + F キーで「ファイル」メニューを選択)、Tab キーを使用してダイアログ・ボックス中のコントロール間を移動することができます。ここでは、もう 1 つのナビゲーションであるキーボード・ショートカットについて説明します。IBM SPSS Modeler インターフェースには、その他のショートカットがあります。

表 13. 一般的なキーボード・ショートカット

ショートカット・キー	関数
Ctrl+1	タブのあるパネルの最初のタブを表示します。
Ctrl+2	タブのあるパネルの 2 番目のタブを表示します。
Ctrl+A	フォーカスのあるパネルのすべての要素を選択します。
Ctrl+C	選択したテキストをクリップボードにコピーします。
Ctrl+E	カテゴリーとコンセプトビューおよびテキストリンク分析ビューの抽出を起動します。
Ctrl+F	リソース・エディター/テンプレート・エディター の「検索」ツールバーが表示されていない場合は表示し、フォーカスします。
Ctrl+I	カテゴリーとコンセプト・ビューの場合、選択したカテゴリーの「カテゴリー定義」ダイアログ・ボックスを起動します。クラスター・ビューの場合、選択したクラスターの「クラスター定義」ダイアログ・ボックスを起動します。
Ctrl+R	リソース・エディター/テンプレート・エディター の「複数のキーワードを追加」を開きます。
Ctrl+T	リソース・エディター/テンプレート・エディター で「タイプのプロパティ」ダイアログ・ボックスを開き、新しいタイプを作成します。
Ctrl+V	クリップボードの内容を貼り付けます。
Ctrl+X	リソース・エディター/テンプレート・エディター から選択した項目を切り取ります。
Ctrl+Y	ビューの最後のアクションをやり直します。
Ctrl+Z	ビューの最後のアクションを取り消します。
[F1] キー	ヘルプを表示するか、ダイアログ・ボックスでは、その項目のコンテキスト ヘルプを表示します。
[F2] キー	テーブルのセルの編集モードを有効にしたり、無効にしたりします。
[F6] キー	アクティブなビューの主なパネル間でフォーカスを移動します。
F8	フォーカスをパネルの分割バーに移動し、サイズを変更します。
F10	メインの「ファイル」メニューを展開します。
上方向矢印、下方向矢印	分割バーが選択されているときに、ウィンドウを垂直方向にサイズ変更します。

表 13. 一般的なキーボード・ショートカット (続き)

ショートカット・キー	関数
左方向矢印、右方向矢印	分割バーが選択されているときに、ウィンドウを水平方向にサイズ変更します。
Home、End	分割バーが選択されているときに、ウィンドウを最大サイズまたは最小サイズに変更します。
タブ	ウィンドウ、パネル、ダイアログ・ボックスの次の項目に移動します。
Shift+F10	項目のコンテキスト・メニューを表示します。
Shift+Tab	ウィンドウ、またはダイアログ・ボックスの前の項目に移動します。
Shift + 矢印	編集モード (F2) のとき、編集フィールドの文字を選択します。
Ctrl+Tab	ウィンドウの次のメイン領域にフォーカスを移動します。
Shift + Ctrl + Tab	ウィンドウの前のメイン領域にフォーカスを移動します。

ダイアログ・ボックスのショートカット

ダイアログ・ボックスを使用している場合、いくつかのショートカットおよびスクリーン・リーダー キーが役立ちます。ダイアログ・ボックスに入力すると、Tab キーを押して、最初のコントロールにフォーカスし、スクリーン・リーダーを起動する必要があります。特殊なキーボード・ショートカットおよびスクリーン・リーダーのショートカットの詳細について、次の表で説明しています。

表 14. ダイアログ・ボックスのショートカット

ショートカット・キー	関数
タブ	ウィンドウ、またはダイアログ・ボックスの次の項目に移動します。
Ctrl+Tab	テキスト・ボックスから次の項目に移動します。
Shift+Tab	ウィンドウ、またはダイアログ・ボックスの前の項目に移動します。
Shift + Ctrl + Tab	テキスト・ボックスから前の項目に移動します。
スペース キー	フォーカスのあるコントロールまたはボタンを選択します。
Esc	変更をキャンセルして、ダイアログ・ボックスを閉じます。
Enter	変更を確認して、ダイアログ・ボックスを閉じます (「OK」ボタンと同じ)。テキスト・ボックスで作業している場合、まず Ctrl + Tab を押して、テキスト・ボックスから移動する必要があります。

第 8 章 コンセプトとタイプの抽出

インタラクティブ・ワークベンチを起動するストリームを実行する場合、抽出はストリームのテキスト・データに実行されます。この抽出の最終結果は、一連のコンセプト、タイプ、そして TLA パターンが言語リソースにある場合はパターンとなります。抽出結果ペインでコンセプトおよびタイプを表示および処理できます。詳しくは、5 ページの『抽出の方法』のトピックを参照してください。

抽出結果を調整する場合、言語リソースを変更し、再抽出できます。詳しくは、89 ページの『抽出結果の調整』のトピックを参照してください。抽出プロセスは、結果の抽出および構成方法を指定する「抽出」ダイアログ・ボックスのリソースおよびパラメーターによって異なります。抽出結果を使用して、カテゴリ一定義の (全部ではない場合) 大部分を定義できます。

抽出結果: コンセプトとタイプ

抽出プロセスで、すべてのテキスト・データがスキャンされ、関連するコンセプトが特定、抽出、そしてタイプに割り当てられます。抽出が完了すると、カテゴリとコンセプト・ビューの左下隅にある抽出結果ペインに結果が表示されます。セッションを初めて起動した場合、ノードで選択した言語リソース・テンプレートを使用して、これらのコンセプトおよびタイプを抽出および構成します。

注: 表示中のペインに結果を表示しきれない場合は、ペインの下部にあるコントロールを使用して前後の結果に移動したり、移動先のページ番号を入力したりすることができます。

抽出されたコンセプト、タイプ、および TLA パターンは、まとめて抽出結果と呼ばれ、カテゴリの記述子、または構築ブロックとして機能します。また、カテゴリ規則でコンセプト、タイプ、およびパターンを使用することもできます。さらに、自動的手法では、コンセプトおよびタイプを使用してカテゴリを作成します。

テキスト・マイニングは、抽出結果をテキスト・データのコンテキストに従ってレビューし、新しい結果を作成するよう調整、そして再評価するインタラクティブ プロセスです。抽出後、結果を表示し、必要に応じて言語リソースを修正することによって、結果を変更する必要があります。抽出結果ペイン、データ・ペイン、「カテゴリ一定義」ダイアログ・ボックス、または「クラスター一定義」ダイアログ・ボックスから直接、リソースをある程度調整できます。詳しくは、89 ページの『抽出結果の調整』のトピックを参照してください。リソース・エディター ビューで直接調整することもできます。詳しくは、74 ページの『リソース・エディター・ビュー』のトピックを参照してください。

調整した後、再抽出して新しい結果を表示できます。最初から抽出結果を調整することによって、再抽出するごとに、カテゴリ一定義で同じ結果を取得することを確認し、データのコンテキストに完全に対応することができます。このようにして、ドキュメント/レコードをより正確で、繰り返し可能な方法で、カテゴリ一定義に割り当てます。

コンセプト

抽出プロセスで、テキスト・データをスキャンして分析し、テキスト内の関心のあるまたは関連する 1 つの単語 (election または peace など) や句 (presidential election, election of the president、または peace treaties など) を特定します。これらの単語や句を、まとめて「キーワード」と呼びます。言語リソースを使用して、関連キーワードを抽出し、類似したキーワードをコンセプトと呼ばれる代表語でグループ化します。

コンセプト名の上にマウスポインタを置くと、コンセプトの基本キーワードのセットが表示されます。これにより、コンセプト名とそのコンセプトにグループ化された数行のキーワードを示すツールヒントが表示されます。これらの基本キーワードには、抽出された複数形/単数形のキーワード、置換キーワード、Fuzzy Grouping のキーワードなどのほか、テキスト内にあったかどうかに関係なく、言語リソースで定義された類義語が含まれます。コンセプト名を右クリックし、コンテキスト・メニューのオプションを選択して、これらのキーワードをコピーしたり、基本キーワードの完全セットを表示したりできます。

デフォルトでは、ドキュメント数 (「ドキュメント」列) に従って、降順で並べられます。コンセプトが抽出されると、それらをタイプに割り当てて、同様のコンセプトをグループ化します。それらはこのタイプに従ってカラー・コード化されます。色は、リソース・エディター のタイプのプロパティで定義されます。詳しくは、189 ページの『キーワード辞書』のトピックを参照してください。

カテゴリ定義でコンセプト、タイプ、またはパターンが使用されている場合、アイコンが並べ替え可能な「投入」列に表示されます。。

データ型

タイプは、コンセプトの意味上のグループ化です。コンセプトが抽出されると、それらをタイプに割り当てて、同様のコンセプトをグループ化します。<Location>、<Organization>、<Person>、<Positive>、<Negative> など、いくつかのビルトインのタイプが IBM SPSS Modeler Text Analytics に付属しています。例えば、<Location> のタイプは、地理的なキーワードや地名をグループ化します。このタイプは、chicago、paris、および tokyo などのコンセプトに割り当てられます。多くの言語の場合、キーワード辞書にはないが、テキストから抽出されたコンセプトは、自動的に <Unknown> のタイプとなります。詳しくは、190 ページの『ビルトインのタイプ』のトピックを参照してください。

タイプ ビューを選択すると、デフォルトでは抽出したタイプがグローバルな頻度の高い順に表示されます。タイプは区別できるように、カラー・コード化されます。色は、タイプのプロパティの一部です。詳しくは、191 ページの『キーワード辞書の作成』のトピックを参照してください。独自のタイプを作成することもできます。

パターン

テキスト・データからパターンを抽出することもできます。ただし、リソース・エディター にテキスト リンク分析 (TLA) パターン規則を含むライブラリーが必要です。「テキスト リンク分析のパターン抽出を有効にする」 オプションを使用して、IBM SPSS Modeler Text Analytics ノードの設定または「抽出」ダイアログ・ボックスでこれらのパターンの抽出を選択する必要があります。詳しくは、149 ページの『第 11 章 テキスト リンク分析の検証』のトピックを参照してください。

データの抽出

抽出が必要な場合、抽出結果ペインが黄色で表示され、メッセージ「抽出ボタンをクリックしてキーワードを抽出してください」というメッセージが、このウィンドウ枠のツールバーの下に表示されます。

抽出結果がない場合、言語リソースに変更を行い抽出結果を更新する必要がない場合、または抽出結果を保存していない セッションを開く場合は、抽出が必要な場合があります (「ツール」 > 「オプション」)。

注: 抽出結果を「セッション作業を使用」オプションを使用してキャッシュした後、ストリームのソースノードを変更する際、抽出結果を更新する場合にインタラクティブ ワークベンチ セッションが起動したら、新しい抽出を実行する必要があります。

抽出実行中には進行状況が表示されます。抽出している間、抽出エンジンはテキスト・データをすべて読み込み、関連キーワードおよびパターンを特定し、それらを抽出して、タイプに割り当てます。そして、エンジンは、1つの主要なキーワード、コンセプトに類義語のキーワードをグループ化します。プロセスが完了すると、生成されたコンセプト、タイプ、パターンが抽出結果ペインに表示されます。

抽出プロセスにより、一連のコンセプト、タイプ、そして有効な場合はテキスト リンク分析 (TLA) パターンが作成されます。カテゴリとコンセプト・ビューの抽出結果ペインでこれらのコンセプトおよびタイプを表示および処理できます。TLA パターンを抽出した場合、これらはテキスト リンク分析ビューにされます。

注: データセットのサイズと、抽出プロセスを完了するためにかかる時間の間には、関連性があります。上流にサンプル・ノードを追加、またはコンピューターの構成を最適化することをいつでも検討することができます。

データを抽出するには

1. メニューの「ツール」 > 「抽出」を選択します。または、「抽出」 ツールバー・ボタンをクリックします。
2. 「抽出設定」ダイアログの表示を選択すると必ず、ダイアログが表示され、変更を行うことができます。各設定の記述子については、このトピックの後半を参照してください。
3. 「抽出」 をクリックして、抽出プロセスを開始します。抽出が始まると、進捗状況のダイアログ・ボックスが表示されます。抽出後、結果が「抽出結果」ウィンドウに表示されます。デフォルトでは、ドキュメント数 (「ドキュメント」列) に従って、降順で並べられます。

ツールバー・オプションを使用して結果を確認し、結果を並べ替える、結果を絞り込む、または異なるビュー (コンセプト、またはタイプ) に切り替えることができます。言語リソースを処理して、抽出結果を調整することもできます。詳しくは、89 ページの『抽出結果の調整』のトピックを参照してください。

抽出時に発生する可能性がある問題

複数のインタラクティブ ワークベンチ セッションを使用すると、動作が遅くなる可能性があります。SPSS Modeler Text Analytics および SPSS Modeler は、インタラクティブ ワークベンチ セッションを起動するとき共通の Java ランタイム エンジンと共有します。SPSS Modeler セッション中に起動するインタラクティブ ワークベンチ セッションの数によっては、同じセッションを開いた後に閉じる場合であっても、システム メモリーのためにアプリケーションの動作が遅くなる可能性があります。大量のデータを処理する場合や、推奨 RAM 設定 (4 GB) に満たないマシンを使用する場合は、この影響が特に顕著になることがあります。マシンの応答が遅いと感じる場合は、すべての作業を保存して SPSS Modeler をシャットダウンし、アプリケーションを再起動することをお勧めします。推奨メモリー未満のマシンで SPSS Modeler Text Analytics を実行すると、特に大規模なデータ・セットを処理する場合や長時間にわたって作業する場合に、Java のメモリーが不足してシャットダウンしてしまうことがあります。大量のデータを処理する場合は、推奨メモリー設定以上にアップグレードする (または SPSS Modeler Text Analytics Server を使用する) ことを強くお勧めします。

オランダ語、英語、フランス語、ドイツ語、イタリア語、ポルトガル語、スペイン語のテキストの場合

「抽出設定」ダイアログ・ボックスには、基本的な抽出オプションがいくつか表示されます。

テキスト リンク分析のパターン抽出を有効にする: テキスト・データから TLA パターンを抽出するように指定します。また、リソース・エディターのいずれかのライブラリーに TLA パターン規則があることも想定します。このオプションを指定すると、抽出時間が大幅に長くなります。詳しくは、149 ページの『第 11 章 テキスト リンク分析の検証』のトピックを参照してください。

句読点エラーを調整: 抽出時に句読点エラー (不適切な使用方法など) を含むテキストを一時的に正規化し、コンセプトの抽出可能性を向上させます。自由記述式アンケートの回答、電子メール、CRM データなど、テキストが短く品質が悪い場合、またはテキストに略語が多く含まれている場合に特に役立ちます。

文字数が次の最小値以上のときにスペルを調整する: Fuzzy Grouping の手法を適用し、共通してミススペルのある単語またはスペルの近い単語を 1 つのコンセプトにグループ化できるようにします。Fuzzy Grouping アルゴリズムでは、最初の母音を除くすべての母音を一時的に抜き取った後抽出した単語から 2 つ/3 つの子音を抜き取り、それらと比較して、それらが同じで modeling と modelling が同じグループに分けられるかどうかを確認します。ただし、各キーワードが <Unknown> タイプを除いて、別のタイプに割り当てられた場合、Fuzzy Grouping 手法は適用されません。

Fuzzy Grouping を使用する前に必要な、語幹文字数の制限を定義することもできます。キーワード内の語幹文字数は、すべての文字を合計し、活用語尾、複合語キーワードの場合は区切り文字および前置詞を形成する文字を差し引いて計算します。例えば、キーワード exercises の語幹文字数は「exercise」という形式で 8 文字と数えられます。語末の s は活用語尾 (複数形) であるためです。同様に、apple sauce の語幹文字は 10 文字 (「apple sauce」)、そして manufacturing of cars の語幹文字は 16 文字 (「manufacturing car」) となります。この算出方法は、Fuzzy Grouping を適用するべきかどうかを確認するためにのみ使用されますが、単語がどのように一致するかについては影響を与えません。

注: 特定の単語が後で不適切にグループ化されていることが分かった場合、「拡張リソース」タブの **Fuzzy Grouping**: 例外 セクションで 明示的に宣言することによって、単語のペアをこの手法から除外できます。詳しくは、205 ページの『Fuzzy Grouping』のトピックを参照してください。

ユニタームを抽出 単語が複合語の一部でない限り、または名詞、またはスピーチ内の認識できない品詞である場合、このオプションは単一の単語 (ユニターム) を抽出します。

固有表現を抽出 電話番号、セキュリティー番号、時間、日付、通貨、数字、パーセント、電子メールアドレス、HTTP アドレスなどの固有表現を抽出します。「拡張リソース」タブの「固有表現: 設定」セクションで、特定の種類の固有表現を追加したり除外したりできます。不要な固有表現を無効にすることにより、抽出エンジンは処理時間を節約できます。詳しくは、210 ページの『構成』のトピックを参照してください。

大文字アルゴリズム キーワードの最初の文字が大文字である場合、組み込み辞書にない単純キーワードおよび複合キーワードを抽出します。このオプションには、最も適切な名詞を抽出するのに優れた方法があります。

可能な場合は、個人名の一部または全部をグループ化 テキスト内で別々の形式で同時に出現する名前をグループ化します。名前はテキストの始めでは完全な形式で、後は短い形式でのみ参照されるため、この機能が役立ちます。このオプションでは、タイプが <Unknown> のユニタームが、タイプ <Person> の複合キーワードの最後の単語に一致するようにします。例えば、doe があり、最初タイプが <Unknown> である場合、抽出エンジンは、<Person> タイプの複合キーワードに最後の単語として doe が含まれているかどうか (例: john doe) を確認します。ほとんどがユニタームとして抽出されることがないため、人の名前に適用されることはありません。

機能語による倒置を次の値を最大値として考慮する 倒置手法を適用する場合に指定されている場合がある非機能的単語の最大数を指定します。この倒置手法では、活用語尾に関係なく、含まれる非機能的単語 (of や the など) によってお互いに異なる類似した句をグループ化します。例えば、この値を最大 2 単語に設定し、company officials および officials of the company が抽出されたとします。この場合、両方の抽出キーワードは、of the が無視されると同じであるとみなされるため、最終コンセプト・リストに共にグループ化されます。

マルチタームをグループ化するときに派生関係を使用: ビッグデータを処理するときにこのオプションを選択すると、派生規則を使用してマルチタームがグループ化されます。

コンセプト・マップのインデックス・オプション コンセプト・マップを後ですぐに描画できるように、抽出時間にマップの指標を作成することを指定します。インデックスの設定を編集するには、「設定」をクリックします。詳しくは、89 ページの『コンセプト・マップ・インデックスの作成』のトピックを参照してください。

抽出前に常にこのダイアログ・ボックスを表示する: 「ツール」メニューを選択しない限り表示したくない場合、抽出ごとに「抽出設定」ダイアログを表示するかどうか、または抽出設定を編集する場合、抽出ごとに表示するかどうかを尋ねるかどうかを指定します。

抽出結果のフィルタリング

非常に大きなデータセットを処理する場合、抽出プロセスでは、多数の結果が作成される場合があります。多くのユーザーによって、多数の結果が作成されると、結果を効率的に確認することが困難になります。そのため、最も関心のある結果に絞込むために、抽出結果ペインで使用できる「フィルター」ダイアログを使用してこれらの結果をフィルタリングできます。

「フィルター」ダイアログのすべての設定を同時に使用して、カテゴリーに使用できる抽出結果をフィルタリングします。

出現頻度でフィルタリング フィルタリングを実行して、特定のグローバル出現頻度値またはドキュメントの出現頻度の値を持つ結果のみを表示できます。

- グローバル出現頻度は、コンセプトがドキュメントまたはレコードの全体的なセットに出現する回数の合計で、「グローバル」列に表示されます。
- ドキュメント出現頻度は、コンセプトが出現するドキュメントまたはレコードの合計数で、「ドキュメント」列に表示されます。

例えば、コンセプト `nato` が 500 件のレコードに 800 回出現した場合、このコンセプトのグローバル出現頻度は 800 で、ドキュメント出現頻度は 500 となります。

タイプ 特定のタイプに属する結果のみを表示できます。すべてのタイプまたは特定のタイプのみを選択できます。

マッチ テキスト別 **AND** 条件 ここで定義する規則に一致する結果のみを表示できます。「マッチ テキスト」フィールドに一致する文字のセットを入力し、マッチに適用する条件を選択します。

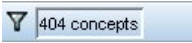


表 15. マッチ・テキストの条件

条件	説明
以下を含む:	文字列が任意の場所で出現する場合、テキストが一致します(デフォルトの選択)。
以下で始まる:	コンセプトまたはタイプが特定のテキストで始まる場合にのみ、テキストが一致します。
以下で終わる:	コンセプトまたはタイプが特定のテキストで終わる場合にのみ、テキストが一致します。
完全一致	文字列全体が、コンセプト名またはタイプ名に一致する必要があります。

抽出結果ペインに表示される結果

フィルタリングに基づいて、結果が「抽出結果」ウィンドウに英語でどのように表示されるかについて、いくつかの例を示します。

表 16. フィルター・フィードバックの例

フィルター・フィードバック	説明
	ツールバーには結果の数が表示されます。テキスト・マッチ フィルターがなく、最大数に達していないため、追加のアイコンは表示されません。
	ツールバーには、結果がフィルターで指定された最大値に制限されていることを示します (この例では 300)。紫のアイコンが表示されている場合、コンセプトの最大数に達したことを示します。アイコンの上にポインタを置くと、詳細が表示されます。
	ツールバーには、マッチ・テキスト・フィルターを使用して、結果が制限されていることを示します。虫めがねのアイコンが表示されます。

結果を絞り込むには

1. メニューの「ツール」 > 「フィルター」を選択します。「フィルター」ダイアログ・ボックスが開きます。
2. 使用するフィルターを選択および調整します。
3. 「OK」 をクリックして、フィルターを適用すると、「抽出結果」ウィンドウに新しい結果が表示されます。

コンセプト・マップの検証

コンセプト・マップを作成して、コンセプトがどのように関連するかを検証できます。1 つのコンセプトを選択し、「マップ」 をクリックすると、コンセプト・マップのウィンドウが開き、選択したコンセプトに関連するコンセプトのセットを検証することができます。含めるタイプ、検索する関係性の種類など、設定を編集して表示するコンセプトを除外できます。

重要: マップを作成する前に、インデックスを生成する必要があります。これには数分かかることがあります。ただし、いったんインデックスを生成すると、再抽出するまで指標を再生成する必要がありません。抽出するごとにインデックスを自動的に生成したい場合は、抽出設定でそのオプションを選択します。詳しくは、82 ページの『データの抽出』のトピックを参照してください。

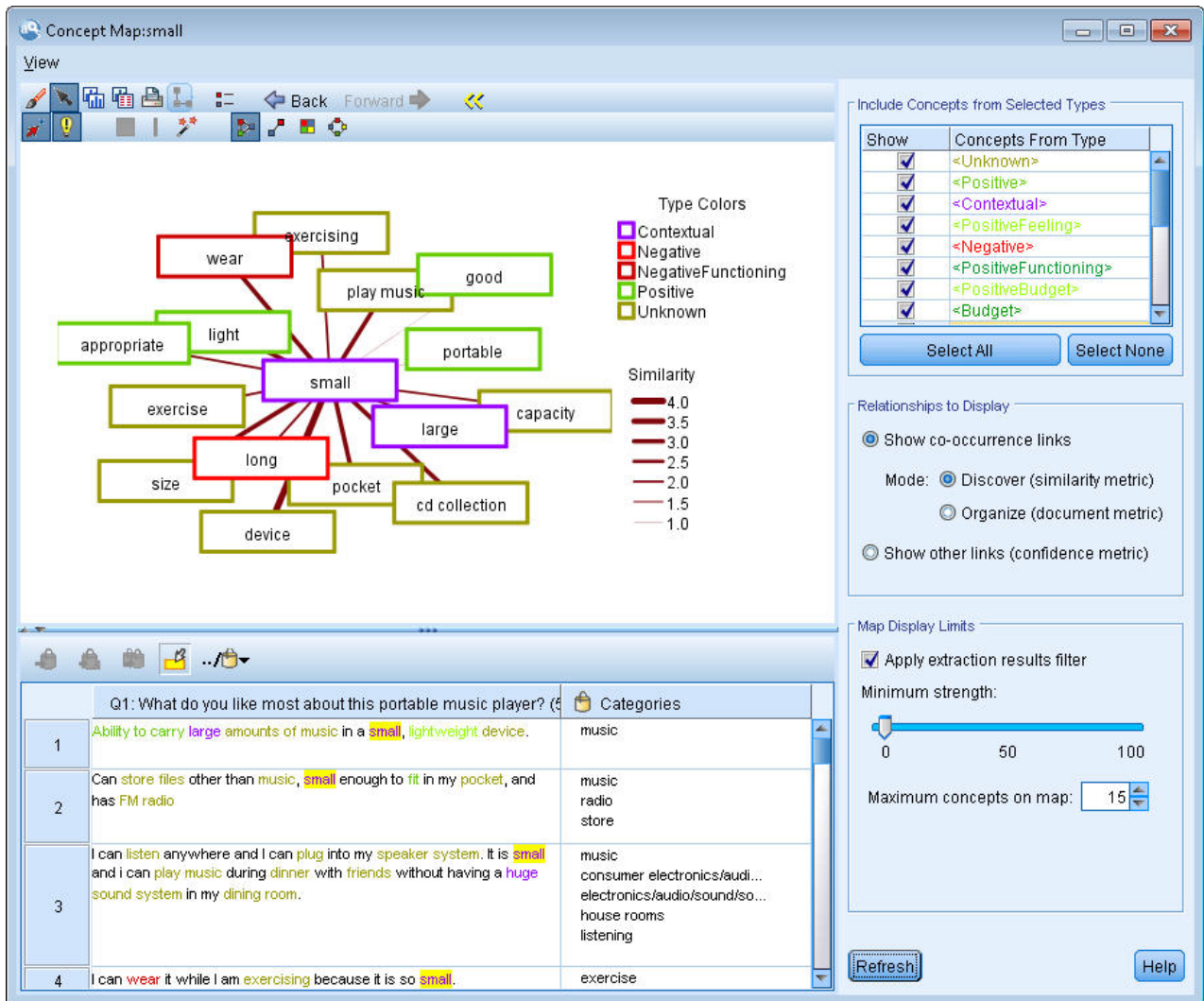


図 27. 選択したコンセプトのコンセプト・マップ

コンセプト・マップを表示するには

1. 「抽出結果」ウィンドウで、1つのコンセプトを選択します。
2. このパネルのツールバーで、「マップ」ボタンをクリックします。マップ・インデックスがすでに生成されている場合、コンセプト・マップが個別のダイアログで開きます。マップ・インデックスが生成されていない、または古い場合、インデックスを再度作成する必要があります。このプロセスには数分かかることがあります。
3. 検証するマップの周辺をクリックします。リンクしたコンセプトをダブルクリックすると、マップが再描画再描画され、ダブルクリックしたコンセプトのリンクしたコンセプトが表示されます。
4. 最上部のツールバーには、以前のマップに戻す、関係の強度に応じてリンクをフィルタリング、出現するコンセプトのタイプや表示する関係の種類を制御するフィルター・ダイアログを開くなど、基本的なマップツールがいくつかあります。2番目のツールバーのラインには、グラフ編集ツールがあります。詳しくは、159ページの『グラフのツールバーおよびパレットの使用』のトピックを参照してください。
5. 検出されるリンクの種類が適切でない場合、マップに右側に表示されたこのマップの設定を確認してください。

マップの設定: 選択したタイプのコンセプトを追加

テーブルの選択されたタイプに属するこれらのコンセプトのみがマップに表示されます。特定のタイプのコンセプトを隠すには、テーブルの該当するタイプの選択を解除します。

マップの設定: 表示すべき関連性

共起リンクを表示: 共起リンクを表示するには、モードを選択します。モードは、リンクの強度がどのように計算されたかに影響を与えます。

- 探索 (類似性メトリック): 類似性メトリックで、2 つのコンセプトが個別に出現する頻度と、同時に出現する頻度を考慮に入れた複雑な計算方法で、リンクの強度を算出します。高い強度の値は、コンセプトのペアは、個別に出現するよりも、同時に出現する頻度が高いことを示します。次の式により、浮動小数点値は整数に変換されます。

$$\text{similarity coefficient} = \frac{(C_{IJ})^2}{(C_I \times C_J)}$$

図 28. 類似度係数の式

この式で、 C_I は、コンセプト I が出現するドキュメントまたはレコードの数です。

C_J は、コンセプト J が出現するドキュメントまたはレコードの数です。

C_{IJ} は、コンセプトのペア I および J が出現するドキュメントまたはレコードの数です。

- 構成 (ドキュメント・メトリック): ドキュメント・メトリックを持つリンクの強度は、共起の未調整度数で決定します。一般的に、2 つのコンセプトがより頻繁に出現すると、同時に出現する確率が高くなります。強度の値が高いと、コンセプトのペアが頻繁に同時に出現します。

他のリンクを表示 (確信度メトリック): 他のリンクの表示を選択することもできます。例えば、セマンティック、派生 (形態)、または、内包 (シンタックス) であり、リンクされたコンセプトからいくつのステップのコンセプトが削除されたかに関連します。このことは、リソースの調整、特に類義語 や曖昧さの回避に役立ちます。このようなグループ化の手法の簡単な説明は、108 ページの『言語学的手法の詳細設定』を参照してください。

注: インデックスを作成するときにこれらのオプションが選択されなかった場合、または関連性が見つからなかった場合、何も表示されません。詳しくは、89 ページの『コンセプト・マップ・インデックスの作成』のトピックを参照してください。

マップの設定: マップ表示制限

抽出結果フィルターの適用: すべてのコンセプトを使用する場合、「抽出結果」ウィンドウでフィルターを使用して、表示内容を制限することができます。このオプションを選択すると、IBM SPSS Modeler Text Analytics が、フィルタリングされたセットを使用して、関連コンセプトを検索します。詳しくは、85 ページの『抽出結果のフィルタリング』のトピックを参照してください。

最小強度: ここで、最小強度を設定します。関連性の強度がこの制限値より低い関連コンセプトは、マップに表示されません。

マップ上の最大コンセプト数: マップに表示する関連性の最大数を指定します。

コンセプト・マップ・インデックスの作成

マップを作成する前に、コンセプトの関連性のインデックスを生成する必要があります。コンセプト・マップを作成する場合、IBM SPSS Modeler Text Analytics は、このインデックスを参照します。このダイアログで手法を選択して、インデックスへの関連性を選択できます。

グループ化手法。1 つまたは複数の手法を選択します。これらの手法の簡単な説明は、「110 ページの『言語学的手法について』」を参照してください。すべてのテキスト言語ですべての手法が使用できるわけではありません。

特定のコンセプトがグループ化されないようにする: 出力の 2 つのコンセプトがグループ化またはペアとないように処理を停止します。コンセプト・ペアを作成または管理するには、「ペアを管理」をクリックします。詳しくは、110 ページの『例外ペアのリンクの管理』のトピックを参照してください。

インデックスの作成には、数分かかる場合があります。ただし、いったんインデックスを生成すると、再抽出するまで、または設定を変更してより多くの関連性を追加しない限り、インデックスを再生成する必要はありません。抽出するごとにインデックスを生成したい場合は、抽出設定でそのオプションを選択します。詳しくは、82 ページの『データの抽出』のトピックを参照してください。

抽出結果の調整

抽出とは反復可能なプロセスで、結果を抽出、確認、変更、および再抽出して結果を更新できます。正常なテキストマイニングおよびカテゴリー化には精度および継続性が不可欠であるため、最初から抽出結果を調整することによって、再抽出ごとに、カテゴリー定義でまったく同じ結果が得られます。このようにして、レコードおよびドキュメントをより正確で、繰り返し可能な方法で、カテゴリーに割り当てます。

抽出結果は、カテゴリーを作成するための要素となります。これらの抽出結果を使用してカテゴリーを作成すると、1 つまたは複数のカテゴリー記述子に一致するテキストが含まれる場合、レコードおよびドキュメントが自動的にカテゴリーに割り当てられます。言語リソースを調整する前にカテゴリー化を開始できますが、開始前に少なくとも 1 回、抽出結果を確認しておくが役立ちます。

結果を確認すると、抽出エンジンが異なる方法で処理する必要のある要素が見つかる場合があります。以下のような例があります。

- **認識されない類義語。** 賢い、知的、頭脳明晰、博識など、類義語と考えられるいくつかのコンセプトが見つかり、抽出結果に個別のコンセプトとしてすべて表示されたとします。知的、頭脳明晰、博識がすべて代表コンセプト賢いの名ですべてグループ化されるよう、類義語定義を作成できます。こうすることにより、これらのコンセプトをすべて賢いとグループ化し、グローバル出現頻度も高くなります。詳しくは、90 ページの『類義語の追加』のトピックを参照してください。
- **ミスタイプ・コンセプト:** 抽出結果のコンセプトがあるタイプに出現し、別のタイプに割り当てたい場合があります。また、抽出結果に 15 種類の野菜のコンセプトがあり、それらすべてを <Vegetable> という新しいタイプに追加したい場合もあります。多くの言語の場合、キーワード辞書にはないが、テキストから抽出されたコンセプトは、自動的に <Unknown> のタイプとなります。コンセプトをタイプに追加できます。詳しくは、91 ページの『コンセプトのタイプへの追加』のトピックを参照してください。
- **重要でないコンセプト:** 抽出されたコンセプトで非常に頻度の高い、つまり多くのレコードまたはドキュメントで見つかる場合があります。ただし、このコンセプトは分析には重要でないと見なされます。このコンセプトを抽出から除外できます。詳しくは、92 ページの『コンセプトの抽出からの除外』のトピックを参照してください。

- **不正な合致:** 特定のコンセプトを含むレコードまたはドキュメントを確認する場合、faculty (能力) と facility (施設) のように 2 つの単語が誤ってグループ化されているのを発見する場合があります。この合致は Fuzzy Grouping という、内部アルゴリズムによるものであり、2 つまたは 3 つの子音および母音を一時的に無視して、一般的なスペルミスグループ化します。これらの単語を無視する必要のある単語のペアのリストに追加できます。詳しくは、205 ページの『Fuzzy Grouping』のトピックを参照してください。
- **未抽出コンセプト:** 特定のコンセプトが抽出されるのを期待しているにもかかわらず、レコードまたはドキュメント テキストを確認しているときに一部の単語または句が抽出されていないことに気づく場合があります。これらの単語は重要でない動詞または形容詞である場合が多くあります。ただし、抽出されなかった単語または句をカテゴリ定義として使用したい場合があります。コンセプトを抽出するために、キーワードをキーワード辞書に強制投入できます。詳しくは、93 ページの『単語を抽出に強制投入』のトピックを参照してください。

こうした変更の多くは、1 つまたは複数の要素を選択して右クリックし、コンテキスト・メニューを使用することによって、抽出結果ペイン、データ・ペイン、「カテゴリ定義」ダイアログ・ボックス、または「クラスター定義」ダイアログ・ボックスから直接実行できます。

変更を行った後、ペインの背景色が変わり、変更を表示するには再抽出が必要であることを示します。詳しくは、82 ページの『データの抽出』のトピックを参照してください。大きなデータ・セットを使用している場合は、1 つ変更するたびにキーワードを再抽出するのではなく、いくつかの変更を加えてから再抽出したほうが効率的です。

注: リソース・エディター・ビューで抽出結果を作成するために使用する、編集可能な言語リソースのセット全体を表示できます (「表示」>「リソース・エディター」)。これらのリソースは、このビューにライブラリーおよび辞書の形式で表示されます。ライブラリーおよび辞書内で直接コンセプトおよびタイプをカスタマイズできます。詳しくは、179 ページの『第 15 章 ライブラリーの使用』のトピックを参照してください。

類義語の追加

類義語は、同じ意味を持つ 2 つ以上の単語と関連があります。類義語はキーワードとその短縮形をまとめるのにもよく使用されます。またよくつづり間違いが起こる語を正しい書き方のもので置き換えるのにも使用されます。類義語を使用すると、代表コンセプトの頻度が高くなり、テキスト・データ内のさまざまな方法で表示される類似した情報を見つけやすくなります。

製品に付属する言語リソース・テンプレートおよびライブラリーには、事前定義された多くの類義語が含まれています。ただし、認識されていない類義語が見つかった場合、次回抽出するときに認識されるよう、類義語を定義することができます。

まず、代表コンセプトまたは主要キーワードを決定します。代表コンセプトは、最終的な結果ですべての類義語のキーワードをグループ化したい単語または句です。抽出時、類義語は、この代表コンセプトの下でグループ化されます。次に、このコンセプトのすべての類義語を特定します。代表コンセプトは、最終的な抽出で、すべての類義語と置き換えられます。類義語となるためにはキーワードは抽出されなければなりません。ただし、代表コンセプトを、置換えを行うために抽出する必要はありません。例えば、知的という単語を賢いという単語に置き換えたい場合、知的が類義語となり、賢いが代表コンセプトとなります。

新しい類義語定義を作成する場合、新しい代表コンセプトが辞書に追加されます。その後、類義語をその代表コンセプトに追加する必要があります。類義語を作成または編集する場合、これらの変更が リソース・エディター の類義語辞書に記録されます。これらの類義語辞書の内容全体を表示したい場合、またはかな

りの数の変更を行いたい場合、リソース・エディター で直接作業することが必要な場合があります。詳しくは、197 ページの『類義語辞書』のトピックを参照してください。

新しい類義語は、リソース・エディター ビューのライブラリー・ツリーの表示された最初のライブラリーに自動的に保存されます。デフォルトでは、これがローカル・ライブラリーになります。

注: 類義語定義を探していて、コンテキスト メニューがなく、リソース・エディターで直接見つからなくても、内部の Fuzzy Grouping 手法によって、合致が発生している場合があります。詳しくは、205 ページの『Fuzzy Grouping』のトピックを参照してください。

新しい類義語を作成するには

1. 抽出結果ペイン、データ・ペイン、「カテゴリー定義」ダイアログ・ボックス、または「クラスター定義」ダイアログ・ボックスで、新しい類義語を作成したいコンセプトを選択します。
2. メニューから「編集」>「類義語に追加」>「新規」を選択します。「類義語の作成」ダイアログ・ボックスが開きます。
3. 「代表語」テキスト・ボックスに、代表語を入力します。これが、すべての類義語がグループ化されるコンセプトとなります。
4. さらに類義語を追加したい場合、「類義語」リスト・ボックスにそれらを入力します。辞書エディターの区切り文字（デフォルトでは「,」（コンマ））を使用して、各類義語を分けます。詳しくは、76 ページの『オプション:「セッション」タブ』のトピックを参照してください。
5. 「OK」をクリックし、変更を適用します。ダイアログ・ボックスが閉じて、抽出結果ペインの背景色が変わり、変更を表示するには再抽出が必要であることを示します。変更が複数ある場合には、再抽出する前に変更を終わらせてください。

類義語を追加するには

1. 抽出結果ペイン、データ・ペイン、「カテゴリー定義」ダイアログ・ボックス、または「クラスター定義」ダイアログ・ボックスで、既存の類義語定義に追加したいコンセプトを選択します。
2. メニューから「編集」>「類義語に追加」を選択します。メニューには、一番最後に作成された代表語がリストの一番上に表示されます。選択したコンセプトを追加したい類義語の名前を選択します。探している類義語が表示されたら、その名前を選択します。選択したコンセプトがその類義語の定義に追加されます。メニューに類義語が表示されない場合、「もっと表示」を選択すると「すべての類義語」ダイアログ・ボックスが表示されます。
3. 「すべての類義語」ダイアログ・ボックスで、リストを自然な並び順（作成順）の昇順または降順で並べ替えることができます。選択したコンセプトを追加したい類義語の名前を選択し、「OK」をクリックします。ダイアログ・ボックスが閉じ、コンセプトが類義語の定義に追加されます。

コンセプトのタイプへの追加

抽出を実行している場合、共通点を持つキーワードをグループ化するために、抽出されたコンセプトがタイプに割り当てられます。IBM SPSS Modeler Text Analytics は、多くのビルトインのタイプに付属しています。詳しくは、190 ページの『ビルトインのタイプ』のトピックを参照してください。多くの言語の場合、キーワード辞書にはないが、テキストから抽出されたコンセプトは、自動的に <Unknown> のタイプとなります。

結果を確認しているとき、あるタイプに割り当てたいコンセプトがいくつか異なるタイプに割り当てられていたり、単語のグループが実際は新しいタイプに割り当てられていたりすることが分かる場合があります。こうした場合、コンセプトを別のタイプに割り当てなおすか、まとめて新しいタイプを作成する場合があります。

例えば、自動車に関連する調査データを扱っており、車のさまざまな領域に焦点を当ててのカテゴリー化に関心があるとします。<Dashboard> というタイプを作成し、車のダッシュボードにある計測器およびノブに関連するすべてのコンセプトをグループ化することができます。そして、gas gauge、gas gauge、radio、および odometer などのコンセプトを、この新しいタイプに割り当てることができます。

また例えば、大学に関連する調査データ、そして <Organization> ではなく <Person> というタイプとしての抽出 Johns Hopkins (大学) を扱っているとします。この場合、このコンセプトを <Organization> タイプに追加することができます。

タイプを作成、またはコンセプトをタイプのキーワード・リストに追加する場合、これらの変更は リソース・エディター の言語リソース・ライブラリーのキーワード辞書に記録されます。これらのライブラリーの内容を表示したい場合、またはかなりの数の変更を行いたい場合、リソース・エディター で直接作業することが必要な場合があります。詳しくは、192 ページの『キーワードの追加』のトピックを参照してください。

コンセプトをタイプに追加するには

1. 抽出結果ペイン、データ・ペイン、「カテゴリー定義」ダイアログ・ボックス、または「クラスター定義」ダイアログ・ボックスで、既存のタイプに追加したいコンセプトを選択します。
2. 右クリックしてコンテキスト・メニューを開きます。
3. メニューから「編集」>「タイプに追加」を選択します。タイプ名が見つかったなら、これを選択します。選択したコンセプトを追加したいタイプの名前を選択します。探しているタイプの名前が表示されたら、その名前を選択します。選択したコンセプトがそのタイプに追加されます。メニューに類義語が表示されない場合、「もっと表示」を選択すると「すべてのタイプ」ダイアログ・ボックスが表示されます。
4. 「すべてのタイプ」ダイアログ・ボックスで、リストを自然な並び順 (作成順) の昇順または降順で並べ替えることができます。選択したコンセプトを追加したいタイプの名前を選択し、「OK」をクリックします。ダイアログ・ボックスが閉じ、コンセプトがタイプに追加されます。

新規タイプを作成するには

1. 抽出結果ペイン、データ・ペイン、「カテゴリー定義」ダイアログ・ボックス、または「クラスター定義」ダイアログ・ボックスで、新しいタイプを作成したいコンセプトを選択します。
2. メニューから「編集」>「タイプに追加」>「新規」を選択します。「タイプのプロパティ」ダイアログ・ボックスが開きます。
3. 「名前」テキスト・ボックスにこのタイプの新しい名前を入力し、他のフィールドを変更します。詳しくは、191 ページの『キーワード辞書の作成』のトピックを参照してください。
4. 「OK」をクリックし、変更を適用します。ダイアログ・ボックスが閉じて、抽出結果ペインの背景色が変わり、変更を表示するには再抽出が必要であることを示します。変更が複数ある場合には、再抽出する前に変更を終わらせてください。

コンセプトの抽出からの除外

結果を確認しているときに、自動カテゴリー作成手法によって必要のないコンセプトが抽出または使用されていることが分かることがあります。これらのコンセプトの頻度が非常に高く、分析には全く重要でない場合もあります。この場合、コンセプトを最終的な抽出結果から除外するようにマークすることができます。通常、このリストに追加するコンセプトは、テキストの中で穴埋めとして使われる単語または句で、特に重要な意味を付け加えるようなものではなく、抽出結果を混乱させる場合があります。コンセプトを不要語辞書に追加しておけば、それらが抽出されないようにすることができます。

コンセプトを除外すると、次回抽出するときには除外したコンセプトはそのすべての変化形も含め、抽出結果から消失します。このコンセプトがカテゴリーの記述子として既に出現している場合、再抽出後は 0 カウントでカテゴリー内に残ります。

除外すると、これらの変更が リソース・エディター の不要語辞書に記録されます。除外する定義をすべて表示して直接編集したい場合は、リソース・エディター で直接作業することを推奨します。詳しくは、200 ページの『不要語辞書』のトピックを参照してください。

コンセプトを除外するには

1. 抽出結果ペイン、データ・ペイン、「カテゴリー定義」ダイアログ・ボックス、または「クラスター定義」ダイアログ・ボックスで、抽出から除外したいコンセプトを選択します。
2. 右クリックしてコンテキスト・メニューを開きます。
3. 「不要語に追加」を選択します。コンセプトがリソース・エディターの不要語辞書に追加され、抽出結果ペインの背景色が変わり、変更を表示するには再抽出が必要であることを示します。変更が複数ある場合には、再抽出する前に変更を終わらせてください。

注: 除外した単語は、リソース・エディターのライブラリー・ツリーの表示された最初のライブラリーに自動的に保存されます。デフォルトでは、これがローカル・ライブラリーになります。

単語を抽出に強制投入

抽出後、データ・ペインでテキストデータを確認するとき、一部の単語または句が抽出されていないことが分かる場合があります。これらの単語は重要でない動詞または形容詞である場合が多くあります。ただし、抽出されなかった単語または句をカテゴリー定義として使用したい場合があります。

これらの単語および句を抽出したい場合、キーワードをタイプ・ライブラリーに強制投入できます。詳しくは、195 ページの『キーワードの強制』のトピックを参照してください。

重要: 辞書のキーワードを強制することが、絶対に確実というわけではありません。強制することによって、キーワードを辞書に明示的に追加している場合でも、再抽出後に抽出結果ペインに出現しない場合、あるいは出現しても宣言したとおりのものではない場合があります。抽出のプロセスにおいてある語や語句がより長い語句の一部として既に抽出されてしまっている場合や、語が品詞（名詞、動詞、形容詞、前置詞など）に分割されている場合が考えられます。これを回避するために、全体（複合語なし）マッチ・オプションをキーワード辞書のこのキーワードに適用します。詳しくは、192 ページの『キーワードの追加』のトピックを参照してください。

第 9 章 テキストデータの 카테고리化

カテゴリとコンセプト・ビューで、テキスト内の主要なキーワード、情報、属性をキャプチャする高いレベルのコンセプトまたはトピックを示すカテゴリを作成できます。

IBM SPSS Modeler Text Analytics のリリース 14 の時点では、カテゴリは階層的な構造を持つことができ、すなわち、サブカテゴリを含むことができ、また、そのサブカテゴリにもそれ自身のサブカテゴリを更に下の階層に向かって持たせることができます。製品内にこのような階層的なカテゴリを構築することが可能であるだけでなく、階層的なカテゴリを持ち、以前はコード・フレームと呼ばれていた、定義済みのカテゴリ構造をインポートすることも可能です。

実際に、階層カテゴリにより、1個または複数のサブカテゴリを持つツリー構造を構築して、例えば異なるコンセプトやトピックの分野の項目をより正確にグループ化することができます。レジャー活動に関して簡単な例を挙げることができます。「時間があればどんな活動がしたいですか?」という質問に対する答えとして、「スポーツ」、「日曜大工」、「釣り」などをトップ カテゴリに設定し、「スポーツ」の下の階層に「球技」、「水泳」などを設定できます。

カテゴリは、コンセプト、タイプ、パターンおよびカテゴリ規則などの一連の記述子で構成されています。また、これらの記述子を共に使用して、ドキュメントまたはレコードが指定されたカテゴリに属するかどうかを特定します。ドキュメントまたはレコード内のテキストをスキャンして、テキストが記述子に一致するかどうかを確認することができます。一致が見つかった場合は、ドキュメントまたはレコードはそのカテゴリに割り当てられます。このプロセスを、カテゴリ化といいます。

カテゴリとコンセプト・ビューの 4 つのパネルに表示されたデータを使用して、カテゴリを処理、作成および視覚的に検証することができます。また、「表示」メニューからその名前を選択して隠したり表示したりできます。

- カテゴリ・ペイン: このパネルでカテゴリを作成し、管理します。詳しくは、96 ページの『カテゴリ・ペイン』のトピックを参照してください。
- 抽出結果ペイン。このパネルで抽出したコンセプトおよびタイプを検証および処理します。詳しくは、81 ページの『抽出結果: コンセプトとタイプ』のトピックを参照してください。
- 視覚化ペイン。このパネルでカテゴリについて、またカテゴリがどのように相互作用するかを視覚的に検証します。詳しくは、155 ページの『カテゴリ・グラフおよび図表』のトピックを参照してください。
- データ・ペイン。このパネルでの選択に対応するドキュメントおよびレコード内に含まれるテキストを検証および確認できます。詳しくは、104 ページの『データ・ペイン』のトピックを参照してください。

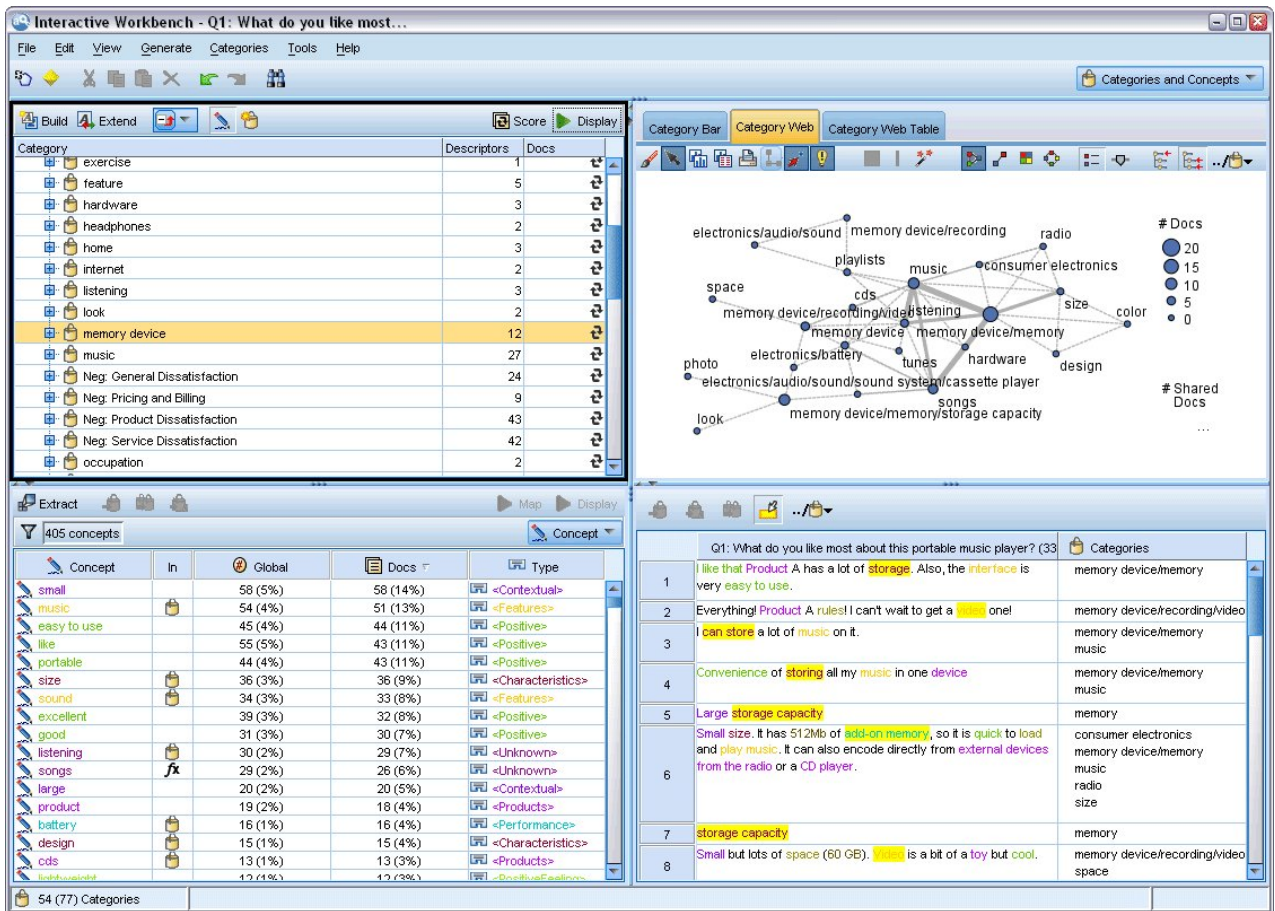


図 29. カテゴリーとコンセプト・ビュー

テキスト分析パッケージ (TAP) の一連のカテゴリーからはじめることができますが、独自のカテゴリーを作成する、または定義済みカテゴリー・ファイルからインポートする必要もあります。コンセプト、タイプ、パターンの抽出結果を使用する製品の自動化手法の頑健なセットを使用して、カテゴリーおよびそれらの記述子を自動的に作成することができます。データに関する追加の洞察を使用して、カテゴリーを手動で作成することもできます。ただし、インタラクティブ・ワークベンチを使用してのみ、カテゴリーを手動で作成し、調整できます。詳しくは、24 ページの『テキストマイニング・ノード: 「モデル」タブ』のトピックを参照してください。抽出結果をカテゴリーに手動でドラッグ・アンド・ドロップして、カテゴリー定義を作成できます。カテゴリー規則をカテゴリーに追加し、独自の事前定義済みカテゴリーを使用して、これらのカテゴリーまたは空のカテゴリーの品質を向上させることができます。

それぞれの手法は、特定の種類のデータおよび状況に適していますが、同じ分析で手法を組み合わせると、全範囲のドキュメントまたはレコードをキャプチャーすると役に立つ場合があります。またカテゴリー化を行う際には、言語リソースに対して変更を加えたほうがよい場合もあります。

カテゴリー・ペイン

カテゴリー・ペインでは、カテゴリーを作成および管理できます。このペインはカテゴリーとコンセプト・ビューの左上隅に表示されます。テキスト・データからコンセプトとタイプを抽出した後、内包関係のコンセプト、共起など、自動的な方法を使用して、または手動で作成してカテゴリーの作成を開始できます。詳しくは、106 ページの『カテゴリーの作成』のトピックを参照してください。

カテゴリを作成または更新するごとに、ドキュメントまたはレコードが「スコア」 ボタンをクリックするとスコアリングされ、いずれかのテキストが任意のカテゴリ内のデスクリプターと合致するか確認されます。一致が見つかった場合は、ドキュメントまたはレコードはそのカテゴリに割り当てられます。この最終結果は、ドキュメントまたはレコードのすべてではなくともその多くが、カテゴリの記述子に基づいて、カテゴリに割り当てられます。

注: 表示中のペインにカテゴリを表示しきれない場合は、ペインの下部にあるコントロールを使用して前後のカテゴリに移動したり、移動先のページ番号を入力したりすることができます。

カテゴリ・ツリー一覧

このパネルのツリー一覧には、一連のカテゴリ、サブカテゴリ、記述子が表示されます。ツリーには、各ツリー項目の情報を示す列があります。表示できるのは次の列です。

- **コード** 各カテゴリのコード値を表示します。この列は、デフォルトでは非表示になっています。「表示」 > 「カテゴリ・ペイン」メニューで、この列を表示することができます。
- **カテゴリ:** カテゴリ・ツリーには、カテゴリおよびサブカテゴリの名前が表示されます。また、記述子のツールバー・アイコンをクリックすると、一連の記述子も表示されます。
- **記述子:** その定義を構成する記述子の数が表示されます。この度数には、サブカテゴリの記述子数は含まれません。「カテゴリ」列に記述子名が表示されている場合、度数は表示されません。ツリー内の記述子自体の表示と非表示を行うには、「表示」 > 「カテゴリ・ペイン」 > 「すべての記述子」メニューを使用します。
- **ドキュメント** スコアリング後、この列には、該当するカテゴリとすべてのサブカテゴリにカテゴリ化されているドキュメントまたはレコードの数が表示されます。つまり、5つのレコードが記述子に基づいて上位カテゴリに合致し、7つの異なるレコードがその記述子に基づいてサブカテゴリに合致する場合、上位カテゴリのドキュメント数の合計は、この2つの数値の合計、この場合は12となります。ただし、同じレコードが上位カテゴリとそのサブカテゴリに合致する場合、度数は11となります。

カテゴリがない場合でも、テーブルには2つの行が表示されます。「すべてのドキュメント」、ドキュメントまたはレコードの総数が表示されます。2番目の行「未カテゴリ化」には、カテゴリがされていないドキュメント/レコード数が表示されます。

パネルの各カテゴリについて、小さい黄色のバケツのアイコンがカテゴリ名の前に表示されます。カテゴリをダブルクリックを選択するか、メニューで「表示」 > 「カテゴリ定義」を選択すると、カテゴリ定義ダイアログボックスが開き、記述子と呼ばれるすべての要素が表示されます。記述子はコンセプト、タイプ、パターン、カテゴリ・ルールなどの定義を決定します。詳しくは、103ページの『カテゴリとは』のトピックを参照してください。デフォルトでは、カテゴリ・ツリー一覧には、カテゴリの記述子は表示されません。「カテゴリ定義」ダイアログ・ボックスではなくツリーで直接記述子を表示する場合、ツールバーの鉛筆のアイコンが表示された切り替えボタンをクリックします。この切り替えボタンを選択すると、ツリーが展開され、記述子が表示されます。

カテゴリのスコアリング

カテゴリ・ツリー一覧の「ドキュメント」列には、特定のカテゴリにカテゴリ化されているドキュメント数またはレコード数が表示されます。数値が過去のものまたは計算されていない場合、アイコンがその列に表示されます。パネル ツールバーの「スコア」ボタンをクリックして、ドキュメント数を再計算することができます。大きいデータセットを使用する場合、スコアリング・プロセスには時間がかかる場合があります。

ツリー内のカテゴリの選択

ツリー内で選択すると、横グループのカテゴリのみ選択できます。つまり、上位レベルのカテゴリを選択すると、サブカテゴリは選択できません。または指定されたカテゴリの 2 つのサブカテゴリを選択すると、同時に別のカテゴリのサブカテゴリを選択できません。不連続なカテゴリを選択すると、以前の選択内容が失われます。

データ・ペインおよび視覚化ペインの表示

テーブルで 1 行を選択すると、「表示」 ボタンをクリックして、視覚化ペインおよびデータ・ペインを更新して、選択内容に対応する情報を表示します。パネルが表示されない場合、「表示」 をクリックしてパネルを開きます。

カテゴリの調整

カテゴリ化を行っても、最初から完全な結果が得られるとは限りません。削除したいカテゴリや、他のカテゴリとまとめたいカテゴリもあるでしょう。また、抽出結果を確認して、役立つと思われるいくつかのカテゴリが作成されていないことが分かる場合があります。その場合、結果を手動で変更し、特定の状況に対して結果を調整することができます。詳しくは、139 ページの『カテゴリの編集および調整』のトピックを参照してください。

カテゴリ作成の方法と戦略

また抽出していない、または抽出結果が古い場合、カテゴリ作成方法または拡張方法のいずれかを使用すると、抽出についてのプロンプトが自動的に表示されます。手法を適用した後、カテゴリにグループ化したコンセプトおよびタイプはその他の手法で構築したカテゴリに使用できます。つまり、再利用しないことを選択しないかぎり、複数のカテゴリのコンセプトを表示することができます。

最適なカテゴリを作成するために、次のことを確認してください。

- カテゴリ作成の方法
- カテゴリ作成の戦略
- カテゴリ作成のヒント

カテゴリ作成の方法

すべてのデータセットが一意であるため、カテゴリ作成方法の数やそれらを適用する順序は、時間によって変わる場合があります。また、テキストマイニングの目標が、データセットによって異なる場合があるため、それぞれの手法を検証して、指定したテキストデータにとってどの手法が最良の結果を生み出すかを確認する必要があります。自動的手法では、データを完全にカテゴリ化できません。そのため、データに合った 1 つまたは複数の自動的手法を見つけ、適用することをお勧めします。

事前作成されたカテゴリ・セットを持つテキスト分析パッケージ (TAP、*.tap) を使用するほか、次の方法を組み合わせて回答をカテゴリに分類することもできます。

- **自動作成手法:** いくつかの言語ベースおよび頻度ベースのカテゴリオプションを使用して、カテゴリを自動的に作成できます。詳しくは、106 ページの『カテゴリの作成』のトピックを参照してください。
- **自動拡張手法:** いくつかの言語学的手法を使用し、記述子を追加および拡張することによって既存のカテゴリを展開し、より多くのレコードをキャプチャーすることができます。詳しくは、116 ページの『カテゴリの拡張』のトピックを参照してください。

- 手動による手法: ドラッグアンドドロップなど、手動による手法がいくつかあります。詳しくは、119 ページの『手作業でのカテゴリーの作成』のトピックを参照してください。

カテゴリー作成の方略

次のリストの方略は包括的ではありませんが、カテゴリーの作成方法について、いくつかのキーワードが用意されています。

- テキスト マイニング・ノードを定義する場合、テキスト分析パッケージからカテゴリーを選択し、いくつかの作成済みカテゴリーを使用して分析を開始します。これらのカテゴリーは、テキストを最初から十分にカテゴリーに分類することができます。ただし、カテゴリーを追加する場合、カテゴリー作成設定（「カテゴリー」>「設定を行う」）を編集することができます。詳細設定: 言語学的手法ダイアログを開き、カテゴリー入力オプションの未使用の抽出結果を選択し、追加カテゴリーを作成します。
- ノードを定義するときに、インタラクティブ・ワークベンチのカテゴリーとコンセプト・ビューのTAPからカテゴリーを選択します。次に、未使用のコンセプトまたはパターンを適切なカテゴリーにドラッグ・アンド・ドロップします。そして、編集した既存のカテゴリーを展開（「カテゴリー」>「カテゴリーを展開」）し、既存のカテゴリー記述子に関連するより多くの記述子を取得します。
- 言語学的手法の詳細設定を使用して、自動的にカテゴリーを作成します（「カテゴリー」>「カテゴリーを作成」）。生成されるカテゴリーが適切なものとなるまで、記述子を削除、カテゴリーを削除、または同様のカテゴリーを結合してカテゴリーを手動で調整します。また、元来「可能な場合ワイルドカードを使用して一般化」オプションを使用しないでカテゴリーを作成する場合、「一般化」オプションをオンにして「カテゴリーを展開」を使用し、自動的にカテゴリーを簡略化することもできます。
- 非常に説明的なカテゴリー名および注釈を持つ事前定義済みカテゴリー・ファイルをインポートします。また、元来そのオプションを選択しないでインポートし、カテゴリー名から記述子をインポートまたは生成する場合、後で「カテゴリーの拡張」ダイアログを使用して、「カテゴリー名から生成された記述子を使用して空白のカテゴリを拡張する」オプションを選択することができます。そして、これらのカテゴリーの 2 回目の展開を行います。今回はグループ化手法を使用します。
- コンセプトまたはコンセプト・パターンを頻度によって並べ替え、最も関心のあるコンセプトまたはコンセプト・パターンをカテゴリー・ペインにドラッグ・アンド・ドロップすることによって、カテゴリーの最初のセットを手動で作成します。カテゴリーの最初のセットを作成したら、展開機能（「カテゴリー」>「カテゴリーを展開」）を使用して、すべてのカテゴリーを展開し、その他の関連記述子を含めてより多くのレコードに一致するように選択したカテゴリーを調整します。

これらの手法を適用した後、作成されたカテゴリーを確認、手動による手法を使用して小さな調整を行い、誤分類を削除、または欠損したと考えられるレコードまたは単語を追加することをお勧めします。また、さまざまな手法を使用すると、重複したカテゴリーを作成する場合もあるため、必要に応じてカテゴリーを結合または削除することもできます。詳しくは、139 ページの『カテゴリーの編集および調整』のトピックを参照してください。

カテゴリー作成のヒント

より良いカテゴリーを作成できるよう、方法を決定できるヒントをいくつか確認できます。

カテゴリー-to-ドキュメント比率のヒント

ドキュメントおよびレコードが割り当てられるカテゴリーは、少なくとも 2 つの理由で、質的テキスト分析で相互に排他的である場合はあまりありません。

- 1 つめの理由は、一般的に、テキスト ドキュメントまたはレコードが長いほど、表されるキーワードおよび意見がより明確なものとなることです。そのため、ドキュメントまたはレコードに複数のカテゴリーが割り当てられるという機会が大幅に増えます。

- 2 つめの理由は、論理的に分けられていないテキスト ドキュメントまたはレコードをグループ化および解釈するさまざまな方法があるということです。 回答者の政治的な信念に関する、自由記述式の質問を含んだ調査の場合、「リベラル」および「保守的」、または「共和党」および「民主党」のようなカテゴリーのほか、「社会的にリベラル」、「財政的に保守的」など、より特定のなカテゴリーも作成できます。 これらのカテゴリーは相互に重複部分がなかったり、すべてをカバーしたりしている必要はありません。

作成するカテゴリー数のヒント

カテゴリーは、データから直接作成する必要があります。データについて何か興味深いものがあった場合、カテゴリーを作成してその情報を示すことができます。一般的に、作成するカテゴリーに数について、推奨される上限はありません。ただし、カテゴリーをあまりに多く作成すると、管理できない場合があります。次の 2 つの原則が適用されます。

- **カテゴリー度数:** カテゴリーが役立つものになるには、最低限のドキュメントまたはレコードが必要です。1 つまたは 2 つのドキュメントには非常に興味深いものが含まれている場合がありますが、それが 1,000 件のドキュメントのうちの 1 つまたは 2 つである場合、含まれる情報は、実際に役立つほど母集団の中では頻繁ではありません。
- **複雑さ:** 作成したカテゴリーが多いほど、分析が完了した後確認および要約する必要がある情報が多くなります。ただし、カテゴリー数が多すぎる場合、複雑さが増しても、役立つ詳細は増えません。

カテゴリー数が多すぎることを判断する規則、またはカテゴリーあたりの最小レコード数を決定する規則はありません。個々の状況の必要性に応じて分析者が判断する必要があります。

しかしながら、基本的なアドバイスはあります。まずカテゴリーの個数は多すぎてもいいませんが、分析の早い段階においては、カテゴリーが少なすぎるよりは多すぎるほうがいいでしょう。比較的類似したカテゴリーをまとめるほうが、ケースを分けて新しいカテゴリーに細分化するよりも簡単なので、多くのカテゴリーからより少ない個数のカテゴリーになるように作業していくほうが、一般に容易だといえます。テキストマイニングの反復性およびこのソフトウェアプログラムを使用して達成できる容易さにより、より多くのカテゴリーを作成することが、開始時点では推奨されます。

最適な記述子の選択

次に、カテゴリーに最適な記述子 (コンセプト、タイプ、TLA パターンおよびカテゴリー規則) の選択または作成におけるガイドラインをいくつか示します。記述子とは、カテゴリーの構築ブロックです。ドキュメントまたはレコードの一部またはすべてのテキストが記述子に合致する場合、ドキュメントまたはレコードはカテゴリーに合致します。

記述子が抽出されたコンセプトまたはパターンを含まないまたは対応しない場合、ドキュメントまたはレコードに合致しません。そのため、次で説明しているように、コンセプト、タイプ、パターン、およびカテゴリー規則を使用します。

コンセプトが、それ自体だけでなく複数/単数形、類義語、およびスペルの変異形にいたる一連の基本キーワードも示すため、コンセプト自体は、記述子、または記述子の一部として使用する必要があります。指定されたコンセプトの基本キーワードについての詳細を知るには、カテゴリーとコンセプト・ビューの抽出結果ペインのコンセプト名をクリックします。コンセプト名でマウスポインタを停止すると、ツールヒントが表示され、そこに最後の抽出時にテキストで検出された基本キーワードが表示されます。すべてのコンセプトに基本キーワードがあるわけではありません。例えば、car と vehicle は類義語ですが car がコンセプトとして、vehicle が基本キーワードとして抽出された場合、vehicle を含むドキュメントまたはレコードに合致するため、記述子には car だけを使用します。

記述子としてのコンセプトとタイプ

コンセプト (または基本キーワードのいずれか) を含むすべてのドキュメントまたはレコードを検出する場合、そのコンセプトを記述子として使用します。この場合、正確なコンセプト名で十分であるため、より複雑なカテゴリ規則を使用する必要はありません。意見を抽出するリソースを使用する場合、文の真の意味を抽出する TLA パターン抽出時にコンセプトが変更される場合があるので注意してください (TLA に関する次の項の例を参照してください)。

例えば、「*Apple and pineapple are the best*」のような、各個人の好きな果物を示す調査の回答によって、apple と pineapple が抽出されます。コンセプト apple をカテゴリに記述子として追加すると、コンセプト apple (または基本キーワードのいずれか) を含むすべての回答がそのカテゴリに合致します。

ただし、とにかく apple について言及する回答を知りたい場合、カテゴリ規則を * apple * のように作成すると、apple、apple sauce、または french apple tart のようなコンセプトを含む回答もキャプチャできます。

また、<Fruit> のようにタイプを記述子として直接指定することによって、同じようにタイプ指定されたコンセプトを含むすべてのドキュメントまたはレコードをキャプチャすることもできます。タイプには * は使用できませんので注意してください。

詳しくは、81 ページの『抽出結果: コンセプトとタイプ』のトピックを参照してください。

記述子としてのテキスト リンク分析 (TLA) パターン

より詳細で微妙なアイデアをキャプチャする場合、TLA パターン結果を記述子として使用します。テキストが TLA 抽出中に分析されると、テキスト全体 (ドキュメントまたはレコード) を処理するのではなく、一度に 1 文、または 1 句ずつテキストが処理されます。1 つの文の全部分を考慮することによって、TLA は意見、2 つの要素間の関係性、または否定的な表現を特定して、真の意味を理解できます。コンセプト・パターンまたはタイプ・パターンを記述子として使用できます。詳しくは、151 ページの『タイプ・パターンおよびコンセプト・パターン』のトピックを参照してください。

例えば、「*the room was not that clean*」というテキストがある場合、次のようなコンセプトが抽出されます: room and clean。ただし、抽出設定で TLA 抽出が有効になっている場合、TLA で clean が否定文で使用されており、実際は not clean に対応し、コンセプト dirty の類義語であることを検出できます。ここで、記述子であるコンセプト clean がこのテキストに合致しますが、清潔さについて示すその他のドキュメントまたはレコードのキャプチャできることが確認できます。そのため、このテキストに合致し、より適切な記述子となるため、dirty を出力コンセプトに指定した TLA コンセプト・パターンを使用することをお勧めします。

記述子としてのカテゴリ・ビジネス規則

カテゴリ規則とは、抽出したコンセプト、タイプ、パターン、およびブール型演算子を使用した論理式に基づいて、ドキュメントまたはレコードを自動的に分類するステートメントです。例えば、「このカテゴリに、アルゼンチンではなく、大使館という抽出したコンセプトを含む」という内容を意味する式を作成することができます。

カテゴリ規則をカテゴリの記述子として記述および使用し、&、|、および !() によってさまざまなアイデアを表現することができます。ブール値。これらの規則のシンタックスおよびそれらの作成、編集方法の詳細は、「121 ページの『カテゴリ規則の使用』」を参照してください。

- & (AND) ブール演算子を含むカテゴリ規則を使用して、2 つ以上のコンセプトが出現するドキュメントまたはレコードを検出します。& 演算子でつながった 2 つ以上のコンセプトは、同じ文またはフ

レーズで出現する必要はありませんが、同じドキュメントまたはレコードのどこかに出現するとカテゴリーに合致すると見なされます。例えば、記述子にカテゴリー規則 `food & cheap` を作成すると、テキストに `food` と `cheap` の両方が含まれているため、`food` が `cheap` という名詞ではないにもかかわらず、「*the food was pretty expensive, but the rooms were cheap*」というテキストを含むレコードに合致します。

- `!()` を含むカテゴリー規則を使用します (NOT) ブール演算子が、いくつかのコンセプトまたはタイプのいずれかを含むドキュメントまたはレコードを検出します。コンテキストではなく、単語に基づいて関連すると思われる情報がグループ化されないようにします。例えば、カテゴリー規則 `<Organization> & !(ibm)` を記述子として作成すると、「*SPSS Inc. was a company founded in 1967*」というテキストには合致しますが、「*the software company was acquired by IBM*」には合致しません。
- `|` (OR) ブール演算子を含むカテゴリー規則を使用して、いくつかのコンセプトまたはタイプのいずれかを含むドキュメントまたはレコードを検出します。例えば、カテゴリー規則 `(personnel|staff|team|coworkers) & bad` を記述子として作成すると、これらの名詞がコンセプト `bad` と共に検出されるドキュメントまたはレコードと合致します。
- カテゴリー規則にタイプを使用すると、その規則はより一般的になり、より展開しやすくなる場合があります。例えば、ホテルのデータを扱っている場合、ホテルのスタッフに対して顧客がどう思っているかについて、非常に興味があります。関連するキーワードには、`receptionist`、`waiter`、`waitress`、`reception desk`、`front desk` (受付、ウェイター、ウェイトレス、受付デスク、フロント・デスク) などがあります。この場合、`<HotelStaff>` という新しいタイプを作成し、上記のキーワードをすべてこのタイプに追加します。 `[* waitress * & nice]`、`[* desk * & friendly]`、`[* receptionist * & accommodating]` のようなすべての種類のスタッフに 1 つのカテゴリー規則を作成することができますが、タイプ `<HotelStaff>` を使用して 1 つの、より一般的なカテゴリー規則を作成して `[<HotelStaff> & <Positive>]` の形式でホテル・スタッフに対して好意的な意見の回答をすべてキャプチャできます。

注： カテゴリー規則に TLA パターンを含む場合、これらの規則に `+` と `&` を両方使用できます。詳しくは、123 ページの『カテゴリー規則内の TLA パターンの使用』のトピックを参照してください。

記述子であるコンセプト、TLA、またはカテゴリー規則がどのように合致するかについての例

次の例では、記述子としてコンセプト、カテゴリー規則、TLA パターンを使用するとドキュメントまたはレコードがどのようにカテゴリー化されるのかについて説明します。次のような 5 つのレコードがあります。

- A: "*awesome restaurant staff, excellent food and rooms comfortable and clean.*" (素晴らしいレストラン・スタッフ、食べ物もおいしく、部屋は快適で清潔)
- B: "*restaurant personnel was awful, but rooms were clean.*" (レストラン・スタッフはひどいが、部屋はきれいだった)
- C: "*Comfortable, clean rooms.*" (快適で清潔な部屋)
- D: "*My room was not that clean.*" (私の部屋はきれいではなかった)
- E: "*Clean.*" (きれいだった)

レコードに「`clean`」という単語を含み、この情報をキャプチャーしたいため、次の表に示す記述子のいずれかを作成します。キャプチャーしようとしている核心に基づき、ある種類の記述子を別の記述子に対して使用すると、どのように異なる結果を生成するかを確認できます。

表 17. レコードが記述子に合致する例：

記述子	A	B	C	D	E	説明
clean	match	match	match	match	match	記述子は抽出されたコンセプトです。TLA のないレコード D も含めてすべてのレコードがコンセプト clean を含み、TLA 規則によって「not clean」が「dirty」を意味することは自動的に認識されません。
clean + .	-	-	-	-	match	記述子は clean を示す TLA パターンです。clean が TLA 抽出時に関連するコンセプトなしで抽出されたレコードとのみ合致します。
[clean]	match	match	match	-	match	記述子は、それ自体またはその他の規則に clean を含む TLA 規則を探すカテゴリ規則です。clean を含む TLA 出力が、clean が room のような別のコンセプトにリンクするかどうか、別のスロット位置にあるかどうかに関係なく、検出されたすべてのレコードに合致します。

カテゴリとは

カテゴリは、密接に関連したコンセプト、オピニオン、または属性のグループのことを言います。短いことば（ラベル）を付けて、カテゴリの内容が簡単にわかるようにしておくと便利です。

例えば、新しい洗濯用洗剤について消費者からのアンケートの回答を分析する場合、製品の香りを示すすべての回答を含む「香り」というラベルの付いたカテゴリを作成できます。ただし、そのようなカテゴリは良い香りと感じた消費者と、香りが強いと感じた消費者とを差別化するものではありません。IBM SPSS Modeler Text Analytics は適切なリソースを使用する場合意見の抽出ができるため、2 つのカテゴリを他に作成して、「香りが良い」と答えた回答者と「香りは好みではない」と答えた回答者を特定することができます。

カテゴリとコンセプト・ビューウィンドウの左上のパネルのカテゴリ・ペインで、カテゴリを作成したり作業することができます。各カテゴリは、1 つまたは複数の記述子で定義されます。記述子は、カテゴリを定義するために使用されているコンセプト、タイプ、パターンおよびカテゴリ規則です。

指定のカテゴリを構成する記述子を表示する場合、カテゴリ・ペインのツールバーの鉛筆のアイコンをクリックし、ツリーを展開して記述子を表示します。また、カテゴリを選択して「カテゴリ定義」ダイアログ・ボックスを開きます（「表示」>「カテゴリ定義」）。

内包関係のコンセプトなどのカテゴリ作成手法を使用してカテゴリを自動的に作成する場合、その手法ではコンセプトおよびタイプを記述子として使用し、カテゴリを作成します。TLA パターンを抽出すると、これらのパターンまたはパターンの一部をカテゴリ記述子として追加することもできます。詳しくは、149 ページの『第 11 章 テキスト リンク分析の検証』のトピックを参照してください。そしてクラスターを作成すると、クラスターのコンセプトを新しいまたは既存のカテゴリに追加することができます。最後に、カテゴリ規則を手動で作成して、カテゴリの記述子として使用することができます。詳しくは、121 ページの『カテゴリ規則の使用』のトピックを参照してください。

カテゴリのプロパティ

記述子のほかカテゴリには、カテゴリの名前を変更、ラベルまたは注釈を追加するために編集できるプロパティもあります。

以下のようなプロパティがあります。

- 「名前」 この名前は、デフォルトでツリーに表示されます。自動化の手法でカテゴリを作成した場合、名前は自動的に付けられます。
- ラベル: ラベルの使用は、他の製品または他のテーブルまたはグラフで使用する、より重要なカテゴリ記述子を作成する場合に役立ちます。「名前の代わりにラベルを表示」を選択すると、インターフェース内でカテゴリを特定する際にラベルが使用されます。
- コード: コード番号は、このカテゴリのコード値に対応します。
- 注釈: このフィールドで各カテゴリの短い説明を追加できます。「カテゴリを作成」ダイアログでカテゴリを生成した場合、この注釈にメモが自動的に追加されます。テキストを選択し、メニューから「カテゴリ」>「注釈に追加」を選択して、データ・ペインからサンプル テキストを注釈に直接追加することもできます。

データ・ペイン

カテゴリを作成した場合、作業しているテキスト・データを確認したい場合があります。例えば、640 件のドキュメントがカテゴリ化されているカテゴリを作成する場合、実際にどのようなテキストが書かれているのかを確認するため、これらのドキュメントの一部またはすべてに目を通したい場合があります。右下のデータ・ペインでレコードまたはドキュメントを確認することができます。デフォルトで表示されない場合は、メニューから「表示」>「パネル」>「データ」を選択してください。

データ・ペインには、特定の表示制限に応じて、カテゴリ・ペイン、抽出結果ペイン、または「カテゴリ定義」ダイアログ・ボックスの選択に該当するドキュメントまたはレコードごとに 1 行ずつ表示されます。デフォルトでは、データ・ペインに表示されるドキュメントまたはレコード数が制限され、データをより迅速に表示できるようになります。ただし、これは「オプション」ダイアログ・ボックスで調整できます。処理しているデータセットが非常に大きい場合は、カテゴリを表示するオプションをオフにすると表示速度が向上する場合があります。詳しくは、76 ページの『オプション: 「セッション」タブ』のトピックを参照してください。

注: 表示中のペインにレコードを表示しきれない場合は、ペインの下部にあるコントロールを使用して前後のレコードに移動したり、移動先のページ番号を入力したりすることができます。

データ・ペインの表示および更新

データ・ペインでは、大きなデータセットの自動データ更新には時間がかかるため、自動的に表示の更新は行われません。そのため、このビューの別のパネルまたは「カテゴリ定義」ダイアログ・ボックスで選択すると、「表示」をクリックしてデータ・ペインの内容を更新します。

テキスト・ドキュメントまたはレコード

テキスト・データがレコードの形式で、テキストの長さが比較的短い場合、データ・ペインのテキスト・フィールドには、テキスト・データの全体が表示されます。ただし、レコードおよび大きいデータセットを処理している場合、テキスト・フィールドの列にはテキストの一部が表示され、右側のテキスト・プレビュー・ペインを開くと、テーブルで選択したレコードの大部分またはすべてが表示されます。テキスト・デー

タが個別ドキュメントの形式の場合、データ・ペインには、ドキュメントのファイル名が表示されます。ドキュメントを選択すると、テキスト・プレビュー・ペインには選択したドキュメントのテキストが表示されます。

色および強調表示

データを表示すると、該当するドキュメントまたはレコードのコンセプトおよび記述子は色付きで強調表示され、テキスト内のコンセプトおよび記述子を特定しやすくなります。カラー・コードは、コンセプトが属するタイプに対応します。カラーコード化された項目上でマウス・ポインタを停止させて、項目が抽出されたコンセプトと、項目が割り当てられたタイプを表示することもできます。抽出されていないテキストは、黒で表示されます。通常、こうした抽出されていない単語は接続詞（「および」または「と」）、代名詞（「私」または「彼ら」）および動詞（「いる」、「持つ」、または「取る」）のケースが多くあります。

データ・ペインの列

テキスト・フィールドの列は常に表示されていますが、その他の列も表示できます。その他の列を表示するには、メニューで「表示」 > 「データ・ペイン」を選択し、データ・ペインに表示したい列を選択します。表示できるのは次の列です。

- 「テキスト・フィールド名」(#)/ドキュメント コンセプトおよびタイプが抽出されたテキスト・データの列を追加します。データがドキュメントにある場合、列名は「ドキュメント」となり、ドキュメント・ファイル名または完全パスのみが表示されます。これらのドキュメントのテキストを表示するには、テキスト・プレビュー・ペインを表示する必要があります。データ・ペインの行数が、この列名の後のカッコ内に表示されます。読み込みの速度向上のために使用される「オプション」ダイアログの制約により、一部のドキュメントまたはレコードが表示されない場合があります。最大値に達すると、数値の後に「- 最大」と表示されます。詳しくは、76 ページの『オプション：「セッション」タブ』のトピックを参照してください。
- カテゴリー レコードが属するカテゴリーがそれぞれ表示されます。この列を表示する場合、最新の情報を示すため、データ・ペインの更新に少し時間がかかる場合があります。
- 適合度順位 1 つのカテゴリーの各レコードの順位が表示されます。この適合度順位は、カテゴリー内の他のレコードと比較して、レコードがカテゴリーにどれだけ適合しているかを示します。カテゴリー・ペイン (左上のパネル) でカテゴリーを選択すると、順位が表示されます。詳しくは、『カテゴリーの関連性』のトピックを参照してください。
- カテゴリーの個数 レコードが割り当てられているカテゴリー数が表示されます。

カテゴリーの関連性

より良いカテゴリーを作成するため、各カテゴリーのドキュメントまたはレコードの関連性のほか、ドキュメントまたはレコードが属するすべてのカテゴリーの関連性を確認できます。

カテゴリーのレコードに対する関連性

データ・ペインにドキュメントまたはレコードが表示されると、それらすべてのカテゴリーが「カテゴリー」列に表示されます。ドキュメントまたはレコードが複数のカテゴリーに含まれる場合、この列のカテゴリーは、関連性が最も大きなものから小さなもの順に表示されます。最初に表示されたカテゴリーは、このドキュメントまたはレコードに最も対応すると考えられます。詳しくは、104 ページの『データ・ペイン』のトピックを参照してください。

レコードのカテゴリーに対する関連性

カテゴリーを選択すると、データ・ペインの「適合度順位」に各レコードの関連性が表示されます。この適合度順位は、ドキュメントまたはレコードが選択したカテゴリーにどれだけ適合しているかを、そのカテゴリーのほかのレコードと比較して示します。単一のカテゴリーのレコードの順位を確認するには、左上のカテゴリー・ペインでそのカテゴリーを選択します。ドキュメントまたはレコードの順位が列に表示されます。デフォルトでは、この列は表示されませんが、列が表示されるよう選択することができます。詳しくは、104 ページの『データ・ペイン』のトピックを参照してください。

レコードの順位が低いほど、このレコードの、選択したカテゴリーに対する適合度または関連性が大きくなり、1 が最も適合度が高くなります。複数のレコードの関連性が同じ場合、それぞれが同じ順位で表示され、その後に関連性が同じであることを示す等号 (=) が追加されます。例えば、1=、1=、3、4 などのようになります。このカテゴリーに最もマッチするレコードが 2 つあることを意味します。

ヒント: 最も関連性の高いレコードのテキストをカテゴリーの注釈に追加して、カテゴリーの説明をより分かりやすくすることができます。テキストを選択し、メニューから「カテゴリー」>「注釈に追加」を選択して、データ・ペインからテキストを直接追加します。

カテゴリーの作成

テキスト分析パッケージのカテゴリーがある場合でも、さまざまな言語学的手法および出現頻度に基づく手法を使用して、カテゴリーを自動的に作成することもできます。「カテゴリー作成設定」ダイアログを使用して、自動的な言語学的手法および出現頻度に基づく手法を適用し、コンセプトまたはコンセプト・パターンのいずれかよりカテゴリーを作成することができます。

一般的に、カテゴリーはさまざまな記述子 (タイプ、コンセプト、TLA パターン、カテゴリー規則) で構成されます。自動化されたカテゴリー作成方法を使用してカテゴリーを作成する場合、作成されたカテゴリーは、選択した入力に応じてコンセプトまたはコンセプトのパターンから命名され、それぞれ一連の記述子を使用します。これらの記述子は、カテゴリー規則またはコンセプトの形式で、手法によって発見されたすべての関連コンセプトを含む場合があります。

カテゴリーを作成した後、カテゴリー・ペインで確認またはグラフや図表を使用して検討することにより、カテゴリーについて多くのことを学ぶことができます。手作業で若干の調整を行い、間違っただけの分類を削除したり、把握されなかったレコードや語を追加することもできます。手法を適用した後、カテゴリーにグループ化したコンセプト、種類、パターンは他の手法にも使用できます。また、さまざまな手法を使用すると、重複したカテゴリーまたは不適切なカテゴリーを作成する場合もあるため、カテゴリーを結合または削除することもできます。詳しくは、139 ページの『カテゴリーの編集および調整』のトピックを参照してください。

重要: 以前のリリースでは、共起規則および類義語規則は、大カッコで囲まれていました。このリリースの場合、大カッコはテキスト リンク分析パターン結果を示します。その代わりに、共起規則や類義語規則は、(スピーカー・システム|スピーカー) のように、カッコで囲まれます。

カテゴリーを作成するには

1. メニューの「カテゴリー」>「カテゴリーを作成」を選択します。プロンプトが表示されないよう設定している場合、メッセージ・ボックスが表示されます。
2. 今すぐ作成するか、左記に設定を編集するかを選択します。
 - 「今すぐ作成」をクリックすると、現在の設定でカテゴリーの作成が開始されます。デフォルトで選択されている設定で、十分カテゴリー化プロセスを開始できます。カテゴリー作成のプロセスが開始し、進捗状況のダイアログが表示されます。
 - 「編集」をクリックして、ビルド設定を確認し、変更します。

注: 表示できる最大カテゴリ数は 10,000 です。この数に到達したか超過した場合は警告が表示されます。その場合は、「カテゴリの作成」オプションまたは「カテゴリの拡張」オプションを変更して作成カテゴリの数を少なくする必要があります。

入力

カテゴリは、タイプ・パターンまたはタイプのいずれかより派生した記述子から作成されます。表内で、カテゴリ作成プロセスに使用する各タイプまたはパターンを選択できます。

タイプ・パターン :タイプ・パターンを選択すると、タイプおよびコンセプトではなくパターンからカテゴリが独自に作成されます。このように、選択したタイプ・パターンに属するコンセプト・パターンを含むレコードまたはドキュメントがカテゴリされます。そのため、表で <Budget> タイプ・パターンおよび <Positive> タイプ・パターンを選択した場合、cost & <Positive> または rates & excellent などのカテゴリを作成することができます。

自動カテゴリ作成でタイプ・パターンを入力として使用すると、その手法によってカテゴリ構造を形成する複数の方法が特定される場合があります。技術的には、カテゴリを作成する適切な方法はありませんが、分析により適した構造がある場合があります。この場合の出力をカスタマイズするために、タイプを優先的に指定できます。作成されたすべての上位レベルのカテゴリは、ここで選択したタイプのコンセプトのみに由来します。すべてのサブカテゴリには、このタイプのテキスト・リンク パターンが含まれています。このタイプを「パターン・タイプ:フィールドにより構造カテゴリ」で選択すると、テーブルが更新され、選択されたタイプを含む適用パターンのみを表示します。多くの場合、<Unknown> が事前に選択されています。この結果、タイプ <Unknown> を含むすべてのパターンが選択されます。表には、タイプ、レコード数またはドキュメント数 (**Doc.** の数) の最も多いものから降順で表示されます。

タイプ :タイプを選択すると、カテゴリは選択したタイプに属するコンセプトから作成されます。そのため、表で <Budget> タイプを選択した場合、cost および price は <Budget> に割り当てられたコンセプトであるため、cost または price のようなカテゴリを作成できます。

デフォルトでは、最も多いレコードまたはドキュメントをキャプチャーするタイプのみが選択されます。このように事前選択すると、最も関心の高いタイプにすばやく焦点をあて、関心の低いカテゴリが作成されないようにすることができます。表には、タイプ、レコード数またはドキュメント数 (**Doc.** の数) の最も多いものから降順で表示されます。意見ライブラリーのタイプは、デフォルトではタイプ テーブルで選択されていません。

入力の選択は、取得するカテゴリに影響します。入力としてタイプを選んだ場合は、明確に関連付けられたコンセプトをより簡単に見ることができます。例えば、タイプを入力として使用してカテゴリを作成する場合、果物 というカテゴリを、リンゴ、梨、柑橘類、オレンジなどのコンセプトとともに取得できます。他方、タイプ・パターンを入力として選び、パターンとして例えば <不明> + <肯定的> を選択した場合には、果物 + <肯定的> というカテゴリに、果物 + おいしい や リンゴ + 良いなどの 1、2 種類の果物を伴ったものを取得することになるでしょう。この第 2 の結果は、2 つのコンセプト・パターンを示すのみです。それは、他の果物の出現が必ずしも肯定的に評価されるものとは限らないからです。また、現在手元にあるテキスト・データについてはこれで十分であったとしても、異なるドキュメント・セットを使用する経時的な研究においては、柑橘類 + 肯定的のような別の記述子を手動で追加したり、タイプを使いたいと考えることもあるでしょう。タイプだけを入力として使用することは、すべての可能な果物を見つけ出すのに役に立ちます。

手法

すべてのデータセットが一意であるため、手法の数やそれらを適用する順序は、時間によって変わる場合があります。テキストマイニングの目標が、データセットによって異なる場合があるため、それぞれの手法を検証して、指定したテキスト・データにとってどの手法が最良の結果を生み出すかを確認する必要があります。

これらの設定について詳しく知らなくても、使用することができます。デフォルトでは、最も一般的で平均的な設定がすでに選択されています。そのため、高度な設定のダイアログを省略して、カテゴリーをすぐに作成することができます。同様に、ここで変更を行うと、最新の設定が常に保持されるため、設定ごとに設定ダイアログに戻る必要はありません。

言語学的手法または出現頻度に基づく手法のいずれかを選択し、「詳細設定」ボタンをクリックして、選択した手法の設定を表示します。自動的手法では、データを完全にカテゴリー化できません。そのため、データに合った 1 つまたは複数の自動的手法を見つけ、適用することをお勧めします。言語学的手法および出現頻度に基づく手法を同時に使用して作成することはできません。

- 高度な言語学的手法。詳しくは、『言語学的手法の詳細設定』を参照してください。
- 高度な出現頻度に基づく手法。詳しくは、115 ページの『出現頻度に基づく手法の詳細設定』を参照してください。

言語学的手法の詳細設定

カテゴリーを作成する場合、「派生関係のコンセプトの語幹」、「内包関係のコンセプト」、「セマンティック・ネットワーク」(英語テキストのみ)、および「共起規則」など、さまざまな詳細言語カテゴリー作成手法から選択することができます。これらの手法は個別に、またはそれぞれを組み合わせることでカテゴリーを作成することができます。

すべてのデータセットが一意であるため、手法の数やそれらを適用する順序は、時間によって変わる場合がありますので、注意してください。テキストマイニングの目標が、データセットによって異なる場合があるため、それぞれの手法を検証して、指定したテキスト・データにとってどの手法が最良の結果を生み出すかを確認する必要があります。自動的手法では、データを完全にカテゴリー化できません。そのため、データに合った 1 つまたは複数の自動的手法を見つけ、適用することをお勧めします。

以下は、「詳細設定: 言語学的手法」ダイアログ・ボックスにある領域とフィールドです。

入力および出力

カテゴリー入力: カテゴリーが作成される内容を選択します。

- 未使用の抽出結果: 既存のカテゴリーで使用されていない抽出結果からカテゴリーを作成できます。レコードが、複数のカテゴリーと合致する傾向が最も小さくなり、作成されるカテゴリーの数が制限されます。
- すべての抽出結果: 抽出結果のいずれを使用してもカテゴリーを作成できます。カテゴリーがないまたは少ない場合に最も役立ちます。

カテゴリー出力: 作成されるカテゴリーの一般的な構造を選択します。

- サブカテゴリーによる階層: サブカテゴリーおよびサブサブカテゴリーの作成を有効にします。作成できる最大数のレベル(「作成するレベルの最大個数」フィールド)を選択して、カテゴリーの深度を設定できます。3 を選択すると、カテゴリー内にサブカテゴリーを作成でき、またこれらのサブカテゴリー内にもサブカテゴリーを作成できます。

- フラットなカテゴリー (単一レベルのみ):1 レベルのみのカテゴリーを作成できます。つまり、サブカテゴリーは生成できません。

グループ化手法

使用できるそれぞれの手法は、特定の種類のデータおよび状況に適していますが、同じ分析で手法を組み合わせて、全範囲のドキュメントまたはレコードをキャプチャーすると役に立つ場合があります。複数のカテゴリーのコンセプトを表示したり、重複するカテゴリーを見つけることができます。

内包関係のコンセプト: この手法では、一方の共通の文字列である単語を含むかどうかに基づき、マルチタームのコンセプト (複合語) をグループ化することによってカテゴリーを作成します。例えば、コンセプト seat (シート) は、コンセプト safety seat (セーフティ シート)、seat belt (シート・ベルト)、および seat belt buckle (シート・ベルトのバックル) とグループ化されます。詳しくは、112 ページの『内包関係のコンセプト』のトピックを参照してください。

セマンティック・ネットワーク: 各コンセプトの考えられる意味を、単語の関係の拡張インデックスから特定することによって開始し、関連するコンセプトをグループ化することによってカテゴリーを作成します。この手法は、コンセプトがセマンティック・ネットワークに認識され、あまり曖昧でない場合に最も適しています。テキストに、ネットワークが認識していない特殊な用語または専門用語が含まれている場合はあまり役に立ちません。例えば、コンセプト granny smith apple は、granny smith と横の関係があるため、gala apple および winesap apple とグループ化されます。また別の例では、コンセプト animal は、その下位語である cat および kangaroo とグループ化されます。このリリースでは、英語テキストにのみ使用できます。詳しくは、113 ページの『セマンティック・ネットワーク』のトピックを参照してください。

注: 「最大検索距離」オプションを使用できるのは、「セマンティック ネットワーク」を選択した場合のみです。

最大検索距離: カテゴリー作成前に手法による検索の距離を選択します。ただし、これらの結果はノイズが少なく、またリンクや関連性が大きくなります。値が大きいほど、取得する結果は多くなります。ただし、これらの結果の信頼性または関連性が弱くなります。このオプションはすべての手法にグローバルに適用されますが、共起とセマンティック・ネットワークに対する効果は最も大きくなります。

特定のコンセプトがグループ化されないようにする: 出力の 2 つのコンセプトがグループ化またはペアとにならないように処理を停止します。コンセプト・ペアを作成または管理するには、「ペアを管理」をクリックします。詳しくは、110 ページの『例外ペアのリンクの管理』のトピックを参照してください。

可能な場合ワイルドカードを使用して一般化: アスタリスク (*) のワイルドカードを使用して、製品が一般的な規則をカテゴリーに生成することができますようになります。例えば、[アップル タルト + .] や [アップル ソース + .] などの複数の記述子を作成する代わりに、[アップル * + .] のようにワイルドカードを使用します。ワイルドカードを使用して一般化すると、以前と同じように、ちょうど同じ数のレコードまたはドキュメントを取得する場合があります。ただし、このオプションには、数の縮小やカテゴリーの記述子の簡略化という利点があります。また、このオプションを使用すると、新しいテキスト・データ (例: 長期的/周期的研究) にこれらのカテゴリーを使用してより多くのレコードまたはドキュメントをカテゴリー化する機能を拡大します。

カテゴリーを作成するその他のオプション

適用するグループ化手法を選択するほか、次のように、その他の作成オプションを編集することができます。

最大数の上位レベルカテゴリが作成されました。このオプションを使用して、「カテゴリを作成」をクリックすると作成できるカテゴリ数を制限します。この値を高く設定し、関心の低いカテゴリを削除すると、よりよい結果が生成される場合があります。

記述子の最小値と/もしくはカテゴリごとのサブカテゴリ: カテゴリが作成するために含む必要のある記述子数およびサブカテゴリ数の最小値を定義します。多くのレコードまたはドキュメントをキャプチャしないカテゴリの作成が制限されます。

同じ記述子が複数のカテゴリに含まれることを許可する: このオプションを選択すると、記述子を次の作成される複数のカテゴリに使用できるようにします。項目が一般的にまたは「自然に」2つ以上のカテゴリになり、より良い品質のカテゴリを作成するため、このオプションが一般的に選択されます。このオプションを選択しない場合、複数のカテゴリのレコードの重複が少なくなり、データのタイプによっては、これが望ましい設定となります。ただし、多くのデータ・タイプでは、記述子を1つのカテゴリに制限すると、品質またはカテゴリの範囲が損なわれます。例えば、car seat manufacturer というコンセプトがあったとします。このオプションを指定すると、このコンセプトは、テキスト car seat に基づいてあるカテゴリに、また manufacturer というテキストに基づいて別のカテゴリに使用されます。ただし、このオプションが選択されていない場合、2つのカテゴリを取得できますが、コンセプト car seat manufacturer は、car seat および manufacturer がそれぞれ出現するレコード数など、いくつかの要素に基づいて、最も一致するカテゴリにのみ、記述子として使用されます。

次の方法で重複するカテゴリ名を解決 名前が既存のカテゴリと同じ新規カテゴリまたはサブカテゴリの処理方法を選択します。新規カテゴリ (およびその記述子) と同じ名前を持つ既存カテゴリとを結合できます。あるいは、既存のカテゴリに重複する名前があった場合、カテゴリの作成をスキップすることもできます。

例外ペアのリンクの管理

カテゴリ作成、クラスタリングおよびコンセプト・マッピングの間、既知の関連性によって内部アルゴリズムが語をグループ化します。2つのコンセプトが対応しないように、またはお互いにリンクしないようにするためには、「カテゴリ作成詳細設定」 ダイアログ、「クラスターの作成」 ダイアログ、および「コンセプト マップ インデックス設定」 ダイアログでこの機能をオンにして、「ペアを管理」 ボタンをクリックします。

表示される「リンクの例外を管理」 ダイアログでコンセプト・ペアを追加、編集または削除できます。1行につき1つのペアを入力します。ここでペアを入力すると、カテゴリ作成または拡張時のグループ化、クラスタリング、コンセプト・マッピングが行われなくなります。必要に応じて、単語を正確に入力します。例えば、単語のアクセント付きバージョンがアクセントのないバージョンとは同じではありません。

例えば、hot dog および dog がグループ化されていないことを確認したい場合、ペアをテーブル内の各行に追加できます。

言語学的手法について

カテゴリを作成または拡張する場合、「派生関係のコンセプトの語幹」、「内包関係のコンセプト」、「セマンティック・ネットワーク」(英語のみ)、および「共起規則」など、さまざまな詳細言語カテゴリ作成手法から選択することができます。これらの手法は個別に、またはそれぞれを組み合わせることでカテゴリを作成することができます。

これらの設定について詳しく知らなくても、使用することができます。デフォルトでは、最も一般的で平均的な設定がすでに選択されています。必要に応じて、高度な設定のダイアログを省略して、カテゴリをす

ぐに作成または拡張することができます。同様に、ここで変更を行うと、最後に使用した設定が記憶されているため、設定ごとに設定ダイアログに戻る必要はありません。

ただし、すべてのデータセットが一意であるため、手法の数やそれらを適用する順序は、時間によって変わる場合がありますので、注意してください。テキストマイニングの目標が、データセットによって異なる場合があるため、それぞれの手法を検証して、指定したテキスト・データにとってどの手法が最良の結果を生み出すかを確認する必要があります。自動的手法では、データを完全にカテゴリー化できません。そのため、データに合った 1 つまたは複数の自動的手法を見つけ、適用することをお勧めします。

カテゴリーを作成する、主な自動化言語法は、次のとおりです。

- 派生関係のコンセプトの語幹: コンセプト・コンポーネントが形態的に関連するかどうかを分析するとき、コンセプトを取得し、そのコンセプトに関連するその他のコンセプトを検索することによって、カテゴリーを作成します。詳しくは、『派生関係のコンセプトの語幹』のトピックを参照してください。
- ないほう関係のコンセプト: コンセプトを取得し、そのコンセプトを含むその他のコンセプトを見つけることによって、カテゴリーを作成します。詳しくは、112 ページの『内包関係のコンセプト』のトピックを参照してください。
- セマンティック・ネットワーク: 各コンセプトの考えられる意味を、単語の関係の拡張インデックスから特定することによって開始し、関連するコンセプトをグループ化することによってカテゴリーを作成します。詳しくは、113 ページの『セマンティック・ネットワーク』のトピックを参照してください。このオプションは、英語テキストにのみ使用できます。
- 共起: 新しいカテゴリーを作成、カテゴリーを拡張するために、または別のカテゴリー手法の入力として使用できる共起規則を作成します。詳しくは、114 ページの『共起規則』のトピックを参照してください。

派生関係のコンセプトの語幹

派生関係のコンセプトの語幹による手法は、コンセプト・コンポーネントが形態的に関連するかどうかを分析するとき、コンセプトを取得し、そのコンセプトに関連するその他のコンセプトを検索することによって、カテゴリーを作成します。コンポーネントは単語です。コンセプトの各コンポーネントの末尾 (接尾辞) を確認し、それらから派生するその他のコンセプトを見つけることによって、コンセプトのグループ化を試みます。単語がお互いに派生している場合、共有するか、意味の上で近い傾向にあります。末尾を特定するために、内部の言語固有の規則が使用されます。末尾を識別するには、内部の言語固有の規則が使用されます。例えば、コンセプト `opportunities to advance` は、コンセプト `opportunity for advancement` および `advancement opportunity` とグループ化されます。

いかなる種類のテキストにも派生関係のコンセプトの語幹を使用できます。それ自体によって作成されるカテゴリーはごくわずかであり、各カテゴリーに含まれるコンセプトも少数です。各カテゴリーのコンセプトは、類義語または状況的に関連しています。手動でカテゴリーを作成する場合でも、このアルゴリズムを作成すると役立ちます。見つかった類義語は、特に関心のあるコンセプトの類義語である場合があります。

注: コンセプトを明示的に指定することにより、コンセプトがグループ化されないようにすることができます。詳しくは、110 ページの『例外ペアのリンクの管理』のトピックを参照してください。

キーワードのコンポーネント化および活用の無効化

派生関係のコンセプトの語幹による手法または内包関係のコンセプトの手法が適用されると、キーワードはまずコンポーネント (単語) に分割され、コンポーネントの活用が無効化されます。手法が適用されると、コンセプトおよびそれらに関連したキーワードが読み込まれ、スペース、ハイフン、アポストロフィなどの区切り文字に基づいて、コンポーネントに分割されます。例えば、`system administrator` というキーワードは、`{administrator, system}` のように、コンポーネントに分割されます。

ただし、元のキーワードの一部は使用できず、ストップワードとして参照されます。英語の場合、こうした無視できるコンポーネントには、a、and、as、by、for、from、in、of、on、or、the、to、およびwithなどがあります。

例えば、キーワード examination of the data のコンポーネントは {data, examination} のようになり、of および the は無視できるとみなされます。また、コンポーネント・セットでは、コンポーネントの順序は意味を持ちません。それらすべては同じコンポーネント・セット {child, cough, relief} を持つため、次の3つのキーワードは同等とすることができます。cough relief for child、child relief from a cough、relief of child cough キーワードのペアが同等のものとして特定されるごとに、対応するコンセプトを結合して、すべてのキーワードを参照する新しいコンセプトを形成します。

また、キーワードのコンポーネントは活用している場合があるため、言語固有の規則が内部的に適用され、複数形など、活用した変異形にかかわらず、同等のキーワードを特定します。このようにして、活用のない単数形が level であるため、キーワード level of support および support levels を同等のものとして特定することができます。

派生関係のコンセプトの語幹の機能

キーワードがコンポーネント化され、活用がなくなった場合 (前セクションを参照)、派生関係のコンセプトの語幹アルゴリズムがコンポーネントの末尾、または接尾辞を分析し、コンポーネントの語幹を検索して、そのコンセプトを、同じ、または類似したルートを持つ他のコンセプトとグループ化します。末尾は、テキスト言語に特有の言語派生規則を使用して特定されます。例えば、接尾辞 ical で終わるコンセプト・コンポーネントは、同じ語源を持ち、接尾辞 ic で終わるコンセプトから派生するという、英語のテキストの派生規則があるとします。この規則 (および活用の無効化) を使用すると、アルゴリズムはコンセプト epidemiologic study および epidemiological studies をグループ化できます。

キーワードはすでにコンポーネント化され、(in や of などの) 無視できるコンポーネントが特定されているため、派生関係のコンセプトの語幹アルゴリズムは、コンセプト studies in epidemiology を epidemiological studies とグループ化することもできます。

一連のコンポーネント派生関係の規則は、このアルゴリズムでグループ化されるコンセプトの大半が類義語となるように選ばれました。例えば、epidemiologic studies、epidemiological studies、studies in epidemiology という3つのコンセプトはすべて同義のキーワードです。完全性を高くするために、一部の派生規則を使用すると、アルゴリズムによって、状況的に関連するコンセプトをグループ化できます。例えば、アルゴリズムは empire builder や empire building などのコンセプトをグループ化できます。

内包関係のコンセプト

内包関係のコンセプトの手法は、コンセプトを取得し、語彙系列のアルゴリズムを使用してカテゴリを作成し、その他のコンセプトに含まれるコンセプトを特定します。コンセプト内の単語が別のコンセプトの部分集合である場合、規定となるセマンティックの関係を反映します。内包関係のコンセプトは、どんな種類のテキストにも使用できる強力な手法です。

この方法は、セマンティック・ネットワークと組み合わせるとうまく動作しますが、個別に使用することもできます。ドキュメントまたはレコードにドメイン固有の用語または専門用語が多く含まれている場合、内包関係のコンセプトを使用するとより良い結果が出ます。これは事前に類義語辞書を使用して、特別なキーワードが適切に抽出・グループ化されるように調整してある場合に、特にいい結果が得られます。

内包関係のコンセプトの機能

内包関係のコンセプト・アルゴリズムを適用する前に、キーワードはコンポーネント化され、活用がありません。詳しくは、111ページの『派生関係のコンセプトの語幹』のトピックを参照してください。次

に、内包関係のコンセプト・アルゴリズムはコンポーネント・セットを分析します。各コンポーネント・セットについて、アルゴリズムは最初のコンポーネント・セットの部分集合である別のコンポーネント・セットを検索します。

例えば、コンポーネント・セットが {breakfast, continental} のコンセプト continental breakfast があり、コンポーネント・セットが {breakfast} のコンセプト breakfast がある場合、アルゴリズムは、continental breakfast は breakfast の一種であると結論付け、これらをともにグループ化します。

より大きな例では、抽出結果ペインにコンセプト seat が表示され、このアルゴリズムを適用する場合、safety seat、leather seat、seat belt、seat belt buckle、infant seat carrier、および car seat laws のようなコンセプトは該当するカテゴリーにグループ化されます。

キーワードはすでにコンポーネント化され、(in や of などの) 無視できるコンポーネントが特定されているため、内包関係のコンセプト・アルゴリズムは、コンセプト advanced spanish course にコンセプト course in spanish が含まれると認識します。

注: コンセプトを明示的に指定することにより、コンセプトのグループ化を禁止することができます。詳しくは、110 ページの『例外ペアのリンクの管理』のトピックを参照してください。

セマンティック・ネットワーク

このリリースでは、セマンティック・ネットワーク手法は、英語テキストにのみ使用できます。

この手法では、単語の関係の組み込みネットワークを使用してカテゴリーを作成します。このため、キーワードが具体的で、あまりあいまいでなければ、この手法を使用すると、非常に良い結果を生成することができます。ただし、この手法が非常に技術的/専門的なコンセプト間に多くのつながりを見つけることを期待することができません。こうしたコンセプトを処理する場合、内包関係のコンセプトおよび派生関係のコンセプトの語幹による手法がより有用な場合があります。

セマンティック・ネットワークの機能

セマンティック・ネットワーク手法は、既知の単語の関係を利用して、類義語または下位語のカテゴリーを作成します。下位語は、1 つのコンセプトがある種の 2 番目のコンセプトである場合、階層の関係性があり、ISA リレーションシップとも呼ばれます。例えば、animal がコンセプトである場合、動物の種類である cat、kangaroo は animal の下位語となります。

類義語および下位語の関係性のほか、セマンティック・ネットワークの手法では、<Location> タイプからコンセプト間の部分的なリンクおよび全体のリンクを検証します。例えば、この手法ではコンセプト normandy、provence、および france を、ノルマンディおよびプロバンスは、フランスの一部であるため、1 つのカテゴリーにグループ化します。

セマンティック・ネットワークは、セマンティック・ネットワークの各コンセプトの考えられる意味を特定することから始めます。コンセプトが類義語または下位語として特定されると、1 つのカテゴリーにグループ化されます。例えば、この手法を使うと、次の 3 つのコンセプトからなる 1 つのカテゴリーを作成されます。生食用リンゴ、デザートของリンゴ、およびグラニー・スミス。なぜならば、セマンティック・ネットワークには次のような情報が含まれるからです。1) デザートのリンゴは生食用リンゴの類義語であり、2) グラニー・スミスは生食用リンゴの一種である (生食用リンゴの下位語という意味)。

個別にみると、多くのコンセプト、特にユニタームがあいまいです。例えば、コンセプト buffet は食事の種類、あるいは家具を表す場合があります。一連のコンセプトに meal、furniture、および buffet がある

場合、アルゴリズムは meal または furniture のいずれかによる buffet のグループ化を選択するよう強制します。アルゴリズムによる選択は、レコードまたはドキュメントのコンテキストにおいては適切でない場合があります。

セマンティック・ネットワーク手法は、特定の種類のデータによる内包関係のコンセプトにおいて優れています。セマンティック・ネットワークと内包関係のコンセプトでは、apple pie が pie の一種であることを認識しますが、tart も pie の一種であることを認識できるのはセマンティック・ネットワークだけです。

セマンティック・ネットワークは、他の手法を組み合わせることで機能します。例えば、セマンティック・ネットワーク手法と内包関係のコンセプトの手法を選択し、セマンティック・ネットワークによりコンセプト teacher をコンセプト tutor とグループ化した (tutor は teacher の一種であるため) とします。内包アルゴリズムはコンセプトを graduate tutor と tutor にグループ分けし、結果として、2つのアルゴリズムが共同してアウトプット カテゴリーを作成します。アウトプットカテゴリーには、tutor, graduate tutor, and teacherが含まれます。

セマンティック・ネットワークのオプション

この手法では、さまざまな追加設定が重要である場合があります。

- 「最大検索距離」を変更します。カテゴリー作成前に手法による検索の距離を選択します。検索範囲を拡大すると、それぞれの共起の最低相似値が下がります。結果として、複数の共起規則が作成される場合がありますが、相似値の低いものは多くの場合さほど重要ではありません。

例えば、距離に応じて、Danish pastry から coffee roll (上位) まで、そして bun (祖父母) および bread まで上方に検索します。

作成されるカテゴリーが大きすぎる、あるいはあまりに多くのものがグループ化されていると感じられる場合は、検索距離を短縮すれば、より小さなカテゴリーを作成でき、作業がしやすくなります。

重要: 誤ったグループ化と行くと結果に大きな悪影響をおよぼす場合があるため、この手法を手法する場合は、オプション 語幹文字数が次の最小値以上のときにスペルを調整する (「抽出」ダイアログ・ボックスまたはノードの「エキスパート」タブで定義) を適用せず、Fuzzy Grouping を行うことをお勧めします。

共起規則

共起規則を使用すると、ドキュメントまたはレコードのセット内で強い関連を持つコンセプトを見つけ、グループ化することができます。ドキュメントおよびレコードでコンセプトが共に頻繁に見つかる場合、共起はおそらくカテゴリー定義の値のものである基底となる関連を反映します。新しいカテゴリーを作成、カテゴリーを拡張するために、または別のカテゴリー手法の入力として使用できる共起規則を作成します。レコードのあるセット内で頻繁に同時に現れ、他のレコードのいずれにも個別にあまり現れない場合、2つのコンセプトは強力に共起します。この手法を使用して、少なくとも数百のドキュメントまたはレコードを持つ大きなデータセットを使用して良い結果を作成することができます。

例えば、多くのレコードに price という単語と availability という単語が含まれている場合、これらのコンセプトを共起規則 (price & available) にグループ化することができます。別の例で、コンセプト peanut butter、jelly、sandwich が個別で現れるより頻繁に同時に現れる場合、これらはコンセプト共起規則 (peanut butter & jelly & sandwich) にグループ化されます。

重要: 以前のリリースでは、共起規則および類義語規則は、大カッコで囲まれていました。このリリースの場合、大カッコはテキスト リンク分析パターン結果を示します。その代わりに、共起規則や類義語規則は、(スピーカー・システム|スピーカー) のように、カッコで囲まれます。

共起規則の機能

この手法では、ドキュメントまたはレコードをスキャンし、同時に現れる傾向のある 2 つ以上のコンセプトを探します。ドキュメントまたはレコードのあるセット内で頻繁に同時に現れ、他のドキュメントまたはレコードのいずれにも個別にあまり現れない場合、2 つ以上のコンセプトは強力に共起します。

共起するコンセプトが見つかった場合、カテゴリ規則が形成されます。これらの規則は、& ブール型演算子を使用してつながっている 2 つ以上のコンセプトで構成されています。これらの規則は、規則内のコンセプトのセットがドキュメントまたはレコードですべて共起する場合、自動的にドキュメントまたはレコードをカテゴリに自動的に分類するという、論理ステートメントです。

共起規則のオプション

共起規則の手法を使用している場合、作成される規則に影響を与えるいくつかの設定を調整できます。

- 「最大検索距離」を変更します。共起規則の手法で共起を検索する距離を選択してください。検索距離を増やすと、各共起に必要な最小類似度値が低くなります。その結果、多数の共起規則が作成されますが、多くの場合、類似度値の低い共起規則の有意性は低くなります。検索距離を減らすと、必要とされる最小類似度値が高くなります。その結果、作成される共起規則は少なくなります。有意性が高く（より強力に）なります。
- 最小ドキュメント。コンセプトの特定のペアが共起として認識されるために、そのペアが含まれている必要があるレコードまたはドキュメントの最小数。このオプションの値を小さくすると、共起の検出数が大きくなります。この値を大きくすると、共起の検出数は少なくなります。共起の有意性は高くなります。例えば、「apple」と「pear」という 2 つのコンセプトが共に 2 件のレコードで見つかり、他のレコードではどちらのコンセプトも見つからないとします。「最小ドキュメント数」をデフォルト値の 2 に設定すると、共起規則手法によってカテゴリ規則（「apple and pear」）が作成されます。この値を 3 にすると、規則は作成されなくなります。

注: 小規模なデータ・セット（応答数が 1000 未満）の場合、デフォルト設定では共起が何も見つからない場合があります。その場合は、検索距離の値を大きくしてください。

注: コンセプトを明示的に指定することにより、コンセプトのグループ化を禁止することができます。詳しくは、110 ページの『例外ペアのリンクの管理』のトピックを参照してください。

出現頻度に基づく手法の詳細設定

単純で機械的な出現頻度に基づく手法に基づいて、カテゴリを作成できます。この手法を使用して、指定されたレコードまたはドキュメントの数を超えて見つかった項目（タイプ、コンセプト、またはパターン）ごとに 1 つずつカテゴリを作成できます。また、あまり頻繁に出現しない項目すべてに 1 つカテゴリを作成できます。度数ごとに、テキスト全体の出現数の合計に対し、対象の抽出されたコンセプト（およびその類義語）、タイプ、またはパターンを含むレコードまたはドキュメントの数を参照します。

一般的または重要な回答を示すため、頻繁に出現する項目をグループ化すると、関心のある結果を生成できます。他の手法が適用されると、この手法が未使用の抽出結果に非常に役立ちます。他にカテゴリが存在しない場合、別のアプリケーションが抽出直後にこの手法を実行し、結果を編集して、関心のないカテゴリを削除し、これらのカテゴリを拡張して、より多くのレコードまたはドキュメントに一致するようにします。詳しくは、116 ページの『カテゴリの拡張』のトピックを参照してください。

この手法を使用する代わりに、抽出結果ペインのレコードまたはドキュメント数の多い順にコンセプトまたはコンセプト・パターンを並べ替え。上位のコンセプトまたはコンセプト・パターンをカテゴリ・ペインにドラッグ・アンド・ドロップして、対応するカテゴリを作成することができます。

以下は、「詳細設定: 出現頻度」ダイアログ・ボックスにあるフィールドです。

カテゴリーの記述子を生成: 記述子の入力の種類を選択します。詳しくは、106 ページの『カテゴリーの作成』のトピックを参照してください。

- **コンセプト レベル:**このオプションを選択すると、コンセプトまたはコンセプト・パターンの度数が使用されます。タイプがカテゴリー作成の入力として選択されている場合はコンセプトが使用され、タイプ・パターンが選択された場合はコンセプト・パターンが使用されます。一般的には、この手法をコンセプト・レベルに適用すると、コンセプトおよびコンセプト・パターンがより低いレベルの尺度を示すため、より特定の結果を作成します。
- **タイプ レベル:**このオプションを選択すると、タイプまたはタイプ・パターンの度数が使用されます。タイプがカテゴリー作成の入力として選択されている場合はタイプが使用され、タイプ・パターンが選択された場合はタイプ・パターンが使用されます。この手法をタイプ レベルに適用すると、指定された情報の種類に関してすばやく表示することができます。

独自のカテゴリーを持つ項目の 最小ドキュメント 数。このオプションを使用すると、頻繁に出現する項目からカテゴリーを作成できます。このオプションは、出力をレコード またはドキュメントの少なくとも X 数に含まれる記述子を含むカテゴリーに制限します。ここで、X はこのオプションに入力する値です。

すべての残りの項目を次のカテゴリーにグループ化: このオプションを使用すると、あまり頻繁に出現しないすべてのコンセプトまたはタイプを、選択した名前の付いた「キャッチオール」カテゴリーにグループ化します。デフォルトでは、カテゴリー名はその他です。

カテゴリー入力: 手法を適用するグループを選択します。

- **未使用の抽出結果:**既存のカテゴリーで使用されていない抽出結果からカテゴリーを作成できます。レコードが、複数のカテゴリーと合致する傾向が最も小さくなり、作成されるカテゴリーの数が制限されません。
- **すべての抽出結果:**抽出結果のいずれを使用してもカテゴリーを作成できます。カテゴリーがないまたは少ない場合に最も役立ちます。

次の方法で重複するカテゴリー名を解決: 名前が既存のカテゴリーと同じ新規カテゴリーまたはサブカテゴリーの処理方法を選択します。新規カテゴリー (およびその記述子) と同じ名前を持つ既存カテゴリーとを結合できます。あるいは、既存のカテゴリーに重複する名前があった場合、カテゴリーの作成をスキップすることもできます。

カテゴリーの拡張

拡張は、記述子を自動的に追加または拡張して、既存のカテゴリーを「拡大」するプロセスです。その目的は、本来カテゴリーに割り当てられていなかった関連レコードまたはドキュメントをキャプチャーするより良いカテゴリーを作成することです。

選択した自動グループ化手法では、既存のカテゴリー記述子に関連するコンセプト、TLA パターン、およびカテゴリー規則を特定しようとします。これらの新しいコンセプト、パターン、カテゴリー規則が新しい記述子として追加されるか、既存の記述子に追加されます。拡張するためのグループ化手法には、「派生関係のコンセプトの語幹」、「内包関係のコンセプト」、「セマンティック・ネットワーク」(英語のみ)、および「共起規則」が含まれます。「カテゴリー名から生成された記述子を使用して空白のカテゴリーを拡張する」の手法を使用すると、カテゴリー名の単語を使用して記述子を生成します。そのため、カテゴリー名が記述的であるほど、結果が良いものとなります。

注: カテゴリーを拡張する場合、出現頻度に基づく手法は使用できません。

拡張は、カテゴリをインタラクティブに改善する重要な方法です。次に、カテゴリを拡張する場合の例をいくつか示します。

- カテゴリ・ペインでコンセプト・パターンをドラッグ/ドロップしてカテゴリを作成した後
- 手動でカテゴリを作成し、簡単なカテゴリ規則および記述子を追加した後
- 非常に記述的な名前を持つ 事前定義済みカテゴリ・ファイルをインポートした後
- 選択した TAP に由来するカテゴリを修正した後

カテゴリを複数回使用できます。例えば、非常に記述的な名前を持つ事前定義済みカテゴリ・ファイルをインポートした場合、「カテゴリ名から生成された記述子を使用して空白のカテゴリを拡張する」オプションを使用して拡張子、記述子の最初のセットを取得して、これらのカテゴリを再度拡張します。ただし、複数回拡張すると、記述子が拡張されて幅広くなると、あまりに一般的なカテゴリが生成される場合があります。作成グループ化手法および拡張グループ化手法では類似した基底のアルゴリズムを使用するため、カテゴリの作成後に直接拡張すると、より関心の高い結果の作成は期待できません。

ヒント:

- 拡張を試みるが結果の使用は望まない場合、拡張を行った直後に操作をいつでも取り消す（「編集」>「取り消し」）ことができます。
- プロセス中、規則は個別に作成されるため、ドキュメントの同じセットに正確に一致するカテゴリのカテゴリ規則を 2 つ以上作成します。必要に応じて、カテゴリを確認し、カテゴリ記述子を手動で編集して重複を削除できます。詳しくは、139 ページの『カテゴリ記述子の編集』のトピックを参照してください。

カテゴリを展開するには

1. カテゴリ・ペインで、展開するカテゴリを選択します。
2. メニューの「カテゴリ」>「カテゴリを展開」を選択します。プロンプトが表示されないようオプションを選択している場合、メッセージ・ボックスが表示されます。
3. 今すぐ作成するか、左記に設定を編集するかを選択します。
 - 「今すぐ拡張」をクリックすると、現在の設定でカテゴリの拡張が開始されます。プロセスが開始し、進捗状況のダイアログが表示されます。
 - 「編集」をクリックして、設定を確認し、変更します。

拡張しようとした後、新しい記述子が見つかったカテゴリには、カテゴリ・ペインで「展開」という単語のフラグが立てられ、すばやくカテゴリを特定できます。「展開」というテキストは、再度展開するか、別の方法で編集、またはコンテキスト・メニューを使用してこれらを解除するまで残ったままです。

注: 表示できる最大カテゴリ数は 10,000 です。この数に到達したか超過した場合は警告が表示されます。その場合は、「カテゴリの作成」オプションまたは「カテゴリの拡張」オプションを変更して作成カテゴリの数を少なくする必要があります。

カテゴリの作成時または拡張時に使用できるそれぞれの手法は、特定の種類のデータおよび状況に適していますが、同じ分析で手法を組み合わせ、全範囲のドキュメントまたはレコードをキャプチャーすると役に立つ場合があります。インタラクティブ・ワークベンチで、カテゴリにグループ化されたコンセプトおよびタイプは、次にカテゴリを作成する場合も使用できます。つまり、複数のカテゴリのコンセプトを表示したり、重複するカテゴリを見つけることができます。

以下は、「カテゴリの作成: 設定」ダイアログ・ボックスにある領域とフィールドです。

次による拡張。カテゴリの展開に使用する入力を選択します。

- 未使用の抽出結果:既存のカテゴリーで使用されていない抽出結果からカテゴリーを作成できます。レコードが、複数のカテゴリーと合致する傾向が最も小さくなり、作成されるカテゴリーの数が制限されます。
- すべての抽出結果:抽出結果のいずれを使用してもカテゴリーを作成できます。カテゴリーがないまたは少ない場合に最も役立ちます。

グループ化手法

これらの手法の簡単な説明は、「108 ページの『言語学的手法の詳細設定』」を参照してください。これらの手法には、次のものが含まれています。

- 派生関係のコンセプトの語幹
- セマンティック ネットワーク (英語テキストのみで、「一般化のみ」オプションが選択されている場合は使用されません。)
- 内包関係のコンセプト
- 共起および最小ドキュメント数のサブオプション

これらのタイプは関連する結果を作成しないため、多くのタイプがセマンティック・ネットワークから永続的に除外します。それらのタイプには、<Positive>、<Negative>、<IP>、その他の非言語的タイプなどがあります。

最大検索距離: カテゴリー作成前に手法による検索の距離を選択します。ただし、これらの結果はノイズが少なく、またリンクや関連性が大きくなります。値が大きいほど、取得する結果は多くなります。ただし、これらの結果の信頼性または関連性が弱くなります。このオプションはすべての手法にグローバルに適用されますが、共起とセマンティック・ネットワークに対する効果は最も大きくなります。

特定のコンセプトがグループ化されないようにする: 出力の 2 つのコンセプトがグループ化またはペアとしないように処理を停止します。コンセプト・ペアを作成または管理するには、「ペアを管理」をクリックします。詳しくは、110 ページの『例外ペアのリンクの管理』のトピックを参照してください。

可能な場合、単純に拡張するか、ワイルドカードを使用して記述子を一般化するか、またはその両方を選択します。

- **拡張および一般化:** このオプションは、選択したカテゴリーを拡張し、記述子を一般化します。一般化を選択した場合、アスタリスク (*) のワイルドカードを使用して、製品が一般的なカテゴリー規則をカテゴリーに作成することができるようになります。例えば、[アップル タルト + .] や [アップル ソース + .] などの複数の記述子を作成する代わりに、[アップル * + .] のようにワイルドカードを使用します。ワイルドカードを使用して一般化すると、以前と同じように、ちょうど同じ数のレコードまたはドキュメントを取得する場合があります。ただし、このオプションには、数の縮小やカテゴリーの記述子の簡略化という利点があります。また、このオプションを使用すると、新しいテキスト・データ (例: 長期的/周期的研究) にこれらのカテゴリーを使用してより多くのレコードまたはドキュメントをカテゴリー化する機能を拡大します。
- **拡張のみ:** 一般化せずにカテゴリーを展開します。手動で作成したカテゴリーには「拡張のみ」オプションを選択し、「拡張および一般化」オプションを使用して同じカテゴリーをもう一度展開すると便利です。
- **一般化のみ:** 別の方法でカテゴリーを展開せずに、記述子を一般化します。

注: このオプションを選択すると、「セマンティック・ネットワーク」オプションが無効になります。これは、「セマンティック・ネットワーク」オプションは説明が拡張されるときのみ有効になるためです。

カテゴリーを拡張するその他のオプション

適用するグループ化手法を選択するほか、次のように、その他のオプションを編集することができます。

記述子を拡張する場合の最大項目数: 項目 (コンセプト、タイプおよびその他の式) で記述子を拡張する場合、単一の記述子に追加できる項目の最大数を定義します。この制限値を 10 に設定すると、最大 10 件の追加項目を既存の記述子に追加できます。10 件を超える項目を追加しようとする場合、10 番目の項目が追加されると、新しい項目の追加を停止します。そうすることにより、記述子のリストが短くなりますが、最も関心の高い項目が最初に使用されたことを保障するものではありません。「可能な場合ワイルドカードを使用して一般化」 オプションを使用して、品質を落とすことなく拡張のサイズを縮小することが必要な場合があります。このオプションは、ブール値 & (AND) または ! (NOT) を含む記述子にのみ適用されます。

サブカテゴリーも展開: 選択したカテゴリー下のサブカテゴリーも展開します。

カテゴリー名から生成された記述子を使用して空白のカテゴリーを拡張する: 記述子が 0 件の、空のカテゴリーにのみ適用されます。カテゴリーにすでに記述子が含まれている場合、この方法では拡張されません。このオプションを選択すると、カテゴリー名を構成する単語に基づいて、各カテゴリーの記述子を自動的に作成しようとします。カテゴリー名をスキャンして、名前の単語が抽出されたコンセプトに一致するかどうかを確認します。コンセプトが認識されると、そのコンセプトを使用して、合致するコンセプト・パターンを検索し、コンセプトとパターンを使用してカテゴリーの記述子を形成します。このオプションを選択すると、カテゴリー名が長く記述的である場合に、最良の結果を作成します。迅速にカテゴリーの記述子を生成し、またカテゴリーはこれらの記述子を含むレコードをキャプチャーすることができます。別の場所からカテゴリーをインポートしたり、長く記述的な名前を使用して手動でカテゴリーを作成する場合に最も役立つオプションです。

記述子を次の形式で生成: このオプションは、先行のオプションがオンの場合のみ適用されます。

- **コンセプト:** 入力テキストから抽出されているかどうかに関係なく、記述子をコンセプトの形式で作成します。
- **パターン:** パターンが抽出されているかどうかに関係なく、記述子をパターンの形式で作成します。

手作業でのカテゴリーの作成

自動カテゴリー作成手法および条件規則エディターを使用してカテゴリーを作成するほか、カテゴリーを手動で作成することもできます。手動で作成する方法は、次のとおりです。

- 要素を 1 つずつ追加する空のカテゴリーを作成する。詳しくは、『カテゴリーの新規作成または名前の変更』のトピックを参照してください。
- キーワード、タイプ、パターンを「カテゴリー」ウィンドウにドラッグする。詳しくは、120 ページの『ドラッグ&ドロップによるカテゴリーの作成』のトピックを参照してください。

カテゴリーの新規作成または名前の変更

空のカテゴリーを作成して、コンセプトおよびタイプをカテゴリーに追加できます。カテゴリーの名前を変更することもできます。

新規の空白カテゴリーを作成するには

1. カテゴリー・ペインに移動します。
2. メニューから、「カテゴリー」>「空白カテゴリーを作成」を選択します。「カテゴリーのプロパティ」ダイアログボックスが開きます。

3. 「名前」フィールドでそのカテゴリーの名前を入力します。
4. 「OK」をクリックすると名前が適用され、ダイアログ・ボックスが閉じます。ダイアログ・ボックスが閉じ、新しいカテゴリー名がパネルに表示されます。

これでこのカテゴリーに追加していくことができます。詳しくは、139 ページの『記述子のカテゴリーへの追加』のトピックを参照してください。

カテゴリーの名前を変更するには

1. カテゴリーを選択して、「カテゴリー」>「カテゴリー名を変更」を選択します。「カテゴリーのプロパティ」ダイアログボックスが開きます。
2. 「名前」フィールドでそのカテゴリーの新しい名前を入力します。
3. 「OK」をクリックすると名前が適用され、ダイアログ・ボックスが閉じます。ダイアログ・ボックスが閉じ、新しいカテゴリー名がパネルに表示されます。

ドラッグ&ドロップによるカテゴリーの作成

ドラッグ&ドロップは手作業の手法であって、アルゴリズムに基づいたものではありません。次のように、「カテゴリー」ウィンドウでカテゴリーを作成できます。

- コンセプト、タイプ、またはパターンを抽出結果ペインからカテゴリー・ペインにドラッグ。
- 抽出されたコンセプトをデータ・ペインからカテゴリー・ペインにドラッグ。
- 行全体をデータ・ペインからカテゴリー・ペインにドラッグ。その行に含まれる抽出されたすべてのコンセプトおよびパターンで構成されたカテゴリーを作成します。

注:抽出結果ペインでは、複数選択をサポートし、複数要素のドラッグ・アンド・ドロップを可能にします。

重要: テキストから抽出されていないコンセプトをデータ・ペインからドラッグ・アンド・ドロップできません。データで検出したコンセプトの抽出を強制する場合、そのコンセプトをタイプに追加する必要があります。その後、抽出を再度実行します。新しい抽出結果には、追加したばかりのコンセプトが含まれません。その結果をカテゴリーに使用できます。詳しくは、91 ページの『コンセプトのタイプへの追加』のトピックを参照してください。

ドラッグ&ドロップでカテゴリーを作成するには：

1. 抽出結果ペインまたはデータ・ペインから、1 つまたは複数のコンセプト、パターン、タイプ、レコード、または一部のレコードを選択します。
2. マウス・ボタンを押したまま、要素を既存のカテゴリーまたはパネルの領域にドラッグして、新しいカテゴリーを作成します。
3. 要素をドロップする領域の上で、マウス・ボタンを離します。要素がカテゴリー・ペインに追加されます。変更が加えられたカテゴリーは背景色が変わります。この色は、「カテゴリー フィールドバック背景」と呼ばれます。詳しくは、76 ページの『オプションの設定』のトピックを参照してください。

注: 結果として得られたカテゴリーには、自動的に名前が付けられます。詳細は、7 章p.114 "カテゴリーのプロパティの編集" を参照してください。詳しくは、119 ページの『カテゴリーの新規作成または名前の変更』のトピックを参照してください。

どのレコードがカテゴリーに割り当てられているかを確認する場合は、カテゴリー・ペインで該当するカテゴリーを選択します。データ・ペインは自動的に更新され、そのカテゴリーのすべてのレコードが表示されます。

カテゴリ規則の使用

カテゴリを作成するには、さまざまな方法があります。それらの方法の 1 つには、キーワードを表すカテゴリ規則を定義することがあります。カテゴリ規則とは、抽出したコンセプト、タイプ、パターン、およびブール型演算子を使用した論理式に基づいて、ドキュメントまたはレコードを自動的に分類するステートメントです。例えば、「このカテゴリに、アルゼンチンではなく、大使館という抽出したコンセプトを含む」という内容を意味する式を作成することができます。

カテゴリ規則の中には、共起 および 派生関係のコンセプトの語幹 などのグループ化手法を用いたカテゴリ作成 (カテゴリ > カテゴリ作成設定 > 詳細設定: 言語学的手法)時に、自動的に作成されるものもありますが、他方、データやコンテキストの独自のカテゴリ理解に従って、カテゴリ・エディターを使用して手動でカテゴリ規則を作成することもできます。各規則は単一のカテゴリに関連付けられるため、規則を満たすドキュメントまたはレコードがそのカテゴリにスコアリングされます。

カテゴリ規則により、より大きな特異性を持つ回答をカテゴリ化できることで、テキストマイニングの結果およびより高度な数量分析の品質および生産性を向上させます。経験およびビジネス情報により、データおよびコンテキストに対する特定の理解を与える場合があります。このように理解することによって、情報をカテゴリに変換し、ブール型ロジックと抽出した要素を結合して、ドキュメントまたはレコードをより効率的かつ正確にカテゴリ化できます。

これらの規則を作成する機能により、ビジネス情報を製品の抽出テクノロジーに重ねて、コード化の精度、効率性および生産性を拡張することができます。

注: 条件規則がどのようにテキストに合致するかについての例は、「127 ページの『カテゴリ規則の例』」を参照してください。

カテゴリ規則シンタックス

カテゴリ規則の中には、共起 および 派生関係のコンセプトの語幹 などのグループ化手法を用いたカテゴリ作成 (カテゴリ > カテゴリ作成設定 > 詳細設定: 言語学的手法)時に、自動的に作成されるものもありますが、他方、カテゴリ・エディターで手動でカテゴリ規則を作成することもできます。各規則は単一のカテゴリの説明であるため、規則を満たすドキュメントまたはレコードが自動的にそのカテゴリにスコアリングされます。




注: 条件規則がどのようにテキストに合致するかについての例は、「127 ページの『カテゴリ規則の例』」を参照してください。

規則を作成または編集する場合、規則を条件規則エディターで開く必要があります。コンセプト、タイプ、またはパターンを追加し、またワイルドカードを使用して一致を拡張することができます。抽出されたコンセプト、タイプ、パターンを使用すると、すべての関連コンセプトを検出することができ便利です。

重要: 一般的なエラーを回避するために、コンセプトを抽出結果ペイン、テキストリンク分析ペインまたはデータ・ペインから直接条件規則エディターにドラッグアンドドロップしたり、使用できる場合はコンテキスト・メニューを使用して追加することを推奨します。

コンセプト、タイプ、パターンが認識されると、テキストの隣にアイコンが表示されます。

表 18. 抽出アイコン

アイコン	説明
	抽出コンセプト
	抽出タイプ
	抽出パターン

条件規則シンタックスおよび演算子

次の表には、条件規則シンタックスを定義する文字を示しています。これらの文字をコンセプト、タイプ、パターンと共に使用して、規則を作成します。

表 19. サポートされるシンタックス

文字	説明
&	「and」ブール型演算子。例えば、次のように a & b には、a と b が含まれます。 - 侵略 & アメリカ合衆国 - 2016 & オリンピック - 良い & リンゴ
	「or」ブール型演算子は包含的演算子で、要素の一部またはすべてが見つかった場合に一致することを意味します。例えば、次のような a b には、a または b が含まれます。 - 攻撃 フランス - コンドミニウム アパート
!()	「not」ブール型演算子。例えば、次のような !(a) には、a が含まれません。 !(良い & ホテル) 、暗殺 & !(オーストリア)、または !(金) & !(銅)
*	使用方法に応じて、1 文字から単語全体にいたるまでの文字を示すワイルドカード 詳しくは、125 ページの『カテゴリ規則におけるワイルドカードの使用』のトピックを参照してください。
()	式の区切り文字。カッコ内の式を最初に評価します。
+	順序特有のパターンを形成するために使用するパターン・コネクタ。この演算子がある場合は、大カッコを使用する必要があります 詳しくは、123 ページの『カテゴリ規則内の TLA パターンの使用』のトピックを参照してください。
[]	カテゴリ規則内部の抽出した TLA パターンに基づいて合致を検出する場合、パターン区切り文字が必要です。ブラケット内の内容は TLA パターンを参照し、単純な共起に基づいてコンセプトまたはタイプに合致しません。この TLA パターンを抽出していない場合、合致は発生しません。詳しくは、123 ページの『カテゴリ規則内の TLA パターンの使用』のトピックを参照してください。パターンではなくコンセプトとタイプの合致に着目している場合は、角カッコを使用しないでください。 注: 古いバージョンの場合、カテゴリ作成手法で生成された共起規則や類義語規則は、大カッコに囲まれていました。すべての新しいバージョンの場合、大カッコは TLA パターンの存在を示します。その代わりに、共起規則による手法や類義語規則による手法で作成された規則は、(スピーカー・システム スピーカー) のように、カッコで囲まれます。

& 演算子および | 演算子は、a & b = b & a および a | b = b | a のように相互的です。

バックスラッシュによる文字のエスケープ

シンタックス文字でもある文字がコンセプトに含まれている場合、その文字の前にバックスラッシュを追加して、規則が正しく解釈されるようにする必要があります。バックスラッシュ (\) 文字を使用して、特別な意味を持つ文字をエスケープします。エディターにドラッグ・アンド・ドロップすると、自動的にバックスラッシュが追加されます。

次の条件規則シンタックス文字を条件規則シンタックスとしてでなく、そのまま扱う場合は、その文字の前にバックスラッシュを追加する必要があります。

& ! | + < > () [] *

例えば、コンセプト `r&d` に「and」演算子 (&) が使用されているため、条件規則エディターで入力する場合は、`r¥d` となるようにバックスラッシュが必要です。

カテゴリ規則内の TLA パターンの使用

テキスト リンク分析パターンを、カテゴリ規則で明示的に定義し、より具体的で文脈上の結果を取得することができます。カテゴリ規則でパターンを定義すると、より単純な抽出結果を省略し、抽出されたテキスト リンク分析パターン結果に基づいたドキュメントおよびレコードのみを合致させます。

重要: カテゴリ規則で TLA パターンを使用してドキュメントを合致させる場合、テキスト リンク分析を有効にして、抽出を実行する必要があります。カテゴリ規則では、そのプロセス時に検出される合致を検索します。テキスト マイニング・ノードの「モデル」タブで TLA 結果の探索を選択していない場合、インタラクティブ・セッションの抽出設定で TLA 抽出を有効にして、再抽出することができます。詳しくは、82 ページの『データの抽出』のトピックを参照してください。

大カッコで区切る。 カテゴリ規則の中で TLA パターンを使用している場合、TLA パターンを大カッコ [] で囲む必要があります。抽出した TLA パターンに基づいて合致を検出する場合、パターン区切り文字が必要です。カテゴリ規則には、タイプ、コンセプト、またはパターンが含まれるため、カッコはカッコ内の内容が抽出された TLA パターンを参照するということを明確にします。この TLA パターンを抽出していない場合、合致は発生しません。カテゴリ・ペインにリンゴ + 良い のようにカッコのないパターンがあった場合、パターンがカテゴリ規則エディター外部でカテゴリに直接追加されたことを示します。例えば、コンセプト・パターンをテキスト リンク分析ビューからカテゴリに直接追加する場合、大カッコは表示されません。ただし、規則内でパターンを使用する場合、[バナナ + !(良い)] のようにカテゴリ規則内の大カッコで囲む必要があります。

パターンで + 記号を使用。 IBM SPSS Modeler Text Analyticsでは、最大 6 つの部分 (スロット) のパターンを作成できます。順序が重要であることを示す場合、[会社1 + 買収 + 会社2] のように、+ 記号を使用して、要素を接続します。ここでは、どの企業が買収するかの意味が変わってしまうため、順序が重要となります。文の構造ではなく、TLA パターン出力の構造がどのようになっているのかによって順序が決まります。例えば、「I love Paris」というテキストがあり、このキーワードを抽出する場合、デフォルトの意見リソースは通常、意見を 2 つの部分で構成されるパターンの 2 番目の位置に配置するため、TLA パターンは [<Positive> + <Location>] ではなく、[paris + like] または [<Location> + <Positive>] となります。そのため、問題を回避するためにパターンを直接記述子として使用すると役に立つ場合があります。ただし、より複雑な表現の一部としてパターンを使用する必要がある場合、合致が検出されるかどうかにおいて、順序が大きな役割を果たすため、テキスト リンク分析ビューのパターン内の要素の順序には特に注意をしてください。

例えば、"I like pineapple" というテキストと、 "I hate pineapple. However, I like strawberries" という 2 つのサンプル テキストがあったとします。表現 like & pineapple はそれはコンセプト式でありテキスト

リンク規則ではない (角カッコでくくられていない) ため、両方のテキストに合致します。 pineapple + like 式は "I like pineapple" とのみ一致します。なぜならば、2 番目のテキストでは、like という語は今度 は strawberries 関連づけられるからです。

パターンによりグループ化: 独自のパターンを使用して規則を簡略化することができます。3 つの式、 cayenne peppers + like、chili peppers + like、および peppers + like の式をキャプチャーするとします。これらを、[* peppers & like] のようにして、単一のカテゴリ規則にグループ化できます。hot peppers + good というもう 1 つの式がある場合、これら 4 つの式を [* peppers + <Positive>] のような規則でグループ化することができます。

パターン内の順序: 出力をよりよく構成するために、製品とともにインストールされたテンプレートに提供されているテキスト リンク分析規則が、文の語順に関係なく、同じ順序で基本パターンを出力しようとします。例えば、テキスト「Good presentations.」を含むレコードおよび「the presentations were good」を含む別のレコードがある場合、いずれのテキストも同じ規則で合致し、出力の順序は、presentation + good や good + presentation ではなく、コンセプトパターン規則のpresentation + good と同じ順序になります。また、例のパターンのような 2 スロット・パターンで、意見ライブラリーのタイプに割り当てられたコンセプトは、デフォルトでは apple + bad のように出力の最後に表示されます。

表 20. パターン・シンタックスおよびブール型演算子の使用

式	ドキュメントまたはレコードが一致する条件
[]	任意の TLA パターンを含む。抽出した TLA パターンに基づいて合致を検出する場合、カテゴリ規則内でパターン区切り文字が必要です。ブラケット内の内容は、単なるコンセプトおよびタイプではなく TLA パターンを参照します。この TLA パターンを抽出していない場合、合致は発生しません。 そのため、パターンを含んでいない規則を作成する場合、!([]) を使用することができます。
[a]	パターン内の位置に関係なく、少なくとも 1 つの要素が a であるパターンを含む。例えば、[deal] は、[deal + good] または [deal + .] に一致します。
[a + b]	コンセプト・パターンを含む。例えば、[deal + good] となります。 注: 他の要素を追加せずにこのパターンをキャプチャーする場合、パターンを使用して規則を作成するのではなく、パターンを直接カテゴリに追加することをお勧めします。
[a + b + c]	コンセプト・パターンを含む。+ 記号は、一致する要素の順序は重要ではありません。例えば、[会社1 + 買収 + 会社2] となります。
[<A> +]	最初のスロットがタイプ <A>、2 番目のスロットがタイプ のパターンを含み、ちょうど 2 つのパターンがある。+ 記号は、一致する要素の順序は重要ではありません。例えば、[<Budget> + <Negative>] となります。 注: 他の要素を追加せずにこのパターンをキャプチャーする場合、パターンを使用して規則を作成するのではなく、パターンを直接カテゴリに追加することをお勧めします。
[<A> &]	タイプ <A> および のタイプ・パターンを含む。例えば、[<Budget> & <Negative>] となります。この TLA パターンは抽出されませんが、そのように記述された場合、[<Budget> + <Negative>][<Negative> + <Budget>] のようになります。+ 記号は、一致する要素の順序は重要ではありません。また、他の要素がパターン内にある場合がありますが、少なくとも <Budget> and <Negative> があります。

表 20. パターン・シンタックスおよびブール型演算子の使用 (続き)

式	ドキュメントまたはレコードが一致する条件
[a + .]	a が唯一のコンセプトであるパターンを含み、そのパターンの他のスロットには何もありません。 例を次に示します。 [deal + .] は、唯一の出力がコンセプト deal であるコンセプト・パターンに一致します。コンセプト deal をカテゴリ記述子として追加した場合、deal を含むすべてのレコードを、deal に関して肯定的な記述を含むコンセプトとして取得します。ただし、[deal + .] を使用すると、deal を示すこれらのレコード・パターン結果のみに合致し、他の関係性または意見は deal + fantastic と合致しません。 注: 他の要素を追加せずにこのパターンをキャプチャーする場合、パターンを使用して規則を作成するのではなく、パターンを直接カテゴリに追加することをお勧めします。
[<A> + <>]	<A> が唯一のタイプであるパターンが含まれている。例えば、[<Budget> + <>] は、唯一の出力がタイプ <Budget> のコンセプトであるパターンに一致します。 注: [price + <>] ではなく、[<Budget> + <>] のように、タイプ・パターンでパターンの + 記号の後に使用する場合にのみ、<> を使用して空のタイプを示すことができます。 注: 他の要素を追加せずにこのパターンをキャプチャーする場合、パターンを使用して規則を作成するのではなく、パターンを直接カテゴリに追加することをお勧めします。
[a + !(b)]	コンセプト a を含み、コンセプト b を含まないパターンを少なくとも 1 つ含む。少なくとも 1 つのパターンを含む必要があります。 例えば、[price + !(high)]、 またはタイプには [!(<Fruit> <Vegetable>) + <Positive>] になります
![(<A> &)]	特定のパターンを含まない。例えば、![(<Budget> & <Negative>)] となります。

注: 条件規則がどのようにテキストに合致するかについての例は、「127 ページの『カテゴリ規則の例』」を参照してください。

カテゴリ規則におけるワイルドカードの使用

ワイルドカードを規則のコンセプトに追加して、マッチング機能を拡張することができます。アスタリスク * ワイルドカードを単語の前および/または後に追加して、コンセプトがどのように一致するかを指定できます。ワイルドカードの使用方法には次の 2 種類があります。

- ワイルドカードの接頭 (尾) 辞: ワイルドカードが接頭辞または接尾辞として使用されます。文字列とアスタリスクの間にスペースはありません。例えば、operat* は、operat、operate、operates、operations、operational などと一致します。
- 単語ワイルドカード: ワイルドカードがコンセプトの前または後ろに使用されます。文字列とアスタリスクの間にスペースがあります。例えば、* operation は、operation、surgical operation、post operation などと一致します。例えば、* operat* * のように、ワイルドカードが接頭 (尾) 辞とともに単語ワイルドカードを使用できます。この場合、operation、surgical operation、telephone operator、operatic aria などと一致します。最後の例のように、範囲があまりに幅広くなったり、不要なマッチをキャプチャーしないよう、ワイルドカードを注意して使用することをお勧めします。

例外:

- ワイルドカードは、単独で使用できません。例えば、(apple | *) は無効です。
- ワイルドカードをタイプ名の一致に使用することはできません。<Negative*> は、タイプ名に合致しません。
- 特定のタイプをワイルドカードで検索されたコンセプトに対する合致から除外することはできません。コンセプトが割り当てられるタイプは自動的に使用されます。

- ワイルドカードは、語の終わりであっても初めであっても (open* account)、または独立した要素であっても (open * account) 語の連鎖の途中に使用することはできません。タイプ名にワイルドカードを使用することはできません。例えば、apple* recipe など、word* word は「applesauce recipe」や他の言葉にも合致しません。ただし、apple* * は、applesauce recipe、apple pie、apple などの言葉に合致します。また、apple * toast など、word* word は 2 つの語の間にアスタリスクを使用しているため、apple cinnamon toast という語には合致しません。ただし、apple* * は、apple cinnamon toast、apple、apple pie などに合致します。

表 21. ワイルドカードの使用方法

式	ドキュメントまたはレコードが一致する条件
*apple	文字で終了し、接頭辞として任意の数の文字を使用しているコンセプトを含む。例:*apple は、apple で終了し、次のように接頭辞を使用します。 - apple - pineapple - crabapple
apple*	文字で開始し、接尾辞として任意の数の文字を使用しているコンセプトを含む。例:*apple は、apple で開始し、次のように接尾辞を使用、または使用しません。 - apple - applesauce - applejack 例えば、apple* & !(pear* quince) には、文字 apple で始まるコンセプトが含まれますが、文字 pear またはコンセプト quince で始まるコンセプトは含まれません。そのため、次のコンセプトとは一致しません。apple & quince 次のコンセプトとは一致します。 - applesauce - apple & orange
product	product という文字を含み、接頭辞または接尾辞、または両方として任意の数の文字を使用しているコンセプトを含む。 例:*product* は、次のコンセプトと一致します。 - product - byproduct - unproductive
* loan	単語 loan を含み、単語の前に別の単語と組み合わせる場合があるコンセプトを含む。例えば、* loan は次のコンセプトと一致します。 - loan - car loan - home equity loan 例えば、[* delivery + <Negative>] は、前半が単語 delivery で終わり、後半にタイプ <Negative> を含むコンセプトを含み、次のコンセプト・パターンと一致します。 - package delivery + slow - overnight delivery + late
event *	単語 event を含み、単語の後に別の単語が続く場合があるコンセプトを含む。例えば、event * は次のコンセプトと一致します。 - event - event location - event planning committee

表 21. ワイルドカードの使用法 (続き)

式	ドキュメントまたはレコードが一致する条件
* apple *	<p>任意の単語で始まり、次に apple で始まる単語が続き、別の単語が続く場合があるコンセプトを含む。* は 0 または n を意味するため、apple にも合致します。例えば、* apple* は次のコンセプトと一致します。</p> <ul style="list-style-type: none"> - gala applesauce - granny smith apple crumble - famous apple pie - apple <p>例えば、[* reservation* * + <Positive>] は、前半が単語 reservation のあるコンセプトを含み (コンセプト内の場所は関係ない)、後半にタイプ <Positive> を含み、次のコンセプト・パターンと一致します。</p> <ul style="list-style-type: none"> - reservation system + good - online reservation + good

注: 条件規則がどのようにテキストに合致するかについての例は、「『カテゴリー規則の例』」を参照してください。

カテゴリー規則の例

それらを表現するために使用するシンタックスに基づき、規則がレコードにどのように合致するかを示すために、次の例について考えてみます。

例のレコード

次の 2 つのレコードがあるとします。

- レコード A: *"when I checked my wallet, I saw I was missing 5 dollars."*
- レコード B: *"\$5 was found at the picnic area, but the blanket was missing."*

次の 2 つの表には、コンセプト・パターンおよびタイプ・パターンのほかにコンセプトおよびタイプに抽出される内容を示しています。

例から抽出されるコンセプトとタイプ

表 22. コンセプトとタイプの抽出例

抽出コンセプト	コンセプトのタイプ
wallet	<Unknown >
missing	<Negative>
USD5	<通貨>
blanket	<Unknown >
picnic area	<Unknown >

例から抽出される TLA パターン

表 23. TLA パターン出力の抽出例

抽出されたコンセプト・パターン	抽出されたタイプパターン	レコード
picnic area + .	<不明> + <>	レコード B
wallet + .	<不明> + <>	レコード A

表 23. TLA パターン出力の抽出例 (続き)

抽出されたコンセプト・パターン	抽出されたタイプパターン	レコード
blanket + missing	<不明> + <否定的>	レコード B
USD5 + .	<通貨> + <>	レコード B
USD5 + missing	<通貨> + <否定的>	レコード A

カテゴリ規則の合致

次の表には、カテゴリ規則エディターに入力できるシンタックスをいくつか示しています。すべての規則が機能するわけではなく、またすべてが同じレコードに合致するわけではありません。異なるシンタックスが合致したレコードにどのように影響するかを確認してください。

表 24. 条件規則のサンプル

条件規則シンタックス	結果
USD5 & missing	抽出コンセプト missing および抽出コンセプト USD5 の両方が含まれているため、レコード A および B のいずれにも合致します。これは以下と同じになります。 (USD5 & missing)
missing & USD5	抽出コンセプト missing および抽出コンセプト USD5 の両方が含まれているため、レコード A および B のいずれにも合致します。これは以下と同じになります。 (missing & USD5)
missing & <通貨>	抽出コンセプト missing およびタイプ <通貨> に合致するコンセプト両方が含まれているため、レコード A および B のいずれにも合致します。これは以下と同じになります。 (missing & <通貨>)
<通貨> & missing	抽出コンセプト missing およびタイプ <通貨> に合致するコンセプト両方が含まれているため、レコード A および B のいずれにも合致します。これは以下と同じになります。 (<通貨> & missing)
[USD5 + missing]	レコード B は USD5 + missing を含む TLA パターン出力を作成していないため、A とは合致しますが B には合致しません (前の表を参照)。これは次の TLA パターン出力と同じになります。 USD5 + missing
[missing + USD5]	抽出された TLA パターン (前の表参照) が、最初の位置に missing を使用して表現した順序に合致しないため、レコード A にも B にも合致しません。これは次の TLA パターン出力と同じになります。 USD5 + missing
[missing & USD5]	レコード B からそのような TLA パターンが抽出されていないため、A とは合致しますが B とは合致しません。& 文字を使用すると、合致時の順序が重要でないことを示すため、この規則では [missing + USD5] または [USD5 + missing] のパターン マッチのいずれかを検索します。レコード A の [USD5 + missing] のみに合致があります。
[missing + <通貨>]	抽出された TLA パターンがこの順序に合致していないため、レコード A にも B にも合致しません。TLA 出力はキーワード (USD5 + missing) またはタイプ (<通貨> + <Negative>) にも基づくため、同等のものはありませんが、コンセプトおよびタイプを組み合わせません。

表 24. 条件規則のサンプル (続き)

条件規則シンタックス	結果
[<通貨> + <否定的>]	TLA パターンがレコード B から抽出されていないため、レコード A に合致しますが B には合致しません。以下の TLA 出力と同じになります。 <通貨> + <否定的>
[<否定的> + <通貨>]	抽出された TLA パターンがこの順序に合致していないため、レコード A にも B にも合致しません。「意見」テンプレートの場合、デフォルトではトピックが意見とともに検出されると、トピック (<Currency>) は最初のスロットに位置し、意見 (<Negative>) は 2 番目のスロットに位置します。

カテゴリー規則の作成

規則を作成または編集する場合、規則を条件規則エディターで開く必要があります。コンセプト、タイプ、またはパターンを追加し、またワイルドカードを使用して一致を拡張することができます。認識されたコンセプト、タイプ、パターンを使用すると、すべての関連コンセプトを検出するため、便利です。例えば、コンセプトを使用すると、そのすべての関連キーワード、複数形、および類義語も規則を満たします。同様に、タイプを使用すると、そのすべてのコンセプトも規則にキャプチャされます。

既存の規則を編集するか、カテゴリー名を右クリックして「条件規則の作成」を選択して、条件規則エディターを開くことができます。

コンテキスト・メニューを使用、ドラッグアンドドロップ、または手動でコンセプト、タイプおよびパターンをエディターに入力します。これらのブール型演算子 (&, !(), |) やブラケットを使用して、条件規則式を形成します。一般的なエラーを回避するために、コンセプトを抽出結果ペインまたはデータ・ペインから直接条件規則エディターにドラッグアンドドロップすることを推奨します。エラーを回避するため、規則のシンタックスには十分注意してください。詳しくは、121 ページの『カテゴリー規則シンタックス』のトピックを参照してください。

注: 条件規則がどのようにテキストに合致するかについての例は、「127 ページの『カテゴリー規則の例』」を参照してください。

規則を作成するには

1. データがまだ抽出されていない、または抽出が過去のものである場合は、抽出してください。詳しくは、82 ページの『データの抽出』のトピックを参照してください。

注: いずれのコンセプトも表示されなくなるような方法で抽出をフィルタリングした場合、カテゴリー規則を作成または編集しようとする、エラー・メッセージが表示されます。このエラー・メッセージを回避するには、コンセプトが表示されるように抽出のフィルターを変更してください。

2. カテゴリー・ペインで、規則を追加するカテゴリーを選択します。
3. メニューから、「カテゴリー」>「条件規則の作成」を選択します。ウィンドウにエディターのカテゴリー条件規則エディターのパネルが開きます。
4. 「条件規則名」フィールドに、規則の名前を入力します。名前を入力しない場合、自動的に式を名前として使用します。規則の名前は、後で変更できます。
5. より大きい式のテキスト・フィールドで、次の作業を実行できます。

- テキストをフィールドに直接入力するか、別のパネルからドラッグアンドドロップします。抽出されたコンセプト、タイプ、パターンのみを使用します。例えば、cats という単語を入力しても単数形の cat のみが抽出結果ペインに表示される場合、エディターは cats を認識できなくなります。単数形には自動的に複数形が含まれる場合があり、そうでない場合はワイルドカードを使用することができます。詳しくは、121 ページの『カテゴリ規則シンタックス』のトピックを参照してください。
 - 規則に追加するコンセプト、タイプ、またはパターンを選択してメニューを使用します。
 - ブール型演算子を規則のリンク要素に追加します。ツールバー・ボタンを使用して、「and」のブール型演算子 &、「or」のブール型演算子 |、「not」のブール型演算子 !()、カッコ ()、パターンのブラケット [] を規則に追加します。
6. 「条件規則をテスト」 ボタンをクリックして、規則が適格であることを確認します。詳しくは、121 ページの『カテゴリ規則シンタックス』のトピックを参照してください。見つかったドキュメントまたはレコードの数は、テキスト「テスト結果」の隣のカッコの中に表示されます。このテキストの右側に、認識された規則の要素またはエラー・メッセージが表示されます。タイプ、パターン、またはコンセプトの隣のグラフィックに赤い疑問符が表示されている場合、要素が既知の抽出に一致しないことを示します。一致しない場合、規則によってレコードは検出されません。
 7. 規則の一部をテストするには、該当する部分を選択して、「選択部分をテスト」 をクリックします。
 8. 問題が見つかった場合は、必要な変更を行い、規則を再度テストします。
 9. 終了したら、「保存して閉じる」をクリックして、規則をもう一度保存し、エディターを閉じます。新しい規則名がカテゴリに表示されます。

規則の編集および削除

規則を作成および保存した後、その規則をいつでも編集することができます。詳しくは、121 ページの『カテゴリ規則シンタックス』のトピックを参照してください。

規則が必要ない場合は、削除することができます。

規則を編集するには

1. 「カテゴリ定義」ダイアログ・ボックスの「記述子」テーブルで、規則を選択します。
2. メニューから、「カテゴリ」>「条件規則の編集」を選択するか、規則名をダブルクリックします。エディターが開き、選択された規則が表示されます。
3. 既存の結果およびツールバー・ボタンを使用して、変更を行います。
4. 規則を再テストして、期待される結果が返されることを確認します。
5. 「保存して閉じる」をクリックして、規則をもう一度保存し、エディターを閉じます。

規則を削除するには

1. 「カテゴリ定義」ダイアログ・ボックスの「記述子」テーブルで、規則を選択します。
2. メニューから「編集」>「削除」を選択します。規則がカテゴリから削除されます。

定義済みのインポートおよびエクスポート

独自のカテゴリを Microsoft Excel (*.xls, *.xlsx) ファイルに保存している場合、それらを IBM SPSS Modeler Text Analytics にインポートできます。

また、使用中のインタラクティブ・ワークベンチ・セッション内のカテゴリを Microsoft Excel (*.xls, *.xlsx) ファイルにエクスポートすることもできます。カテゴリをエクスポートすると、記述子やスコアなど、いくつかの追加情報を含めたり除外したりできます。詳しくは、135 ページの『カテゴリのエクスポート』のトピックを参照してください。

定義済みカテゴリにコードがない場合、または新しいコードが必要な場合、メニューから「カテゴリ」>「カテゴリを管理」>「コードを自動生成」を選択して、カテゴリ・ペインでカテゴリ・セットの新しいコードのセットを自動的に生成できます。これにより、既存のコードがすべて削除され、自動的に番号が指定し直されます。

定義済みカテゴリのインポート

定義済みカテゴリを IBM SPSS Modeler Text Analytics にインポートできます。インポートする前に、定義済みカテゴリ・ファイルが Microsoft Excel (*.xls, *.xlsx) ファイルであり、サポート可能な形式のいずれかの構造であることを確認してください。形式を自動的に検知するよう選択することもできます。次の形式がサポートされています。

- フラット リスト形式: 詳しくは、132 ページの『フラット・リスト形式』のトピックを参照してください。
- コンパクト形式: 詳しくは、133 ページの『コンパクト形式』のトピックを参照してください。
- インデント化形式: 詳しくは、134 ページの『インデント形式』のトピックを参照してください。

定義済みカテゴリをインポートするには

1. インタラクティブ・ワークベンチ・メニューから、カテゴリ > カテゴリを管理 > 定義済みカテゴリのインポート を選択します。定義済みカテゴリのインポート・ウィザードが表示されます。
2. 「参照」ドロップダウン・リストから、ファイルを投入するドライブとフォルダーを選択します。
3. リストからファイルを選択します。「ファイル名」テキスト・ボックスにファイルの名前が表示されます。
4. リストから、定義済みカテゴリを含むワークシートを選択します。ワークシート名が「ワークシート」フィールドに表示されます。
5. データ形式の選択を開始するには、「次へ」をクリックします。
6. ファイルの形式を選択するか、自動的に形式を検知しようとするオプションを選択します。自動検知は、最も一般的な形式に最も適しています。
 - フラット リスト形式: 詳しくは、132 ページの『フラット・リスト形式』のトピックを参照してください。
 - コンパクト形式: 詳しくは、133 ページの『コンパクト形式』のトピックを参照してください。
 - インデント化形式: 詳しくは、134 ページの『インデント形式』のトピックを参照してください。
7. その他のインポート・オプションを定義するには、「次へ」をクリックします。形式の自動検知を選択した場合、最後の手順に進みます。
8. このワークシートの 1 行または複数行が列見出しまたはその他の外部情報である場合、「インポート開始行」オプションでインポート開始場所である行番号を選択します。例えば、カテゴリ名が行 7 で始まる場合、ファイルを正しくインポートするためには、このオプションで行番号 7 を入力する必要があります。
9. ファイルにカテゴリ・コードが含まれている場合、オプション「カテゴリ コードを含む」を選択します。これにより、ウィザードがデータを正しく認識します。

10. カラーコード化されたセルおよび凡例を確認し、データが正しく特定されるようにします。ファイルで検出されたエラーは赤で表示され、形式プレビュー・テーブルの下に表示されます。不正な形式が選択された場合、戻って別の形式を選択します。ファイルを修正する必要がある場合、変更を行い、ファイルをもう一度選択してウィザードを再起動します。ウィザードを終了する前にすべてのエラーを修正してください。
11. インポートされる一連のカテゴリおよびサブカテゴリを確認し、これらのカテゴリの記述子の作成方法を定義するには、「次へ」をクリックします。
12. テーブルでインポートされる一連のカテゴリを確認します。記述子として表示されるはずのキーワードが表示されない場合、インポート時に認識されなかったことが考えられます。正しく接頭辞が使用され、正しいセルに表示されていることを確認してください。
13. セッションの既存のカテゴリの処理方法を選択します。
 - すべての既存のカテゴリを置き換え: 既存のカテゴリすべてを削除し、新しくインポートされたカテゴリが代わりに使用されます。
 - 既存のカテゴリに追加 カテゴリをインポートし、既存カテゴリと共通カテゴリを結合します。既存のカテゴリに追加する場合、重複カテゴリの処理方法を決定する必要があります。オプション「結合」を選択すると、カテゴリ名を共有する場合、インポートされたカテゴリは既存のカテゴリと結合されます。オプション「インポートから除外」は、同じ名前が存在する場合、カテゴリのインポートを禁止します。
14. 「キーワードを記述子としてインポート」は、関連するカテゴリの記述子としてデータで特定されるキーワードをインポートするオプションです。
15. 「記述子を派生させてカテゴリを拡張」は、カテゴリ、またはサブカテゴリの名前を示す単語、および注釈を構成する単語空記述子を生成するオプションです。単語が抽出結果に合致する場合、記述子としてカテゴリに追加されます。このオプションを選択すると、カテゴリ名または注釈が長く記述的である場合に、最良の結果を作成します。迅速にカテゴリの記述子を生成し、またカテゴリはこれらの記述子を含むレコードをキャプチャーすることができます。
 - 「生成元」フィールドを使用して、記述子が派生するテキスト、名前またはカテゴリおよびサブカテゴリ、注釈内の単語から選択できます。
 - 「形式を指定」フィールドを使用して、これらの記述子をコンセプトまたは TLA パターンの形式で作成することを選択できます。TLA 抽出が行われない場合、ウィザードの「パターン」オプションが無効となります。
16. 定義済みカテゴリが「カテゴリ」パネルにインポートするには、「終了」をクリックします。

フラット・リスト形式

フラット・リスト形式では、階層のない、上位レベルのカテゴリのみがあります。つまり、サブカテゴリやサブネットはありません。カテゴリ名は 1 つの列に表示されます。

この形式のファイルには、次の情報が含まれます。

- オプションのコード列には、各カテゴリを一意に特定する数値が入力されます。データ・ファイルにコードを含むよう指定（「コンテンツ設定」で「カテゴリ コードを含む」オプションを選択）した場合、カテゴリ名の左隣のセルに各カテゴリの一意のコードを含む列がなければなりません。データにコードが含まれず、後でコードをいくつか作成したい場合、後でいつでもコードを生成できます（「カテゴリ」>「カテゴリを管理」>「コードを自動生成」）。
- 必須の「カテゴリ名」列には、カテゴリのすべての名前が入力されています。この列は、この形式を使用してインポートする場合に必要です。
- カテゴリ名のすぐ右にあるオプションの「注釈」セル。この注釈は、カテゴリ/サブカテゴリを説明するテキストで構成されています。

- オプションの「キーワード」は、カテゴリーの記述子としてインポートできます。認識できるようにするために、これらのキーワードは関連するカテゴリー/サブカテゴリー名のすぐ下にあるセルになければならず、「_firearms, weapons / guns」のようにキーワードの前にアンダースコア (_) を追加する必要があります。キーワード・セルには、各カテゴリーの説明に使用する 1 つまたは複数の単語を入力できます。これらの単語は、ウィザードの最後の手順の指定内容に応じて、記述子としてインポートしたり、無視されます。後で、記述子はテキストから抽出された結果と比較されます。合致が見つかった場合、レコードまたはドキュメントがこの記述子を含むカテゴリーにスコアリングされます。

表 25. コード、キーワード、および注釈を含むフラット・リスト形式

列 A	列 B	列 C
カテゴリー・コード (オプション)	カテゴリー名	注釈
	_記述子/キーワード・リスト(オプション)	

コンパクト形式

コンパクト形式は、階層カテゴリーで使用される点を除いて、フラット・リスト形式と同じ構造です。そのため、各カテゴリーおよびサブカテゴリーの階層レベルを定義するには、コード レベル列が必要です。

この形式のファイルには、次の情報が含まれます。

- 必須の「コード レベル」列には、その行の後続の情報の階層の位置を示す番号が入力されます。例えば、値 1、2、3 が指定され、カテゴリーおよびサブカテゴリーの両方がある場合、1 はカテゴリー、2 はサブカテゴリー、3 はサブ-サブカテゴリーを示します。カテゴリーおよびサブカテゴリーのみがある場合、1 はカテゴリーを、2 はサブカテゴリーを示します。カテゴリーの深さの限り続きます。
- オプションのコード列には、各カテゴリーを一意に特定する値が入力されます。データ・ファイルにコードを含むよう指定（「コンテンツ設定」で「カテゴリー コードを含む」オプションを選択）した場合、カテゴリー名の左隣のセルに各カテゴリーの一意のコードを含む列がなければなりません。データにコードが含まれず、後でコードをいくつか作成したい場合、後でいつでもコードを生成できます（「カテゴリー」>「カテゴリーを管理」>「コードを自動生成」）。
- 必須の「カテゴリー名」列には、カテゴリーおよびサブカテゴリーのすべての名前が入力されています。この列は、この形式を使用してインポートする場合に必要です。
- カテゴリー名のすぐ右にあるオプションの「注釈」セル。この注釈は、カテゴリー/サブカテゴリーを説明するテキストで構成されています。
- オプションの「キーワード」は、カテゴリーの記述子としてインポートできます。認識できるようにするために、これらのキーワードは関連するカテゴリー/サブカテゴリー名のすぐ下にあるセルになければならず、「_firearms, weapons / guns」のようにキーワードの前にアンダースコア (_) を追加する必要があります。キーワード・セルには、各カテゴリーの説明に使用する 1 つまたは複数の単語を入力できます。これらの単語は、ウィザードの最後の手順の指定内容に応じて、記述子としてインポートしたり、無視されます。後で、記述子はテキストから抽出された結果と比較されます。合致が見つかった場合、レコードまたはドキュメントがこの記述子を含むカテゴリーにスコアリングされます。

表 26. コードを含むコンパクト形式の例

列 A	列 B	列 C
階層コード レベル	カテゴリー・コード (オプション)	カテゴリー名
階層コード レベル	サブカテゴリー・コード (オプション)	サブカテゴリー名

表 27. コードを含まないコンパクト形式の例

列 A	列 B
階層コード レベル	カテゴリー名
階層コード レベル	サブカテゴリー名

インデント形式

インデント・ファイル形式の場合、コンテンツには階層があります。つまり、カテゴリーと、1 レベルまたは複数レベルのサブカテゴリーがあります。さらに、構造がこの階層を示すよう、インデント化されています。ファイルの各行にはカテゴリーまたはサブカテゴリーが含まれますが、サブカテゴリーはカテゴリーからインデント化され、サブ-サブカテゴリーはサブカテゴリーからインデント化されています。この構造を Microsoft Excel で手動で作成したり、別の製品からエクスポートし Microsoft Excel にインポートした構造を使用することもできます。

- 最上位レベルのカテゴリー・コードおよびカテゴリー名は、それぞれ列 A および列 B に表示されます。またはコードがない場合、カテゴリー名が列 A に表示されます。
- サブカテゴリー・コードおよびサブカテゴリー名は、それぞれ列 B および列 C に表示されます。またはコードがない場合、サブカテゴリー名が列 B に表示されます。サブカテゴリーはカテゴリーのメンバーです。上位レベルのカテゴリーがない場合、サブカテゴリーはありません。

表 28. コードを含むインデント構造

列 A	列 B	列 C	列 D
カテゴリー・コード (オプション)	カテゴリー名		
	サブカテゴリー・コード (オプション)	サブカテゴリー名	
		サブ-サブカテゴリー・コード (オプション)	サブ-サブカテゴリー名

表 29. コードのないインデント構造

列 A	列 B	列 C
カテゴリー名		
	サブカテゴリー名	
		サブ-サブカテゴリー名

この形式のファイルには、次の情報が含まれます。

- オプションのコードは、各カテゴリーまたはサブカテゴリーを一意に特定する数値でなければなりません。データ・ファイルにコードを含むよう指定 (「コンテンツ設定」で「カテゴリー コードを含む」オプションを選択) した場合、カテゴリー名の左隣のセルに各カテゴリーまたはサブカテゴリーの一意のコードがなければなりません。データにコードが含まれず、後でコードをいくつか作成したい場合、後でいつでもコードを生成できます (「カテゴリー」>「カテゴリーを管理」>「コードを自動生成」)。
- 各カテゴリーおよびサブカテゴリーの必須の「名前」。サブカテゴリーは、カテゴリーから、各行の右側に 1 セルインデントされていなければなりません。
- カテゴリー名のすぐ右にあるオプションの「注釈」セル。この注釈は、カテゴリー/サブカテゴリーを説明するテキストで構成されています。

- オプションの「キーワード」は、カテゴリの記述子としてインポートできます。認識できるようにするために、これらのキーワードは関連するカテゴリ/サブカテゴリ名のすぐ下にあるセルになければならず、「_firearms, weapons / guns」のようにキーワードの前にアンダースコア () を追加する必要があります。キーワード・セルには、各カテゴリの説明に使用する 1 つまたは複数の単語を入力できます。これらの単語は、ウィザードの最後の手順の指定内容に応じて、記述子としてインポートしたり、無視されます。後で、記述子はテキストから抽出された結果と比較されます。合致が見つかった場合、レコードまたはドキュメントがこの記述子を含むカテゴリにスコアリングされます。

重要: 1 つのレベルでコードを使用する場合、各カテゴリおよびサブカテゴリにコードを含む必要があります。そうでない場合、インポート・プロセスが失敗します。

カテゴリのエクスポート

また、使用中のインタラクティブ・ワークベンチ・セッション内のカテゴリを Microsoft Excel (*.xls, *.xlsx) ファイル形式にエクスポートすることもできます。エクスポートされるデータは、大部分がカテゴリ・プロパティのカテゴリ・ペインの現在のコンテンツのデータです。そのため、ドキュメントスコア値もエクスポートする場合は、もう一度スコアリングを行うことをお勧めします。

表 30. カテゴリ・エクスポート・オプション

常にエクスポート...	オプションでエクスポート...
<ul style="list-style-type: none"> • ある場合はカテゴリ・コード • カテゴリ (およびサブカテゴリ) 名 • ある場合はコード レベル (フラット/コンパクト形式) • 列見出し (フラット/コンパクト形式) 	<ul style="list-style-type: none"> • ドキュメント スコア • カテゴリの注釈 • 記述子名 • 記述子数

重要: 記述子をエクスポートする場合、それらはテキスト文字列とアンダースコアの接頭辞に変換されます。この製品に再度インポートする場合、パターン、条件規則、単純なコンセプトである記述子と区別する機能が失われます。これらのカテゴリの本製品で再利用する場合、テキスト分析パッケージ (TAP) を使用して、現在定義されているすべての記述子 (使用されているすべてのカテゴリ、コード、言語ソース) を保持することをお勧めします。TAP ファイルは IBM SPSS Modeler Text Analytics と IBM SPSS Text Analytics for Surveys の両方で使用できます。詳しくは、136 ページの『テキスト分析パッケージの使用』のトピックを参照してください。

定義済みカテゴリをエクスポートするには

1. インタラクティブ・ワークベンチ・メニューから、カテゴリ > カテゴリを管理 > カテゴリのエクスポート を選択します。カテゴリのエクスポート・ウィザードが表示されます。
2. 場所を選択して、エクスポートするファイルの名前を入力します。
3. 「ファイル名」テキスト・ボックスに出力ファイルの名前を入力します。
4. カテゴリ・データをエクスポートする形式を選択するには、「次へ」をクリックします。
5. 次のいずれかの形式を選択します。
 - フラット・リスト形式およびコンパクト・リスト形式: 詳しくは、132 ページの『フラット・リスト形式』のトピックを参照してください。フラット・リスト形式には、サブカテゴリはありません。詳しくは、133 ページの『コンパクト形式』のトピックを参照してください。コンパクト・リスト形式には階層カテゴリが含まれています。
 - インデント化形式: 詳しくは、134 ページの『インデント形式』のトピックを参照してください。
6. エクスポートするコンテンツを選択し、提案されたデータを確認するには、「次へ」をクリックします。

7. エクスポート・ファイルの内容を確認します。
8. 注釈または記述子名など、エクスポートする追加内容を選択または選択解除します。
9. カテゴリーをエクスポートするには、「終了」をクリックします。

テキスト分析パッケージの使用

TAP と呼ばれるテキスト分析パッケージは、テキスト回答のカテゴリー化を行うためのテンプレートとして機能します。TAP には多くのレコードを迅速かつ自動的にコード化するために必要な事前に作成されたカテゴリー・セットおよび言語リソースが含まれているため、TAP を使用すると、最小限の介入でテキスト・データをカテゴリー化できます。言語リソースを使用して、テキスト・データを分析およびマイニングし、主要キーワードを抽出します。テキストの主要コンセプトとパターンに基づき、レコードを TAP で選択したカテゴリー セットにカテゴリー化できます。独自の TAP を作成または TAP を更新できます。

TAP は、次の要素で構成されています。

- **カテゴリー セット:** カテゴリー・セットは、定義済みカテゴリー、カテゴリー・コード、各カテゴリーの記述子、カテゴリー・セット全体の名前で構成されています。記述子とは、キーワード安いやパターン高価などの言語学的要素です。記述子を使用してカテゴリーを定義し、テキストがカテゴリー記述子に一致すると、ドキュメントまたはレコードがカテゴリーに投入されます。
- **言語リソース:** 言語リソースは、一連のライブラリー、および主要キーワードおよびパターンを抽出するために調整された高度なリソースです。これらの抽出コンセプトおよびパターンは、レコードをカテゴリー・セットのカテゴリーに投入できる記述子として使用されます。

独自の TAP を作成、TAP を更新、または TAP を読み込むことができます。

TAP を選択してカテゴリー・セットを選択した後、SPSS Modeler Text Analytics でレコードを抽出およびカテゴリー化できます。

注: TAP を作成し、IBM SPSS Text Analytics for Surveys および SPSS Modeler Text Analytics の製品間で交互に使用できます。ただし、SPSS Modeler Text Analytics からテキスト分析パッケージ (TAP) を直接ロードするか IBM SPSS Text Analytics for Surveys から TAP をロードするかに応じて、SPSS Modeler Text Analytics でのルールのスコアリングが異なる場合があります。SPSS Modeler Text Analytics 内で作成された TAP を使用することをお勧めします。IBM SPSS Text Analytics for Surveys で作成された TAP は、異なるバージョンの言語リソースを使用して作成されている可能性があるからです。

テキスト分析パッケージの作成

少なくとも 1 つのカテゴリーといくつかのリソースを含むセッションがある場合、オープン インタラクティブ・ワークベンチ・セッションのコンテンツからテキスト分析パッケージ (TAP) を作成できます。カテゴリーおよび記述子 (コンセプト、タイプ、条件規則または TLA パターン出力) のセットをリソース・エディターで開かれたすべての言語リソースと共に使用して TAP を作成できます。

リソースが作成された言語を表示できます。言語は、テンプレート・エディター または リソース・エディター ビューの「高度なリソース」タブで設定します。

テキスト分析パッケージを作成するには

1. メニューで、「ファイル」>「テキスト分析パッケージ」>「パッケージを作成」を選択します。「パッケージを作成」ダイアログが表示されます。

2. TAP を保存するディレクトリーを参照します。デフォルトでは、TAP は製品のインストール・ディレクトリーの ¥TAP サブディレクトリーに保存されます。
3. 「ファイル名」 フィールドに TAP の名前を入力します。
4. 「パッケージ ラベル」 フィールドにラベルを入力します。ファイル名を入力すると、この名前がラベルとして自動的に表示されますが、このラベルは変更できます。
5. TAP からカテゴリ・セットを除外するには、「追加」 チェックボックスをオフにします。カテゴリ・セットを除外すると、カテゴリ・セットがパッケージに追加されなくなります。デフォルトでは、質問ごとに 1 つのカテゴリ・セットが TAP に追加されます。TAP には、1 つ以上のカテゴリ・セットが必要です。
6. カテゴリ・セットの名前を変更します。「新しいカテゴリ・セット」 列にはデフォルトで一般名が入力されています。一般名はテキスト変数名に Cat_ 接頭辞を追加して生成されます。セルを 1 回クリックすると、名前を編集できます。入力して他の場所をクリックすると、名前の変更が適用されます。カテゴリ・セットの名前を変更すると、TAP のみの名前が変更され、オープン セッションの編集名は変わりません。
7. 必要に応じて、カテゴリ・セット テーブルの右側にある矢印キーを使用して、カテゴリ・セットを並べ替えます。
8. 「保存」 をクリックして、テキスト分析パッケージを作成します。ダイアログ・ボックスがクローズします。

テキスト分析パッケージの読み込み

テキスト マイニング・モデル作成ノードを設定している場合、抽出時に使用するリソースを指定する必要があります。リソース・テンプレートを選択する代わりに、テキスト分析パッケージ (TAP) を選択して、そのリソースだけでなく、カテゴリ・セットをノードにコピーすることができます。

カテゴリ・セットをカテゴリ化の開始ポイントとして使用できるため、カテゴリ・モデルをインタラクティブに作成する場合、TAP が最も重要となります。ストリームを実行すると、インタラクティブ・ワークベンチ・セッションが起動し、このセットのカテゴリがカテゴリ・ペインに表示されます。このように、これらのカテゴリを使用してすぐにドキュメントおよびレコードをスコアリングを行い、これらのカテゴリが適切なものとなるまで調整、作成、拡張を続行します。詳しくは、98 ページの『カテゴリ作成の方法と戦略』のトピックを参照してください。

バージョン 14 以降、「読み込み」 をクリックして TAP を選択すると、この TAP のリソースが定義された言語を表示することもできます。

テキスト分析パッケージを読み込むには

1. テキスト マイニング・モデル作成ノードを編集します。
2. 「モデル」 タブの「リソースのコピー元」で、「テキスト分析パッケージ」を選択します。
3. 「読み込み」 をクリックします。「テキスト分析パッケージの読み込み」ダイアログが表示されます。
4. ノードにコピーしたいリソースおよびカテゴリ・セットを含む TAP の場所を参照します。デフォルトでは、TAP は製品のインストール・ディレクトリーの ¥TAP サブディレクトリーに保存されます。
5. 「ファイル名」 フィールドに TAP の名前を入力します。ラベルが自動的に表示されます。
6. 使用したいカテゴリ・セットを選択します。インタラクティブ・ワークベンチ・セッションで表示されるカテゴリのセットです。手動で、またはカテゴリの作成または展開のオプションを使用して、これらのカテゴリを調整および改善できます。

7. 「読み込み」をクリックして、テキスト分析パッケージの内容をノードにコピーします。ダイアログ・ボックスがクローズします。TAP が読み込まれると、TAP のコピーがノードにコピーされます。そのため、リソースおよびカテゴリーに行った変更は、明示的に更新および再読み込みしないかぎり TAP に反映されます。

テキスト分析パッケージの更新

カテゴリー・セット、言語リソースを改善するか、新しいカテゴリー・セットを作成する場合、テキスト分析パッケージ (TAP) を更新して、これらの改善点を後で再利用しやすくすることができます。これには、TAP に置きたい情報を含むオープンな セッション内にいる必要があります。更新する場合、カテゴリー・セットの追加、リソースの置き換え、パッケージ ラベルの変更、またはカテゴリー・セットの名前変更/並べ替えを選択できます。

テキスト分析パッケージを更新するには

1. メニューで、「ファイル」>「テキスト分析パッケージ」>「パッケージを更新」を選択します。「パッケージの更新」ダイアログが表示されます。
2. 更新したいテキスト分析パッケージを含むディレクトリーを参照します。
3. 「ファイル名」フィールドに TAP の名前を入力します。
4. TAP 内の言語リソースと現在のセッションの言語リソースと置き換えるには、「このパッケージのリソースを使用中のセッションのリソースで置き換える」オプションを選択します。通常、言語リソースはカテゴリー定義の作成に使用される主要キーワードおよびパターンを抽出するために使用されるため、言語リソースの更新は重要です。最新の言語リソースがあれば、レコードのカテゴリー化において最善の結果を得ることができます。このオプションを選択しない場合、パッケージ内にすでにある言語リソースは変更されません。
5. 言語リソースのみを更新するには、かならず「このパッケージのリソースを使用中のセッションのリソースで置き換える」オプションを選択し、TAP 内にすでにある現在のカテゴリー・セットのみを選択してください。
6. オープンなセッションの新しいカテゴリー・セットを TAP に追加するには、追加する各カテゴリーのチェックボックスをチェックしてください。1 つ、複数のカテゴリー・セットを追加できますし、追加しなくてもかまいません。
7. TAP からカテゴリー・セットを削除するには、対応する「追加」チェックボックスをオフにします。改善されたカテゴリー・セットを追加するため、すでに TAP にあるカテゴリー・セットを削除する場合があります。カテゴリー・セットを削除するには、「現在のカテゴリー・セット」列の該当するカテゴリー・セットの「追加」チェックボックスをオフにします。TAP には、1 つ以上のカテゴリー・セットが必要です。
8. 必要に応じて、カテゴリー・セットの名前を変更します。セルを 1 回クリックすると、名前を編集できます。入力して他の場所をクリックすると、名前の変更が適用されます。カテゴリー・セットの名前を変更すると、TAP のみの名前が変更され、オープン セッションの編集名は変わりません。2 つのカテゴリー・セットに同じ名前が付いている場合、重複を修正するまで名前が赤で表示されます。
9. 選択した TAP の内容と結合されたセッションの内容で新しいパッケージを作成するには、「名前を付けて保存」をクリックします。「テキスト分析パッケージの名前を付けて保存」ダイアログが表示されます。次の説明を参照してください。
10. 「更新」をクリックすると、選択した TAP に行われた変更が保存されます。

テキスト分析パッケージを保存するには

1. TAP ファイルを保存するディレクトリーを参照します。デフォルトでは、TAP ファイルはインストール・ディレクトリーの TAP サブディレクトリーに保存されます。

2. 「ファイル名」フィールドに TAP ファイルの名前を入力します。
3. 「パッケージ ラベル」フィールドにラベルを入力します。ファイル名を入力すると、この名前がラベルとして自動的に表示されます。ただし、このラベルの名前は変更できます。ラベルはかならず指定する必要があります。
4. 「保存」 をクリックして新しいパッケージを作成します。

カテゴリーの編集および調整

カテゴリーを作成すると、それらを検証して、何らかの調整を行う必要が常にあります。言語リソースを調整するほか、定義を結合またはクリーンアップしたり、カテゴリー化されたドキュメントまたはレコードの一部をチェックするための方法を探すことによってカテゴリーを確認する必要があります。また、ニュアンスや特徴が分かるようにカテゴリーが定義されるよう、カテゴリーのドキュメントまたはレコードの確認を行うこともできます。

ビルトインで、自動のカテゴリー作成手法を使用して、カテゴリーを作成できますが、これらのカテゴリーに何らかの調整が必要な場合がよくあります。1 つまたは複数の手法を使用した後、多くの新規カテゴリーがウィンドウに表示されます。カテゴリー定義が適切なものとなるまで、カテゴリーのデータを確認して調整を行うことができます。詳しくは、103 ページの『カテゴリーとは』のトピックを参照してください。

カテゴリー調整のオプションがいくつかあります。その多くについては次のページで説明します。

記述子のカテゴリーへの追加

自動的手法を使用した後に、カテゴリー定義のいずれにも使用されなかった抽出結果が作成される場合があります。結果のリストを「抽出結果」ウィンドウで確認する必要があります。カテゴリーに含ませたい要素があった場合、これを既存のあるいは新規のカテゴリーに追加します。

コンセプトまたはタイプをカテゴリーに追加するには

1. 抽出結果ペインおよびデータ・ペインから、新規または既存のカテゴリーに追加する要素を選択します。
2. メニューから、「カテゴリー」>「カテゴリーに追加」を選択します。「すべてのカテゴリー」ダイアログボックスにカテゴリーのセットが表示されます。選択した要素を追加したいカテゴリーを選択します。要素を新規カテゴリーに追加する場合、「新規カテゴリー」を選択します。最初に選択した要素を使用した新しいカテゴリーが、カテゴリー・ペインに表示されます。

カテゴリー記述子の編集






いくつかのカテゴリーを作成すると、各カテゴリーを開いてその定義を構成するすべての記述子を表示できます。「カテゴリー定義」ダイアログ・ボックスで、カテゴリー記述子にあらゆる編集を行うことができます。また、カテゴリーがカテゴリー・ツリーに表示されている場合、ここで処理することもできます。

カテゴリーを編集するには

1. カテゴリー・ペインで編集したいカテゴリーを選択します。
2. メニューから、「表示」>「カテゴリー定義」を選択します。「カテゴリー定義」ダイアログ・ボックスが開きます。
3. 編集したい記述子を選択し、該当するツールバー・ボタンをクリックします。

次の表で、カテゴリー定義を編集できるツールバーボタンについて説明します。

表 31. ツールバー・ボタンおよび説明：

アイコン	説明
	選択した記述子をカテゴリから削除します。
	選択した記述子を新規または既存のカテゴリに移動します。
	選択した記述子を & カテゴリ規則の形式でカテゴリに移動します。詳しくは、121 ページの『カテゴリ規則の使用』のトピックを参照してください。
	選択した各記述子を、独自の新規カテゴリとして移動します。
 「表示」	選択した記述子に従って、データ・ペインおよび視覚化ペインの表示内容を更新します。

カテゴリの移動

カテゴリを別の既存カテゴリに投入または記述子を別のカテゴリに移動したい場合、カテゴリを移動できます。

カテゴリを移動するには

1. カテゴリ・ペインで、別のカテゴリに移動したいカテゴリを選択します。
 2. メニューから、「カテゴリ」>「カテゴリに移動」を選択します。メニューに一連のカテゴリが表示され、リストの上部には最近作成されたカテゴリが表示されます。選択したコンセプトを移動したいカテゴリ名を選択します。
- 探している名前が表示されたら、その名前を選択します。選択した要素がそのカテゴリに追加されます。
 - 名前が表示されない場合、「もっと表示」を選択すると「すべてのカテゴリ」ダイアログ・ボックスが表示され、リストからカテゴリを選択します。

カテゴリのフラット化

カテゴリおよびサブカテゴリを持つ階層カテゴリ構造がある場合、構造をフラットにすることができます。カテゴリをフラット化する場合、そのカテゴリのサブカテゴリのすべての記述子が選択されたカテゴリに移動し、空のサブカテゴリが削除されます。このように、サブカテゴリへの合致に使用されるすべてのドキュメントは、選択したカテゴリにカテゴリ化されます。

カテゴリをフラット化するには

1. カテゴリ・ペインで、フラットにするカテゴリ（上位レベルまたはサブカテゴリ）を選択します。
2. メニューの「カテゴリ」>「カテゴリをフラット化」を選択します。サブカテゴリが削除され、記述子が選択したカテゴリに結合されます。

カテゴリーの結合・組み合わせ

2 つ以上の既存カテゴリーを 1 つの新規カテゴリーに結合したい場合、それらを結合できます。カテゴリーを結合する場合、一般名の付いた新規カテゴリーが作成されます。カテゴリー記述子に使用されているすべてのコンセプト、タイプ、およびパターンがこの新規カテゴリーに移動します。このカテゴリー名は、カテゴリーのプロパティを編集することで後から変更できます。

カテゴリーまたはカテゴリーの一部を結合するには

1. カテゴリー・ペインで、結合したい要素を選択します。
2. メニューの「カテゴリー」>「カテゴリーの結合」を選択します。カテゴリー・プロパティ・ダイアログボックスが表示されるので、新しく作成したカテゴリーの名前を入力します。選択したカテゴリーがサブカテゴリーとして新しいカテゴリーに結合されます。

カテゴリーの削除

カテゴリーを保持しない場合、削除することができます。

カテゴリーを削除するには

1. カテゴリー・ペインで、削除したいカテゴリーを選択します。
2. メニューから「編集」>「削除」を選択します。

第 10 章 クラスターの分析

クラスター・ビューで、コンセプトのクラスターを作成および検討できます（「表示」 > 「クラスター」）。クラスターは、ドキュメント/レコード・セットでこれらのコンセプトが出現する頻度、および共起とも呼ばれる、同じドキュメントで同時に出現する頻度に基づいてアルゴリズムをクラスターリングすることによって生成される関連コンセプトのグループです。クラスター内の各コンセプトは、クラスター内の 1 つ以上の他のコンセプトと共に出現します。カテゴリーの目的は、含まれるテキストが各カテゴリーの記述子（コンセプト、条件規則、パターン）にどのように合致するかに基づいてドキュメントまたはレコードをグループ化することですが、クラスターの目的は共起するコンセプトをグループ化することです。

良いクラスターとは、リンクが強く頻繁に共起するコンセプトを含み、他のクラスターのコンセプトへのリンクが少ないクラスターです。大きなデータセットを扱う場合、この手法の処理時間が大幅に長くなる場合があります。

クラスターリングは、コンセプトのセットを分析し、ドキュメントで頻繁に共起するコンセプトを探すことから始まります。ドキュメント内で共起する 2 つのコンセプトは、コンセプト・ペアと見なされます。次に、クラスターリング・プロセスで、ペアが同時に出現するドキュメント数を各コンセプトが出現するドキュメント数と比較して、各コンセプト・ペアの類似度値を評価します。詳しくは、145 ページの『類似度リンク値の計算』のトピックを参照してください。

最後に、リンク値と「クラスターの作成」ダイアログ・ボックスで定義された設定を集約および考慮し、類似したコンセプトをグループ化します。集約とは、コンセプトを追加、またはクラスターが飽和するまで小さいクラスターを大きいクラスターに結合することです。コンセプトまたは小さいクラスターのさらなる結合によってクラスターが「クラスターの作成」ダイアログ・ボックスの設定（コンセプト、内部リンク、外部リンクの数）を超えると、クラスターが飽和します。クラスターは、クラスター内の他のコンセプトへのリンク数全体が最も大きいクラスター内のコンセプトの名前を使用します。

別のクラスターにより強いリンクがある場合、そして飽和によってコンセプトが出現するクラスターの結合が行われない場合があるため、同じクラスターのすべてのコンセプト・ペアが同時に出現するとは限りません。このため、内部リンクと外部リンクの両方が存在します。

- 内部リンクは、クラスター内のコンセプト・ペア間のリンクです。すべてのコンセプトがクラスター内のお互いのコンセプトにリンクしているわけではありません。ただし、各コンセプトは、クラスター内の 1 つ以上の他のコンセプトにリンクしています。
- 外部リンクは、別のクラスターのコンセプト・ペア間のリンクです（あるクラスターのコンセプトと別のクラスターのコンセプトとの間）。

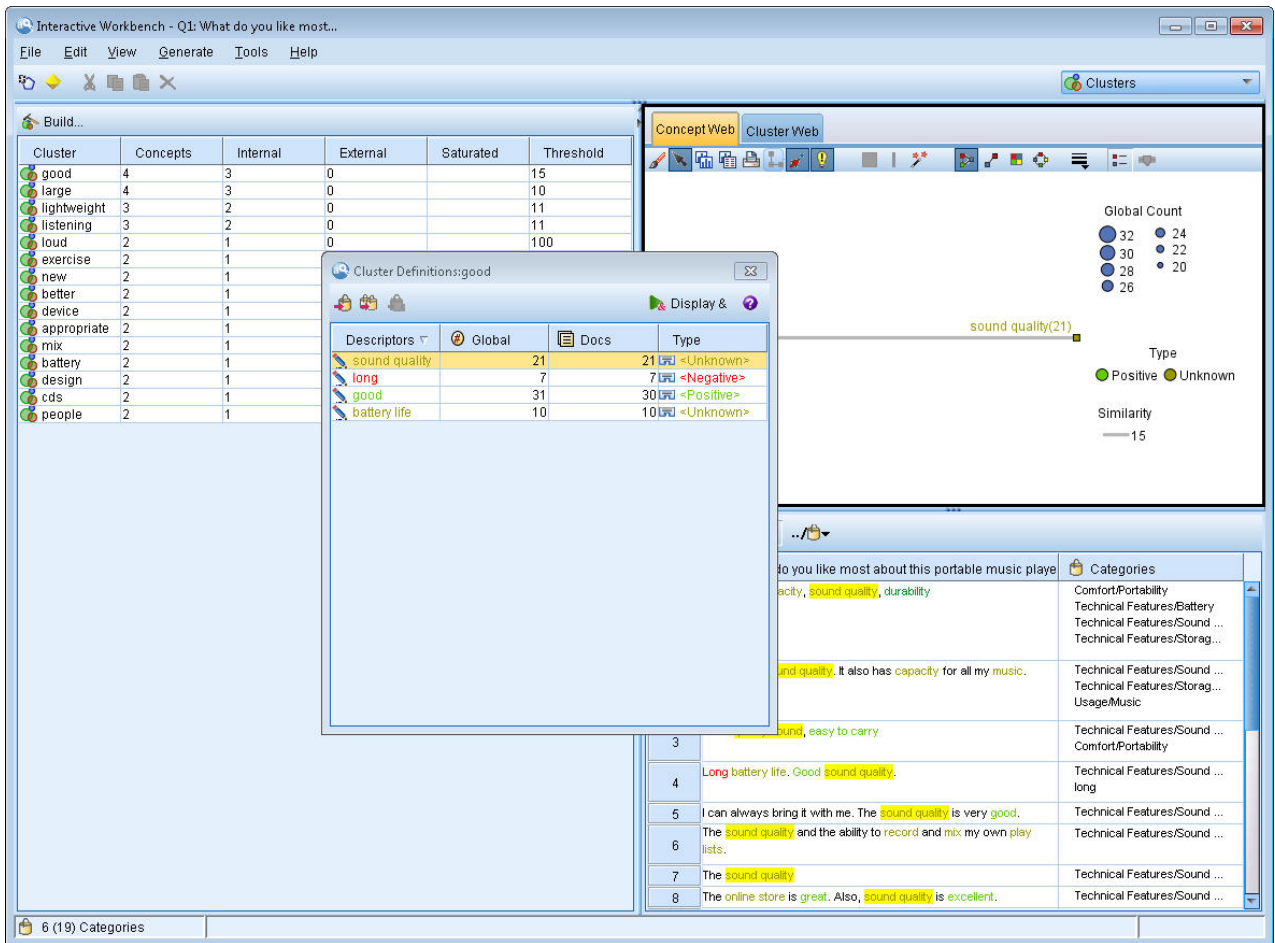


図 30. クラスタ・ビュー

クラスタ・ビューは次の 3 つのパネルで構成され、「表示」メニューから名前を選択して隠したり表示したりできます。

- **クラスタ・ペイン:** このパネルでカテゴリーを作成し、管理できます。詳しくは、146 ページの『クラスタの検証』のトピックを参照してください。
- **視覚化パネル:** このペインで、クラスタについて、またクラスタがどのように相互作用するかを視覚的に探索できます。詳しくは、157 ページの『クラスタ・グラフ』のトピックを参照してください。
- **データ パネル:** このパネルでの選択に対応するドキュメントおよびレコード内に含まれるテキストを検証および確認できます。詳しくは、147 ページの『クラスタ定義』のトピックを参照してください。

クラスタの作成

クラスタ・ビューを初めて開くと、クラスタは表示されません。メニュー (「ツール」>「クラスタの作成」) を使用するか、ツールバーの「作成...」をクリックして、クラスタを作成できます。「クラスタの作成」ダイアログ・ボックスが開き、クラスタ作成の設定および制約を定義できます。

注: 抽出結果がリソースに一致しない場合、このパネルが抽出結果ペインと同じように黄色で表示されます。再抽出を実行して、最新の抽出結果を取得すると、黄色の表示は解除されます。ただし、抽出が実行されるごとにクラスタ・ペインは解除され、クラスタを再作成する必要があります。同様に、あるセッションのクラスタは別のセッションに保存されません。

以下は、「クラスターの作成」ダイアログ・ボックスにある領域とフィールドです。

入力

入力テーブル クラスターは、特定のタイプから派生した記述子から作成されます。テーブル内で、作成プロセスに使用するタイプを選択できます。デフォルトでは、最も多くのレコードまたはドキュメントをキャプチャーするタイプが事前に選択されています。

クラスターのコンセプト: クラスタリングに使用するコンセプトの選択方法を選択します。コンセプトの数を減らすと、クラスタリング・プロセスの速度を向上させることができます。さまざまな上位コンセプト、上位コンセプトの割合を使用して、またはすべてのコンセプトを使用してクラスタリングを行うことができます。

- **数値 (ドキュメント数に基づく):** 「上位のコンセプト」を選択した場合、クラスタリングに考慮するコンセプトの数を入力します。コンセプトは、ドキュメント数の値が最も大きいコンセプトに基づいて選択します。ドキュメント数は、コンセプトが出現するドキュメントまたはレコードの数です。最大値は 150,000 です。
- **パーセンテージ (ドキュメント数に基づく):** 「上位コンセプトの割合」を選択した場合、クラスタリングに考慮するコンセプトの割合を入力します。コンセプトは、ドキュメント数の値が最も大きいコンセプトの割合に基づいて選択します。

出力の制約

作成する最大クラスター数 この値は、生成し、クラスター・ペインに表示される最大クラスター数です。クラスタリング・プロセスで、飽和したクラスターは不飽和クラスターの前に表示されます。つまり、生成される多くのクラスターが飽和します。より多くの不飽和クラスターを表示するために、この設定を飽和クラスターの数より大きい値に変更できます。

クラスター内の最大コンセプト数 この値は、クラスターが含むことができる最大コンセプト数です。

クラスター内の最小コンセプト数 この値は、クラスターを作成するためにリンクする必要のある最小コンセプト数です。

内部リンクの最大数 この値は、クラスターが含むことができる内部リンクの最大数です。内部リンクは、クラスター内のコンセプト・ペア間のリンクです。

外部リンクの最大数 この値は、クラスター外部のコンセプトへのリンクの最大数です。外部リンクは、別のクラスターのコンセプト・ペア間のリンクです。

最小リンク値 この値は、クラスタリングに考慮されるコンセプト・ペアに受け入れられる最小リンク値です。リンク値は、類似度評価式を使用して計算します。詳しくは、『類似度リンク値の計算』のトピックを参照してください。

特定のコンセプトがグループ化されないようにする: 出力の 2 つのコンセプトがグループ化またはペアとにならないように処理を停止します。コンセプト・ペアを作成または管理するには、「ペアを管理」をクリックします。詳しくは、110 ページの『例外ペアのリンクの管理』のトピックを参照してください。

類似度リンク値の計算

コンセプト・ペアが出現するドキュメント数がわかっているだけでは、2 つのコンセプトがどの程度類似しているかはわかりません。この場合、類似度値が役に立ちます。類似度リンク値は、関連性における各コンセプトの共起ドキュメント数を個別ドキュメント数に比較して測定します。類似度を計算する場合、測定の

単位はコンセプトまたはコンセプト数が見つかったドキュメント数です。コンセプトまたはコンセプト・ペアが「少なくとも」1回ドキュメント内に出現した場合、コンセプトまたはコンセプト・ペアがドキュメント内で「見つかった」といえます。コンセプト・グラフのラインの太さを、グラフの類似度リンク値を示すよう選択できます。

アルゴリズムを使用して、最も強いこれらの関連性を明らかにします。つまり、コンセプトがテキスト・データで同時に出現する傾向は、個別に出現する傾向より高くなります。内部的に、アルゴリズムは0から1の類似度係数を生成します。1の値は2つのコンセプトが常に同時に出現し、個別には出現しないことを意味します。類似度係数の結果に100をかけ、最も近い整数に丸められます。類似度係数は、次の図に示された式を使用して計算されます。

$$\text{similarity coefficient} = \frac{(C_{IJ})^2}{(C_I \times C_J)}$$

図 31. 類似度係数の式

この場合、次のようになります。

- C_I は、コンセプト I が出現するドキュメントまたはレコードの数です。
- C_J は、コンセプト J が出現するドキュメントまたはレコードの数です。
- C_{IJ} は、コンセプトのペア I および J が同時に出現するドキュメントまたはレコードの数です。

例えば、5,000 件のドキュメントがあるとします。I および J が抽出されたコンセプト、IJ が I および J のコンセプト・ペアの共起であるとします。次の表では、係数とリンク値を計算する方法を示す2つのシナリオを示しています。

表 32. コンセプトの度数の例

コンセプト/ペア	シナリオ A	シナリオ B
コンセプト:I	20 件のドキュメントに出現	30 件のドキュメントに出現
コンセプト:J	20 件のドキュメントに出現	60 件のドキュメントに出現
コンセプト・ペア:IJ	20 件のドキュメントに共起	20 件のドキュメントに共起
類似度係数	1	0.22222
類似度リンク値	100	22

シナリオ A の場合、コンセプト I と J、および IJ は 20 件のドキュメントに出現します。この場合類似度係数は 1 となり、コンセプトは常に同時に出現します。このペアの類似度リンク値は 100 となります。

シナリオ B の場合、コンセプト I は 30 件のドキュメントに出現し、コンセプト J は 60 件のドキュメント、ペア IJ は 20 件のドキュメントにのみ出現します。その結果、類似度係数は 0.22222 となります。このペアの類似度リンク値は、丸められて 22 となります。

クラスターの検証

クラスターを作成した後、クラスター・ペインで一連の結果を確認できます。各クラスターについて、次の情報がテーブルに表示されます。

- クラスター: クラスターの名前です。クラスターは、内部リンク数が最も多いコンセプトから名前を付けられます。

- **コンセプト:** クラスタ内でのコンセプト数です。詳しくは、『クラスタ定義』のトピックを参照してください。
- **内部:** クラスタ内での内部リンクの数です。内部リンクは、クラスタ内でのコンセプト・ペア間のリンクです。
- **外部:** クラスタ内での外部リンクの数です。外部リンクは、あるクラスタのコンセプトと、別のクラスタのコンセプトとのコンセプト・ペア間のリンクです。
- **飽和:** 記号が表示されている場合、このクラスタが大きく、1 つまたは複数の制約を超えていることを示します。そのため、そのクラスタのクラスタリング・プロセスは終了し、「飽和」していると見なされます。クラスタリング・プロセスの終了時、飽和したクラスタは不飽和クラスタの前に表示されます。つまり、生成される多くのクラスタが飽和します。より多くの不飽和クラスタを表示するためには、「作成する最大クラスタ数」の設定を飽和クラスタの数より大きい値に変更するか、「最小リンク値」の値を小さくできます。詳しくは、144 ページの『クラスタの作成』のトピックを参照してください。
- **閾値:** クラスタ内でのすべての共起コンセプト・ペアについて、クラスタで最も低い類似度リンク値です。詳しくは、145 ページの『類似度リンク値の計算』のトピックを参照してください。閾値の最も大きいクラスタは、そのクラスタのコンセプトの全体の類似度が高く、閾値が小さいクラスタのコンセプトより密接に関連していることを示します。

指定されたクラスタについての詳細を知るには、クラスタを選択すると、右側の視覚化パネルにクラスタを検証するための 2 つのグラフが表示されます。詳しくは、トピック「157 ページの『クラスタ・グラフ』」を参照してください。テーブルの内容を切り取り、別のアプリケーションに貼り付けることができます。

抽出結果がリソースに一致しない場合、このパネルが抽出結果ペインと同じように黄色で表示されます。再抽出を実行して、最新の抽出結果を取得すると、黄色の表示は解除されます。ただし、抽出が実行されるごとにクラスタ・ペインは解除され、クラスタを再作成する必要があります。同様に、あるセッションのクラスタは別のセッションに保存されません。

クラスタ定義

クラスタ・ペインでクラスタを選択し、「クラスタ定義」ダイアログ・ボックスを開くと、クラスタ内のすべてのコンセプトが表示されます（「表示」>「クラスタ定義」）。

選択したクラスタのすべてのコンセプトが「クラスタ定義」ダイアログ・ボックスに表示されます。



「クラスタ定義」ダイアログ・ボックスの 1 つまたは複数のコンセプトを選択し、「表示」をクリックすると、「データ」ウィンドウに「選択したすべてのコンセプトがいっしょに出現する」すべてのレコードまたはドキュメントが表示されます。ただし、クラスタ・ペインでクラスタを選択した場合、データ・ペインにはテキスト・レコードまたはドキュメントは表示されません。「データ」パネルに関する一般情報については、「104 ページの『データ・ペイン』」を参照してください。

このダイアログ・ボックスでコンセプトを選択すると、コンセプト Web グラフも変わります。詳しくは、トピック「157 ページの『クラスタ・グラフ』」を参照してください。同様に、「クラスタ定義」ダイアログ・ボックスで 1 つまたは複数のコンセプトを選択すると、視覚化ペインにこれらのコンセプトの外部リンクおよび内部リンクがすべて表示されます。

列の説明

各記述子を容易に特定できるよう、アイコンが表示されます。





表 33. 列および記述子アイコン

列	説明
記述子	コンセプトの名前
 グローバル	データセット全体にこの記述子が出現する回数を示します。グローバル出現頻度とも呼ばれます。
 ドキュメント	この記述子が出現するドキュメントまたはレコードの数を示します。ドキュメント数とも呼ばれます。
データ型	記述子が属するタイプを示します。記述子がカテゴリ規則である場合、この列にタイプ名は表示されません。

ツールバーの操作

このダイアログ・ボックスから、カテゴリで使用する 1 つまたは複数のコンセプトを選択することもできます。コンセプトを選択する方法はいくつかありますが、クラスター内で共起するコンセプトを選び、カテゴリ規則としてそれらを追加することが最も興味深い方法です。詳しくは、114 ページの『共起規則』のトピックを参照してください。 ツールバー・ボタンを使用して、コンセプトをカテゴリに追加できます。

表 34. コンセプトをカテゴリに追加するツールバー・ボタン

アイコン	説明
	選択したコンセプトを新規または既存のカテゴリに追加します。
	選択したコンセプトを新規または既存のカテゴリに & カテゴリ規則の形式で追加します。詳しくは、121 ページの『カテゴリ規則の使用』のトピックを参照してください。
	選択した各コンセプトを、独自の新規カテゴリとして追加します。
	選択した記述子に従って、データ・ペインおよび視覚化ペインの表示内容を更新します。

注: コンテキスト・メニューを使用して、コンセプトをタイプに、類義語、または不要語項目として追加することもできます。

第 11 章 テキスト リンク分析の検証

テキスト リンク分析 (TLA) ビューでは、テキスト リンク分析パターンを検証できます。テキスト リンク分析はパターンマッチ手法で、パターン規則を定義し、それらをテキスト内の実際の抽出されたコンセプトおよび関連性と比較することができます。

例えば、組織に関するキーワードを抽出しても、重要でない場合があります。TLA を使用して、この組織と他の組織、または組織内の人々間のリンクについて学習することができます。TLA を使用して、製品に関する意見、または遺伝子間の関連性についていくつかの言語で抽出することもできます。

TLA パターン結果を抽出すると、テキスト リンク分析ビューのタイプ・パターン・ペインまたはコンセプト・パターン・ペインでそれらを確認できます。詳しくは、151 ページの『タイプ・パターンおよびコンセプト・パターン』のトピックを参照してください。このビューのデータ・ペインまたは視覚化ペインで TLA パターン結果をさらに検証できます。そしておそらく最も重要なことですが、TLA パターン結果をカテゴリに追加できます。

また TLA パターンの抽出を選択していない場合、「抽出」をクリックして、「抽出設定」ダイアログ・ボックスで「テキスト リンク分析のパターン抽出を有効にする」を選択できます。詳しくは、150 ページの『TLA パターン結果の抽出』のトピックを参照してください。

いくつかの TLA パターン規則が、TLA パターン結果を抽出するために使用するリソース・テンプレートまたはライブラリーで定義されています。IBM SPSS Modeler Text Analytics に付属する特定のリソース・テンプレートで TLA パターンを使用できます。抽出できる関連性およびパターンの種類は、全体的にリソースで定義された TLA 規則によって異なります。独自の TLA 規則を定義できます。パターンは、マクロ、単語リスト、およびブール型質問を形成する単語の空所、または入力テキストと比較される条件規則で構成されています。詳しくは、215 ページの『第 18 章 テキスト リンク規則について』のトピックを参照してください。

TLA パターン規則がテキストに一致する場合、このテキストをパターンとして抽出し、出力データとして再構築できます。そして結果は、テキスト リンク分析ビューのパネルで表示されます。「表示」メニューでパネルの名前を選択して、各パネルを隠したり表示することができます。

- タイプ・パターン・ペインおよびコンセプト・パターン・ペイン: これら 2 つのパネルでパターンを作成し、検証できます。詳しくは、151 ページの『タイプ・パターンおよびコンセプト・パターン』のトピックを参照してください。
- 視覚化ペイン。このパネルで、パターンのコンセプトおよびタイプがどのように相互作用するかを視覚的に検証できます。詳しくは、158 ページの『テキスト リンク分析のグラフ』のトピックを参照してください。
- データ・ペイン。別のパネルでの選択に対応するドキュメントおよびレコード内に含まれるテキストを検証および確認できます。詳しくは、153 ページの『データ パネル』のトピックを参照してください。

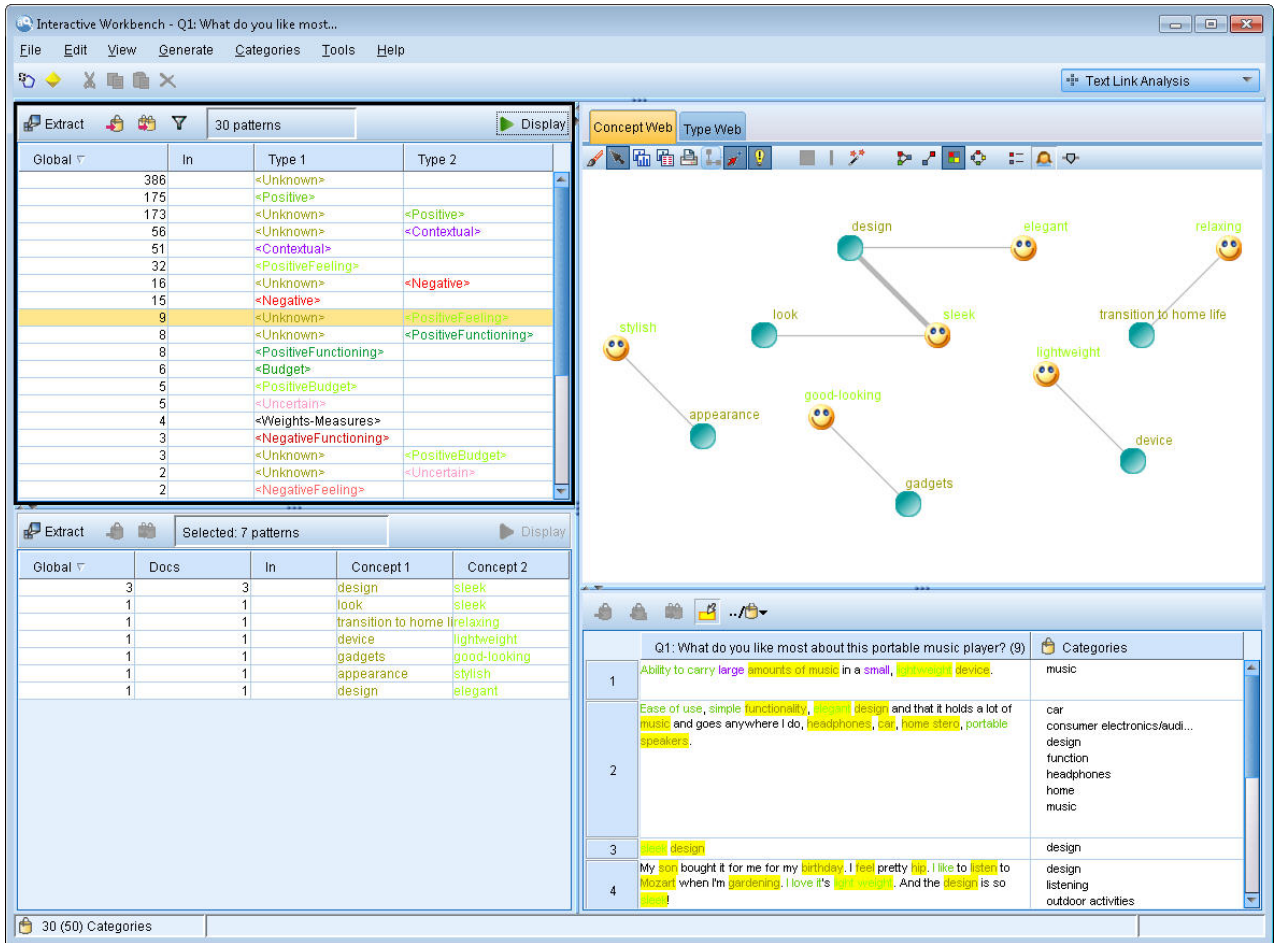


図 32. テキスト リンク分析ビュー

TLA パターン結果の抽出

抽出プロセスにより、一連のコンセプト、タイプ、そして有効な場合はテキスト リンク分析 (TLA) パターンが作成されます。TLA パターンを抽出した場合、これらはテキスト リンク分析ビューにされます。抽出結果がリソースと同期していない場合、パターン・ペインが黄色で表示され、再抽出をすると異なる結果が生成されることを示します。

「テキスト リンク分析のパターン抽出を有効にする」 オプションを使用して、ノードの設定または「抽出」ダイアログ・ボックスでこれらのパターンの抽出を選択する必要があります。詳しくは、82 ページの『データの抽出』のトピックを参照してください。

注:データセットのサイズと、抽出プロセスを完了するためにかかる時間の間には、関連性があります。パフォーマンス統計および推奨事項については、インストール手順を参照してください。上流にサンプル・ノードを追加、またはコンピューターの構成を最適化することをいつでも検討することができます。

データを抽出するには

1. メニューの「ツール」>「抽出」を選択します。または、「抽出」 ツールバー・ボタンをクリックします。

2. 使用するオプションを変更します。このタブで TLA パターン結果を抽出するには、オプション「テキスト リンク分析のパターン抽出を有効にする」を選択し、テンプレートに TLA 規則を使用する必要があります。詳しくは、82 ページの『データの抽出』のトピックを参照してください。
3. 「抽出」をクリックして、抽出プロセスを開始します。

抽出が始まると、進捗状況のダイアログ・ボックスが表示されます。抽出を中断する場合は、「キャンセル」をクリックします。抽出が完了すると、ダイアログ・ボックスが閉じられ、結果がパネルに表示されます。詳しくは、『タイプ・パターンおよびコンセプト・パターン』のトピックを参照してください。

タイプ・パターンおよびコンセプト・パターン

パターンは、2 つの部分、コンセプトとタイプを組み合わせて構成されています。パターンは、特定のサブジェクトに関する意見またはコンセプト間の関連性を探索する場合に最も役立ちます。例えば、競合他社の製品名を抽出しても、重要でない場合があります。この場合、抽出したパターンを参照し、ドキュメントまたはレコードに、製品が良い、悪い、または高いことを示すテキストが含まれている例があるかどうかを確認することができます。

パターンは最大 6 つのタイプまたは 6 つのコンセプトから構成されます。そのため、両方のパターンのワイン同枠の行には、最大 6 つのスロットまたは場所があります。各スロットは、言語リソースで定義されているように、TLA パターンの要素固有の場所に対応しています。インタラクティブ・ワークベンチでは、スロットに値がない場合、スロットはテーブルに表示されません。例えば、最も長いパターン結果に 4 つのスロットがある場合、後半 2 つのスロットは表示されません。詳しくは、215 ページの『第 18 章 テキスト リンク規則について』のトピックを参照してください。

パターン結果を抽出する場合、まずタイプ レベルでグループ化され、コンセプト・パターンに分割されます。そのため、「タイプ パターン」(左上) および「コンセプト パターン」(左下) の 2 つの結果パネルが表示されます。返されたすべてのコンセプト・パターンを表示するには、タイプ・パターンをすべて選択します。一番下のコンセプト・パターンのパネルには、順位の最大値(「フィルター」ダイアログ・ボックスで定義)までのコンセプト・パターンがすべて表示されます。

タイプ・パターン TLA パターン規則を満たす 1 つまたは複数の関連タイプで構成されているパターン規則が表示されます。タイプ・パターンは、<組織名> + <地名> + <肯定的> と表され、特定の場所の組織について、肯定的なフィードバックを提供します。シンタックスは、次のようになります。

```
<Type1> + <Type2> + <Type3> + <Type4> + <Type5> + <Type6>
```

コンセプト・パターン 上の「タイプ・パターン」で現在選択されているすべてのタイプ・パターンのコンセプト・レベルでパターンの結果が表示されます。コンセプト・パターンは、ホテル + パリ + すばらしいなどの構造に従います。シンタックスは、次のようになります。

```
concept1 + concept2 + concept3 + concept4 + concept5 + concept6
```

パターン結果が 6 つ未満の最大スロットを使用する場合、必要な数だけのスロット (または列) が表示されます。<タイプ 1>+<>+<タイプ 2>+<>+<>+<> のように、2 つの入力されたスロットの間の空白のスロットは破棄され、<Type1>+<Type3> のようになります。コンセプト・パターンの場合、コンセプト 1+コンセプト 2 となります。(、は空値を示します)。

カテゴリとコンセプト・ビューの抽出結果と同様、ここで結果を確認できます。これらのパターンを構成するタイプおよびコンセプトに調整を行う場合、カテゴリとコンセプト・ビューの抽出結果ペインまたはリソース・エディターで変更を行うか、パターンを再抽出します。カテゴリ定義でコンセプト、タイプ、

またはパターンが使用されている場合、カテゴリーまたは条件規則のアイコンが「パターン」テーブルまたは「抽出結果」テーブルの「投入」列に表示されます。

注: 表示中のペインに結果を表示しきれない場合は、ペインの下部にあるコントロールを使用して前後の結果に移動したり、移動先のページ番号を入力したりすることができます。

TLA 結果のフィルタリング

非常に大きなデータセットを処理する場合、抽出プロセスでは、多数の結果が作成される場合があります。多くのユーザーによって、多数の結果が作成されると、結果を効率的に確認することが困難になります。ただし、これらの結果をフィルタリングして、最も関心の高い結果に焦点を当てることができます。「フィルター」ダイアログ・ボックスの設定を変更して、表示されるパターンを制限できます。これらの設定はすべていっしょに使用されます。

TLA ビューの「フィルター」ダイアログ・ボックスには以下の領域とフィールドが含まれます。

出現頻度でフィルタリング フィルタリングを実行して、特定のグローバル出現頻度値またはドキュメントの出現頻度の値を持つ結果のみを表示できます。

- グローバル出現頻度は、パターンがドキュメントまたはレコードの全体的なセットに出現する回数の合計で、「グローバル」列に表示されます。
- ドキュメント出現頻度は、パターンが出現するドキュメントまたはレコードの合計数で、「ドキュメント」列に表示されます。

例えば、あるパターンが 500 件のレコードに 300 回出現した場合、このパターンのグローバル出現頻度は 300 で、ドキュメント出現頻度は 500 となります。

マッチ テキスト別 **AND** 条件 ここで定義する規則に一致する結果のみを表示できます。「マッチ テキスト」フィールドに合致する文字のセットを入力し、スロット番号またはそれらのすべてを特定して、コンセプト名またはタイプ名のどちらでこのテキストを検索するかを選択します。合致を適用する条件を選択します (タイプ名の開始と終了を示す各カッコを使用する必要はありません)。ドロップダウン・リストから「**AND**」または「**OR**」を選択して条件規則が両方の文またはいずれかに一致するようにし、最初の文と同じ方法で、2 番目のテキスト・マッチ文を定義します。

表 35. マッチ・テキストの条件

条件	説明
含む	文字列が任意の場所で出現する場合、テキストが一致します(デフォルトの選択)。
開始	コンセプトまたはタイプが特定のテキストで始まる場合にのみ、テキストが一致します。
終了	コンセプトまたはタイプが特定のテキストで終わる場合にのみ、テキストが一致します。
完全一致	文字列全体が、コンセプト名またはタイプ名に一致する必要があります。

パターン・ペインに表示される結果

ソフトウェアの英語版を使用しえいとします。フィルタリングに基づいて、結果が「パターン」ウィンドウにどのように表示されるかについて、いくつか例を示します。



図 33. フィルターの結果の例 1

この例では、フィルターで指定された順位の最大値により、返されるパターン数が制限されていることがツールバーに示されています。紫色のアイコンが表示されている場合、パターンの最大数に達していることを示します。アイコンの上にポインタを置くと、詳細が表示されます。「順位別 **AND** 条件」 フィルターに関する前述の説明を参照してください。

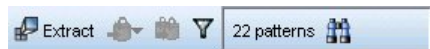


図 34. フィルターの結果の例 2

この例では、マッチ・テキスト・フィルターを使用して、結果が制限されていることがツールバーに示されています (虫めがねのアイコンを参照)。アイコンにポインタを置くと、マッチ・テキストの内容が表示されます。

結果を絞り込むには

1. メニューの「ツール」 > 「フィルター」を選択します。「フィルター」ダイアログ・ボックスが開きます。
2. 使用するフィルターを選択および調整します。
3. 「OK」 をクリックするとフィルターが適用され、新しい結果が表示されます。

データ パネル

テキスト リンク分析パターンを抽出および検証すると、作業しているデータをいくつか確認したい場合があります。例えば、パターンのグループが発見された実際のレコードを確認したい場合があります。右下のデータ・ペインでレコードまたはドキュメントを確認することができます。デフォルトで表示されない場合は、メニューから「表示」 > 「パネル」 > 「データ」を選択してください。

データ・ペインには、特定の表示制限に応じて、ビュー内の選択に該当するドキュメントまたはレコードごとに 1 行ずつ表示されます。デフォルトでは、データ・ペインに表示されるドキュメントまたはレコード数が制限され、データをより迅速に表示できるようになります。ただし、これは「オプション」ダイアログ・ボックスで調整できます。詳しくは、76 ページの『オプション: 「セッション」 タブ』のトピックを参照してください。

注: 表示中のペインに結果を表示しきれない場合は、ペインの下部にあるコントロールを使用して前後の結果に移動したり、移動先のページ番号を入力したりすることができます。

データ・ペインの表示および更新

データ・ペインでは、大きなデータセットの自動データ更新には時間がかかるため、自動的に表示の更新は行われません。そのため、このビューでタイプ・パターンまたはコンセプト・パターンを選択すると、「表示」 をクリックしてデータ・ペインの内容を更新できます。

テキスト・ドキュメントまたはレコード

テキスト・データがレコードの形式で、テキストの長さが比較的短い場合、データ・ペインのテキスト・フィールドには、テキスト・データの全体が表示されます。ただし、レコードおよび大きいデータセットを処理している場合、テキスト・フィールドの列にはテキストの一部が表示され、右側のテキスト・プレビュー・ペインを開くと、テーブルで選択したレコードの大部分またはすべてが表示されます。テキスト・データが個別ドキュメントの形式の場合、データ・ペインには、ドキュメントのファイル名が表示されます。ドキュメントを選択すると、テキスト・プレビュー・ペインには選択したドキュメントのテキストが表示されます。

色および強調表示

データを表示すると、該当するドキュメントまたはレコードのコンセプトおよび記述子は色付きで強調表示され、テキスト内のコンセプトおよび記述子を特定しやすくします。カラー・コードは、コンセプトが属するタイプに対応します。カラーコード化された項目上でマウス・ポインタを停止させて、項目が抽出されたコンセプトと、項目が割り当てられたタイプを表示することもできます。抽出されていないテキストは、黒で表示されます。通常、こうした抽出されていない単語は接続詞（「および」または「と」）、代名詞（「私」または「彼ら」）および動詞（「いる」、「持つ」、または「取る」）のケースが多くあります。

データ・ペインの列

テキスト・フィールドの列は常に表示されていますが、その他の列も表示できます。その他の列を表示するには、メニューで「表示」 > 「データ・ペイン」を選択し、データ・ペインに表示したい列を選択します。表示できるのは次の列です。

- 「テキスト・フィールド名」(#)/ドキュメント コンセプトおよびタイプが抽出されたテキスト・データの列を追加します。データがドキュメントにある場合、列名は「ドキュメント」となり、ドキュメント・ファイル名または完全パスのみが表示されます。これらのドキュメントのテキストを表示するには、テキスト・プレビュー・ペインを表示する必要があります。データ・ペインの行数が、この列名の後のカッコ内に表示されます。読み込みの速度向上のために使用される「オプション」ダイアログの制約により、一部のドキュメントまたはレコードが表示されない場合があります。最大値に達すると、数値の後に「- 最大」と表示されます。詳しくは、76 ページの『オプション：「セッション」タブ』のトピックを参照してください。
- カテゴリー レコードが属するカテゴリーがそれぞれ表示されます。この列を表示する場合、最新の情報を示すため、データ・ペインの更新に少し時間がかかる場合があります。
- 適合度順位 1 つのカテゴリーの各レコードの順位が表示されます。この適合度順位は、カテゴリー内の他のレコードと比較して、レコードがカテゴリーにどれだけ適合しているかを示します。カテゴリー・ペイン (左上のパネル) でカテゴリーを選択すると、順位が表示されます。詳しくは、105 ページの『カテゴリーの関連性』のトピックを参照してください。
- カテゴリーの個数 レコードが割り当てられているカテゴリー数が表示されます。

第 12 章 グラフの視覚化

カテゴリとコンセプト・ビュー、クラスター・ビュー、およびテキスト リンク分析ビューにはすべて、ウィンドウの右上隅に視覚化ペインがあります。このパネルを使用して、データを視覚的に検証することができます。次のグラフおよび図表を使用できます。

- **カテゴリとコンセプト・ビュー:** このビューには、**カテゴリ・バー**、**カテゴリ Web**、および **カテゴリ Web テーブル** の 3 つのグラフと図表があります。このビューでは、「表示」をクリックする場合にのみグラフが更新されます。詳しくは、トピック「『カテゴリ・グラフおよび図表』」を参照してください。
- **クラスター・ビュー:** このビューには、**コンセプト Web グラフ** および **クラスター Web グラフ** の 2 つの Web グラフがあります。詳しくは、トピック「157 ページの『クラスター・グラフ』」を参照してください。
- **テキスト リンク分析ビュー:** このビューには、**コンセプト Web グラフ** および **タイプ Web グラフ** の 2 つの Web グラフがあります。詳しくは、158 ページの『テキスト リンク分析のグラフ』のトピックを参照してください。

グラフの編集に使用するすべての一般的なツールバーおよびパレットの詳細は、オンライン・ヘルプまたはファイル *ModelerSPOnodes.pdf* のグラフの編集に関するセクションを参照してください。このファイルは、製品ダウンロードの一部として使用できます。

カテゴリ・グラフおよび図表

カテゴリを作成する場合、時間をかけてカテゴリ定義、含まれるドキュメントまたはレコード、およびカテゴリの重複を確認することが重要になります。視覚化ペインには、カテゴリに関するいくつかの視点が表示されています。視覚化ペインは、カテゴリとコンセプト・ビューの右上隅に表示されます。表示されない場合、「表示」メニュー（「表示」>「パネル」>「視覚化」）からこのパネルにアクセスできます。

このビューの視覚化ペインには、ドキュメントまたはレコードのカテゴリ化における共通性について 3 つの視点が表示されています。このパネルの図表やグラフを使用して、カテゴリ化の結果を分析したり、カテゴリまたはレポートの調整を行うことができます。カテゴリを調整する場合、このパネルを使用して、カテゴリ定義を確認し、あまりに類似している（ドキュメントまたはレコードの 75% 以上を共有しているなど）またはあまりに異なるカテゴリを明らかにできます。2 つのカテゴリがあまりに似ている場合、2 つのカテゴリの結合することができます。また、一方のカテゴリから特定の記述子を削除して、カテゴリ定義の調整することもできます。

抽出結果ペイン、カテゴリ ペインまたは「カテゴリ定義」ダイアログ ボックスで選択した内容に応じて、このパネルの各タブで、ドキュメント/レコードとカテゴリの間の該当する交互作用を表示できます。それぞれは、同様の情報を、異なる方法で、またはさまざまなレベルの詳細情報とともに表示されます。ただし、現在選択している部分のグラフを更新するには、選択を行ったパネルまたはダイアログ・ボックスのツールバーの「表示」をクリックします。

カテゴリとコンセプト・ビューの視覚化ペインには、次のようなグラフおよび図表が表示されます。

- **カテゴリー棒グラフ:** テーブルおよび棒グラフを使用して、選択に該当するドキュメント/レコードと関連するカテゴリーとの間の重なりを示します。また、棒グラフは、カテゴリー内のドキュメント/レコード数の、ドキュメント/レコード数の合計に対する比率を示します。詳しくは、『カテゴリー棒グラフ』のトピックを参照してください。
- **カテゴリー Web グラフ:** このグラフは、その他のパネルの選択部分に従って、ドキュメント/レコードが属するカテゴリーのドキュメント/レコードの重なりを示します。詳しくは、『カテゴリー Web グラフ』のトピックを参照してください。
- **カテゴリー Web テーブル:** このテーブルは、「カテゴリー Web グラフ」タブと同じ情報をテーブル形式で表示します。このテーブルには 3 つの列があり、列の見出しをクリックするとソートできます。詳しくは、157 ページの『カテゴリー Web テーブル』のトピックを参照してください。

詳しくは、95 ページの『第 9 章 テキストデータのカテゴリー化』のトピックを参照してください。

カテゴリー棒グラフ

このタブには、選択に該当するドキュメント/レコードと関連するカテゴリーとの間の重なりを示すテーブルおよび棒グラフが表示されます。また、棒グラフは、カテゴリー内のドキュメント/レコード数の、ドキュメントまたはレコード数の合計に対する比率も示します。このグラフのレイアウトは編集できません。ただし、列の見出しをクリックして、列をソートすることはできます。

テーブルには、次の列が表示されます。

- **カテゴリー:** 選択したカテゴリーの名前が表示されます。デフォルトでは、選択した中で最も一般的なカテゴリーが最初に表示されます。
- **棒グラフ:** 指定されたカテゴリーのドキュメントまたはレコード数の、ドキュメントまたはレコード数の合計に対する比率を視覚的に表示します。
- **選択 %:** カテゴリーのドキュメントまたはレコード数の合計の、選択部分に表示されたドキュメントまたはレコード数の合計に対する比率に基づいたパーセンテージを示します。
- **ドキュメント:** 指定したカテゴリーの選択部分のドキュメントまたはレコードの数を示します。

カテゴリー Web グラフ

このタブには、カテゴリー Web グラフが表示されます。Web グラフは、その他のパネルの選択部分に従って、ドキュメントまたはレコードが属するカテゴリーのドキュメントまたはレコードの重なりを示します。カテゴリー・ラベルがある場合は、グラフにこれらのラベルが表示されます。このパネルのツールバー・ボタンを使用して、グラフのレイアウト (ネットワーク、サークル、有向、またはグリッド) を選択できます。

Web グラフで、各ノードはカテゴリーを示します。マウスを使用し、パネル内のノードを選択して移動できます。ノードのサイズは、選択部分のカテゴリーのドキュメントまたはレコードの数に基づいた相対的なサイズを示します。カテゴリー間の線の太さと色は、含まれている共通のドキュメントまたはレコードの数を示します。探索的分析モードでノードの上にマウス・ポインタを停止させると、ヒントにカテゴリーの名前 (またはラベル) およびカテゴリー内のドキュメントまたはレコードの全体数が表示されます。

注: デフォルトでは、ノードを移動できるグラフの探索的分析モードが有効化されています。ただし、編集モードに切り替えて、色、フォント、凡例など、グラフのレイアウトを編集できます。詳しくは、159 ページの『グラフのツールバーおよびパレットの使用』を参照してください。

「視覚化データをコピー」ボタンを使用してグラフのデータをコピーする場合に、データをスプレッドシートやテキスト エディターに貼り付けると、データに V1、V2、... V7 の列見出しが付くことがわかります。これらの列には以下の情報が格納されます。

- **V1、V2** これらの値は画面座標 (それぞれ X と Y) に対応します。
- **V3、V5** カテゴリーのコンセプトをリストします。
- **サイズ、V6** コンセプトが見つかったドキュメントの数を示します。
- **V7** 現在は未使用です。

カテゴリー Web テーブル

このタブには、「カテゴリー Web グラフ」タブと同じ情報をテーブル形式で表示されます。このテーブルには次の 3 つの列があり、列の見出しをクリックするとソートできます。

- **度数:** 2 つのカテゴリーで共有している、または共通のドキュメントまたはレコードの数を表示します。
- **カテゴリー 1:** 最初のカテゴリーの名前、そして含まれるドキュメントまたはレコード数の合計がカッコ内に表示されます。
- **カテゴリー 2:** 2 番目のカテゴリーの名前、そして含まれるドキュメントまたはレコード数の合計がカッコ内に表示されます。

クラスター・グラフ

クラスターを作成した後、視覚化ペインの Web グラフで視覚的にクラスターを検証できます。視覚化ペインは、「コンセプト Web グラフ」および「クラスター Web グラフ」の 2 つのクラスター化のパーспекティブを提供します。このパネルで Web グラフを使用して、クラスターリングの結果を分析し、カテゴリーに追加するコンセプトおよび規則を見つけることができます。視覚化ペインは、クラスター・ビューの右上隅にあります。表示されない場合、「表示」メニュー（「表示」>「パネル」>「視覚化」）からこのパネルにアクセスできます。クラスター・ペインでクラスターを選択すると、視覚化ペインに該当するグラフを自動的に表示できます。

注:デフォルトでは、グラフはインタラクティブ/選択モードで、ノードを移動できます。ただし、編集モードで、色、フォント、凡例など、グラフのレイアウトを編集できます。詳しくは、159 ページの『グラフのツールバーおよびパレットの使用』のトピックを参照してください。

クラスター・ビューには 2 つの Web グラフがあります。

- **コンセプト Web グラフ:**このグラフには、選択したクラスター内のすべてのコンセプトおよびクラスター外のリンクしたコンセプトが表示されます。このグラフを使用して、クラスター内のコンセプトがどのようにリンクしているかと、外部リンクを確認することができます。詳しくは、『コンセプト Web グラフ』のトピックを参照してください。
- **クラスター Web グラフ:**このグラフには、選択したクラスターと、表示される選択したクラスター間のすべての外部リンク (点線で表示) を表示します。詳しくは、158 ページの『クラスター Web グラフ』のトピックを参照してください。

詳しくは、143 ページの『第 10 章 クラスターの分析』のトピックを参照してください。

コンセプト Web グラフ

このタブには、クラスター外のリンクしたコンセプトのほか、選択したクラスター内のすべてのコンセプトを表示する Web グラフが表示されます。このグラフを使用して、クラスター内のコンセプトがどのよう

にリンクしているかと、外部リンクを確認することができます。 クラスター内の各コンセプトはノードとして表示され、タイプの色によって色分けされます。 詳しくは、 191 ページの『キーワード辞書の作成』のトピックを参照してください。

クラスター内のコンセプト間の内部リンクが描画され、各リンクの線の太さは、グラフ・ツールバーの選択に応じて、各コンセプト・ペアの共起のドキュメント数または類似度リンク値に直接関連します。 クラスターのコンセプトおよびクラスター外のこれらのコンセプト間の外部リンクも表示されます。

「クラスター定義」ダイアログ・ボックスでコンセプトを選択した場合、コンセプト Web グラフには、これらのコンセプトと、それらに関連する内部リンクおよび外部リンクが表示されます。 選択したコンセプトのいずれかを含まないその他のコンセプト間のリンクは、グラフには表示されません。

注: デフォルトでは、グラフはインタラクティブ/選択モードで、ノードを移動できます。 ただし、編集モードで、色、フォント、凡例など、グラフのレイアウトを編集できます。 詳しくは、 159 ページの『グラフのツールバーおよびパレットの使用』を参照してください。

「視覚化データをコピー」 ボタンを使用してグラフのデータをコピーする場合に、データをスプレッドシートやテキスト エディターに貼り付けると、データに V1、V2、... V7 の列見出しが付くことがわかります。 これらの列には以下の情報が格納されます。

- **V1、V2** これらの値は画面座標 (それぞれ X と Y) に対応します。
- **V3、V6** コンセプト タイプをリストします。
- **V4、V5** コンセプト ラベルを示します。
- **V7** 現在は未使用です。

クラスター Web グラフ

このタブには、選択したクラスターを示す Web グラフが表示されます。 他のクラスター間のリンクのほか、選択したクラスター間の外部リンクがすべて点線で表示されます。 クラスター Web グラフでは、各ノードはクラスター全体を表し、線の太さは、2 つのクラスター間の外部リンク数を示します。

重要: クラスター Web グラフを表示するには、外部リンクを持つクラスターを先に構築する必要があります。 外部リンクは、別々のクラスターにあるコンセプトのペア (あるクラスターのコンセプトと外部の別のクラスターのコンセプトとの間) 間のリンクです。

例えば、2 つのクラスターがあるとします。 クラスター A には 3 つのコンセプト、A1、A2、および A3 があります。 クラスター B には 2 つのコンセプト、B1 および B2 があります。 コンセプト間のリンクは、A1-A2、A1-A3、A2-B1 (外部)、A2-B2 (外部)、A1-B2 (外部)、および B1-B2 です。 クラスター Web グラフでは、線の太さが 3 つの外部リンクを示すことを意味します。

注: デフォルトでは、グラフはインタラクティブ/選択モードで、ノードを移動できます。 ただし、編集モードで、色、フォント、凡例など、グラフのレイアウトを編集できます。 詳しくは、 159 ページの『グラフのツールバーおよびパレットの使用』のトピックを参照してください。

テキスト リンク分析のグラフ

テキスト リンク分析 (TLA) パターンを抽出した後、視覚化ペインの Web グラフで視覚的にパターンを検証できます。 視覚化パネルは、コンセプト (パターン) Web グラフ、およびタイプ (パターン) Web グラフ の 2 つの TLA パターンのパースペクティブを提供します。 このパネルの Web グラフを使用して、パターンを視覚的に表すことができます。 視覚化ペインは、テキスト リンク分析の右上隅にあります。 表

示されない場合、「表示」メニュー（「表示」>「パネル」>「視覚化」）からこのパネルにアクセスできません。選択項目がない場合、グラフ領域が空になります。

注:デフォルトでは、グラフはインタラクティブ/選択モードで、ノードを移動できます。ただし、編集モードで、色、フォント、凡例など、グラフのレイアウトを編集できます。詳しくは、『グラフのツールバーおよびパレットの使用』のトピックを参照してください。

テキスト リンク分析ビューには 2 つの Web グラフがあります。

- **コンセプト Web グラフ**。このグラフには、選択したパターンのすべてのコンセプトを示します。コンセプト グラフの線の幅およびノードのサイズ (タイプ・アイコンが表示されていない場合) には、選択したテーブルのグローバル出現値を示します。詳しくは、『コンセプト Web グラフ』のトピックを参照してください。
- **タイプ Web グラフ**。このグラフには、選択したパターンのすべてのタイプを示します。グラフの線の幅およびノードのサイズ (タイプ・アイコンが表示されていない場合) には、選択したテーブルのグローバル出現値を示します。ノードは、タイプ カラーまたはアイコンによって示されます。詳しくは、『タイプ Web グラフ』のトピックを参照してください。

詳しくは、149 ページの『第 11 章 テキスト リンク分析の検証』のトピックを参照してください。

コンセプト Web グラフ

Web グラフは、現在の選択で示されているすべてのコンセプトを表示します。例えば、3 つの一致コンセプト・パターンがあるタイプ・パターンを選択した場合、このグラフにはリンクしたコンセプトが 3 セット表示されます。コンセプト グラフの線の幅およびノード・サイズは、グローバル頻度値を示します。グラフには、パターンのパネルで選択されたものと同じ情報が表示されます。各コンセプトのタイプは、グラフ・ツールバーの選択内容に応じて、色またはアイコンによって表示されます。詳しくは、『グラフのツールバーおよびパレットの使用』のトピックを参照してください。

タイプ Web グラフ

この Web グラフは、現在の選択の各タイプ・パターンを示します。例えば、2 つのコンセプト・パターンを選択した場合、このグラフには選択したパターンのタイプごとに 1 つのノードおよび同じパターンで見つかったタイプ間のリンクを示します。線の幅およびノード・サイズは、セットのグローバル頻度値を示します。グラフには、パターンのパネルで選択されたものと同じ情報が表示されます。タイプは、グラフに表示されるタイプ名のほか、グラフ ツールバーで選択した内容によって色またはタイプ アイコンによっても識別されます。詳しくは、『グラフのツールバーおよびパレットの使用』のトピックを参照してください。

グラフのツールバーおよびパレットの使用

各グラフにツールバーがあり、グラフに行うさまざまな操作を実行できるいくつかの共通パレットをすぐに使用できます。各ビュー (カテゴリーおよびコンセプト、クラスター、テキスト リンク分析) には、若干異なるツールバーがあります。探索的分析ビュー・モードまたは編集ビュー・モードから選択できます。

探索モードでは、視覚化によって表現されたデータや値を分析的に検討することができます。一方、編集モードでは、視覚化のレイアウトや外観を変更することができます。例えば、フォントや色を自分の組織のスタイル・ガイドに合わせて変更することが可能です。このモードを選択するには、メニューから「表示」>「視覚化パネル」>「編集モード」を選ぶか、ツールバーにあるアイコンをクリックします。

編集モードには、視覚化のレイアウトのさまざまな要素に影響を与えるいくつかのツールバーがあります。使用しないものがある場合は、そのツールバーを非表示にしてダイアログ・ボックスにおけるグラフの表示領域を増やすことができます。ツールバーの選択または選択解除を行うには、「表示」メニューで目的のツールバーまたはパレットの名前をクリックします。

グラフの編集に使用するすべての一般的なツールバーおよびパレットの詳細は、オンライン ヘルプまたはファイル *ModelerSPOnodes.pdf* の視覚化の編集に関するセクションを参照してください。このファイルは、製品ダウンロードの一部として入手できます。

表 36. *Text Analytics* のツールバー・ボタン :











ボタン/リスト	説明
	編集モードを有効化します。編集モードに切り替えて、フォントの拡大、会社のスタイル・ガイドに合った色への変更、またはラベルや凡例の削除など、グラフの外観を変更できます。
	探索的分析モードを有効化します。デフォルトでは、探索的分析モードがオンになっており、グラフの周囲でノードを移動およびドラッグ、そしてグラフ・オブジェクトの上でマウス・ポインタを停止させて、ヒントの詳細情報を表示できます。
	<p>カテゴリとコンセプト・ビューおよびテキスト リンク分析ビューでグラフの Web 表示の種類を選択します。</p> <ul style="list-style-type: none"> • サークル レイアウト どのようなグラフにも適用できる一般的なレイアウト。リンクに方向がないことを想定してグラフをレイアウトし、すべてのノードを同様に扱います。ノードは円の周囲にのみ配置されます。 • ネットワーク レイアウト どのようなグラフにも適用できる一般的なレイアウト。リンクに方向がないことを想定してグラフをレイアウトし、すべてのノードを同様に扱います。ノードは、レイアウト内に自由に配置されます。 • 有向レイアウト 方向のあるグラフにのみ使用できるレイアウト。このレイアウトは、ルート・ノードからリーフ・ノードへのツリー上の構造を作成し、色別に構成します。このレイアウトを使用すると、階層データが適切に表示されます。 • グリッド レイアウト どのようなグラフにも適用できる一般的なレイアウト。リンクに方向がないことを想定してグラフをレイアウトし、すべてのノードを同様に扱います。ノードは領域内のグリッド ポイントにのみ配置されます。
	<p>リンクの太さの表示。グラフで線の太さが示す内容を選択します。これは、クラスター・ビューにのみ適用されます。クラスター Web グラフは、クラスター間の外部リンク数のみを表示します。以下から選択できます。</p> <ul style="list-style-type: none"> • 類似度 線の太さは、2 つのクラスター間の外部リンク数を示します。 • 共起 記述子の共起が出現するドキュメント数を示します。
	凡例を表示する切り替えボタン。ボタンを押さない場合、判例は表示されません。
	タイプの色ではなくグラフ内のタイプのアイコンを表示する切り替えボタン。これは、テキスト リンク分析ビューにのみ適用されます。
	グラフの下にリンク スライダーを表示する切り替えボタン。矢印をスライドして、結果を絞り込むことができます。
	サブカテゴリではなく、選択されたカテゴリの最上位レベルのグラフが表示されます。

表 36. Text Analytics のツールバー・ボタン (続き):

ボタン/リスト	説明
	<p>選択されたカテゴリーの最下位レベルのグラフが表示されます。</p>
	<p>サブカテゴリーの名前を出力内でどのように表示するかを制御します。</p> <ul style="list-style-type: none"> • 完全カテゴリー・パス カテゴリー名と、該当する場合、カテゴリー名とサブカテゴリー名をスラッシュを使用して区切り、上位カテゴリーの完全パスを出力します。 • 省略したカテゴリー パス ただし、省略記号を使用して、該当するカテゴリーの上位カテゴリー数を示します。 • 下位レベルのカテゴリー 完全パスまたは上位カテゴリーを表示せず、カテゴリー名のみを出力します。

第 13 章 セッション・リソース・エディター

IBM SPSS Modeler Text Analytics は、主要キーワードをテキスト・データから迅速にかつ正確にキャプチャーします。この抽出プロセスは、テキストデータからの情報抽出を管理する言語リソースに大きく依存しています。デフォルトでは、これらのリソースはリソース・テンプレートによって決まります。

IBM SPSS Modeler Text Analytics は、言語リソースおよび非言語リソースを含む、ライブラリーおよび高度なリソースの形式で専門的なリソース・テンプレートのセットに付属しており、データの処理方法および抽出方法を定義できます。詳しくは、167 ページの『第 14 章 テンプレートとリソース』のトピックを参照してください。

ノードのダイアログ・ボックスで、テンプレートのリソースをノードに読み込むことができます。一度インタラクティブ・ワークベンチ・セッションに入ると、必要に応じてこのノードのデータ用にこれらのリソースをカスタマイズできます。インタラクティブ・ワークベンチ・セッションの間、リソース・エディタービューでリソースを作業できます。インタラクティブ・セッションが起動すると、ノードにデータおよび抽出結果をキャッシュしていない場合は、ノードのダイアログ・ボックスで読み込まれたリソースを使用して抽出を実行します。

リソース・エディターを使用したリソースの編集

リソース・エディターでは、インタラクティブ・ワークベンチ・セッションの抽出結果 (コンセプト、タイプ、およびパターン) を作成に使用するリソースのセットへのアクセスが用意されています。このエディターは、テンプレート・エディターと非常に類似していますが、リソース・エディターでは、インタラクティブ・ワークベンチ・セッションでリソースを編集するという点で異なります。リソースの作業および実行した他の作業を終了している場合、モデル作成ノードを更新してこの作業を保存し、後続のインタラクティブ・ワークベンチ・セッションで復元できます。詳しくは、78 ページの『モデル作成ノードの更新および保存』のトピックを参照してください。

ノードにリソースを読み込むために使用するテンプレートで直接作業する場合は、テンプレート・エディターを使用することをお勧めします。次のようなリソース・エディター内で実行できる多くのタスクは、テンプレート・エディターと同じように実行されます。

- ライブラリーの使用。詳しくは、179 ページの『第 15 章 ライブラリーの使用』のトピックを参照してください。
- キーワード辞書の作成。詳しくは、191 ページの『キーワード辞書の作成』のトピックを参照してください。
- キーワードを辞書に追加。詳しくは、192 ページの『キーワードの追加』のトピックを参照してください。
- 類義語の作成。詳しくは、198 ページの『類義語の定義』のトピックを参照してください。
- テンプレートのインポートおよびエクスポート。詳しくは、175 ページの『テンプレートのインポートおよびエクスポート』のトピックを参照してください。
- ライブラリーの公開。詳しくは、186 ページの『ライブラリーの公開』のトピックを参照してください。

オランダ語、英語、フランス語、ドイツ語、イタリア語、ポルトガル語、スペイン語のテキストの場合

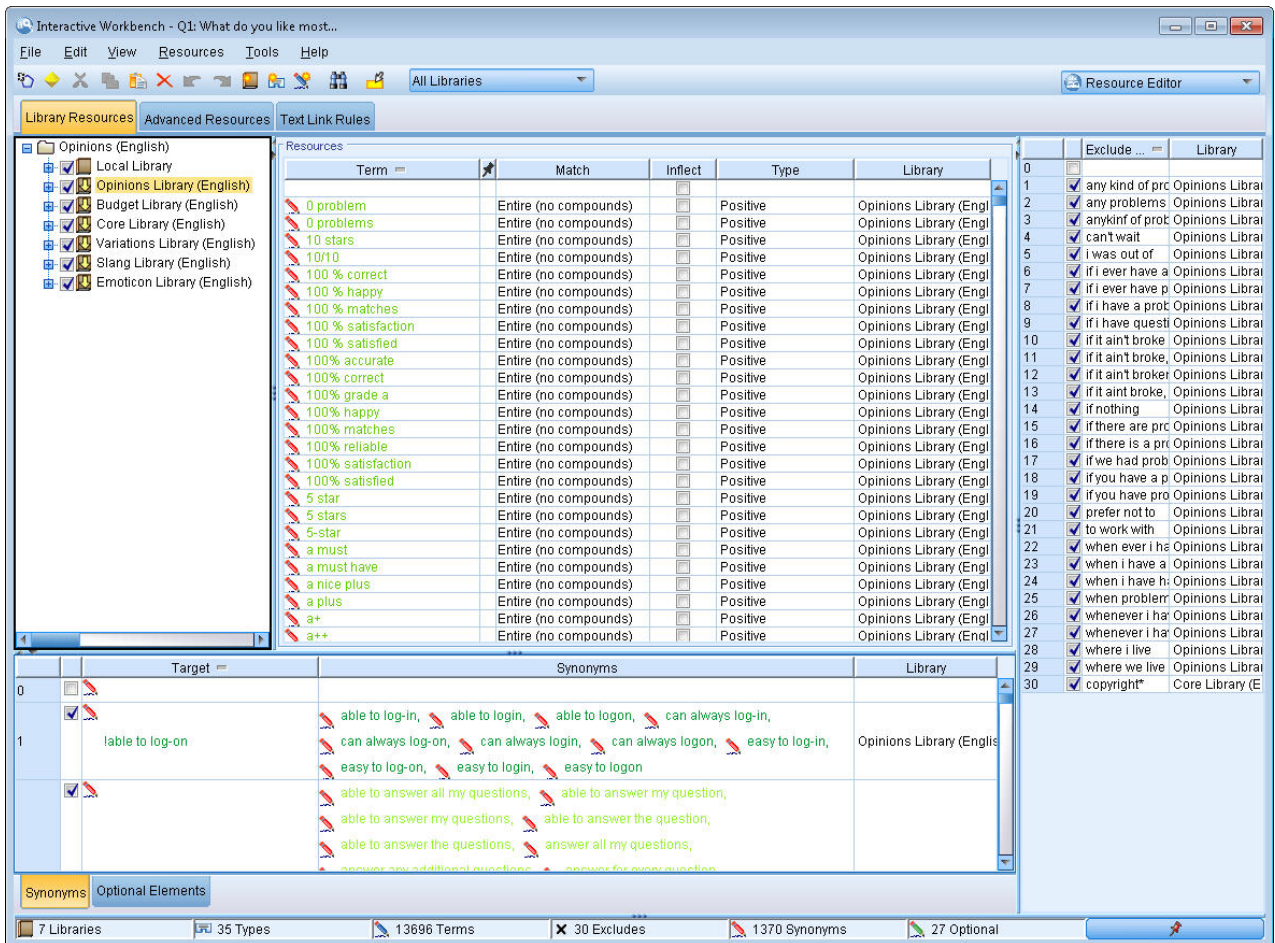


図 35. リソース・エディター・ビュー

テンプレートの作成および更新

リソースに変更を行い、今後それらを再利用したい場合、リソースをテンプレートとして保存することができます。保存する場合は、既存のテンプレート名を使用して保存するのか、新しい名前を付けるのかを選択できます。その後、テンプレートを読み込むと、同じリソースを取得することができます。詳しくは、27 ページの『テンプレートおよび TAP からのリソースのコピー』のトピックを参照してください。

注:ライブラリーを公開して共有することもできます。詳しくは、184 ページの『ライブラリーの共有』のトピックを参照してください。

テンプレートを作成 (または更新) するには

1. リソース・エディター ビューのメニューで、「リソース」>「リソース テンプレートを作成」を選択します。「リソース・テンプレートを作成」ダイアログ・ボックスが開きます。
2. 新しいテンプレートを作成する場合は、「テンプレート名」フィールドに新しい名前を入力します。既存のテンプレートとその時点で読み込まれているリソースで上書きする場合は、テーブルでテンプレートを選択します。
3. 「保存」をクリックして、テンプレートを作成します。

重要: テンプレートはノードで選択されたときに読み込まれ、ストリームが実行されているときは読み込まれないため、最新の変更を取得したい場合にリソース・テンプレートを使用するその他のノードのリソース・テンプレートを再読み込みしてください。詳しくは、173 ページの『読み込み後のノード・リソースの更新』のトピックを参照してください。

リソース・テンプレートの切り替え

現在セッションで読み込まれたリソースを別のテンプレートからのコピーに置き換えたい場合、これらのリソースに切り替えることができます。これによって、現在セッション内に読み込まれているリソースを上書きすることができますリソースを切り替えて定義済みのテキスト リンク分析 (TLA) パターン・ルールを使用する場合、必ずTLA列内でマークされたテンプレートを選択してください。

セッション作業 (カテゴリー、パターン、リソース) を復元したいが他のセッション作業を失わずにテンプレートからリソースのコピーを読み込んで更新したい場合特に、リソースの切り替えが役立ちます。リソース・エディター に内容をコピーしたいテンプレートを選択し、「OK」をクリックします。これにより、このセッションのリソースが置き換えられます。セッションのリソースが置き換えられます。次回インタラクティブ・ワークベンチ・セッションを起動するときにこれらの変更を保持したい場合、セッションの終わりにモデル作成ノードを更新してください。

注: インタラクティブ・セッションで別のテンプレートの内容に切り替える場合、ノードに表示されるテンプレートの名前は、最後に読み込まれ、コピーされたテンプレートの名前となります。これらのリソースまたは他のセッション作業を利用するには、セッションを終了する前にモデル作成ノードを更新し、ノードで「セッション作業を使用」 オプションを選択します。詳しくは、78 ページの『モデル作成ノードの更新および保存』のトピックを参照してください。

リソースを切り替えるには

1. リソース・エディター ビューのメニューで、「リソース」>「リソース テンプレートを切り替え」を選択します。「リソースを切り替え」ダイアログ・ボックスが開きます。
2. テーブルに表示されたテンプレートから、使用したいテンプレートを選択します。
3. 「OK」をクリックして、現在読み込まれているこれらのリソースを中止し、代わりに選択したテンプレートのリソースのコピーを読み込みます。リソースに変更を行い、今後使用するためにライブラリーを保存したい場合、切り替える前にそれらを公開、更新、共有することができます。詳しくは、184 ページの『ライブラリーの共有』のトピックを参照してください。

第 14 章 テンプレートとリソース

IBM SPSS Modeler Text Analytics は、主要キーワードをテキスト・データから迅速にかつ正確にキャプチャーします。この抽出プロセスは、テキスト データからの情報抽出方法を決定するために言語リソースに大きく依存しています。詳しくは、5 ページの『抽出の方法』のトピックを参照してください。リソース・エディター ビューで、これらのリソースを調整できます。

ソフトウェアをインストールすると、専門的なリソース セットも取得します。これらの付属リソースは、特定の言語と特定の応用分野で、数年にわたる調査と調整の結果得られたもので、ユーザーはその恩恵を受けることができます。ただし、これらの付属リソースは、使用するデータの文脈に完全に適合してはいないので、ユーザー側の組織のデータ向けに、これらのリソース テンプレートを編集したり、独自に調整したカスタム ライブラリを作成して使用することができるようになっています。これらのリソースの形式は多岐にわたり、それぞれセッションで使用できます。リソースは、次の中にあります。

- **リソース テンプレート:**テンプレートは、製品への意見といったように、ある特定の領域や文脈に特化したリソースをまとめた形で、ライブラリ、タイプ、および拡張リソースのセットで構成されています。
- **テキスト分析パッケージ (TAP):**テンプレートに保存されているリソースに加え、リソースをもとに作成した専門的カテゴリセットをまとめたテキスト分析パッケージ(TAP)は、カテゴリとリソースをいっしょに保存して再利用することを可能にします。詳しくは、136 ページの『テキスト分析パッケージの使用』のトピックを参照してください。
- **ライブラリ:**ライブラリは、TAP およびテンプレートの構成要素として使用されます。それらは、セッションのリソースに個別に追加できます。各ライブラリはいくつかの辞書で構成され、タイプのリスト、類義語リスト、不要語リストを定義、管理するために使用されます。ライブラリは個別に提供されていますが、テンプレートおよび TAP と一緒にパッケージ化されています。詳しくは、179 ページの『第 15 章 ライブラリを使用』のトピックを参照してください。

注:抽出時、いくつかのコンパイル済み内部辞書も使用されます。これらのコンパイル済み辞書には、コア・ライブラリのタイプを補完する多くの定義が含まれています。これらのコンパイル済み辞書は編集できません。

リソース・エディター を用いることで、抽出結果 (コンセプト、タイプ、およびパターン) を出力する際に使用されるリソース セットへのアクセスが可能となります。リソース・エディター で実行するタスクには、次のような数多くのものがあります。

- **ライブラリを使用。**詳しくは、179 ページの『第 15 章 ライブラリを使用』のトピックを参照してください。
- **キーワード辞書の作成。**詳しくは、191 ページの『キーワード辞書の作成』のトピックを参照してください。
- **キーワードを辞書に追加。**詳しくは、192 ページの『キーワードの追加』のトピックを参照してください。
- **類義語の作成。**詳しくは、198 ページの『類義語の定義』のトピックを参照してください。
- **TAP のリソースの更新。**詳しくは、138 ページの『テキスト分析パッケージの更新』のトピックを参照してください。
- **テンプレートの作成。**詳しくは、164 ページの『テンプレートの作成および更新』のトピックを参照してください。

- テンプレートのインポートおよびエクスポート。詳しくは、175 ページの『テンプレートのインポートおよびエクスポート』のトピックを参照してください。
- ライブラリーの公開。詳しくは、186 ページの『ライブラリーの公開』のトピックを参照してください。

テンプレート エディターとリソース エディターの比較

テンプレート、ライブラリー、およびそれらのリソースを使用および編集する場合、主な方法は 2 つあります。テンプレート・エディター または リソース・エディター で言語リソースの作業ができます。

テンプレート・エディター

テンプレート・エディター を使用すると、インタラクティブ・ワークベンチ・セッションがなく特定のノードまたはストリームから独立しているリソース・テンプレートを作成および編集できます。このエディターを使用して、テキスト リンク分析ノードおよびテキスト マイニング モデル作成ノードにリソース テンプレートを読み込む前に、それらを作成または編集できます。

テンプレート・エディターは、IBM SPSS Modeler のメイン ツールバーで「ツール」>「Text Analytics テンプレート・エディター」メニューからアクセスできます。

リソース・エディター

インタラクティブ・ワークベンチ・セッション内で使用できる リソース・エディター によって、特定のノードおよびデータセットのコンテキストでリソースを使用できます。テキスト マイニング モデル作成ノードをストリームに追加すると、リソース・テンプレートの内容のコピー、またはテキスト分析パッケージ (カテゴリー・セットおよび リソース) のコピーを読み込んで、テキスト マイニングに使用するテキストの抽出方法を制御できます。インタラクティブ・ワークベンチ・セッションを起動すると、カテゴリーの作成、テキスト リンク分析パターンの抽出、カテゴリー・モデルの作成のほか、統合された リソース・エディター ビューでそのセッションのデータのリソースを調整することもできます。詳しくは、163 ページの『リソース・エディターを使用したリソースの編集』のトピックを参照してください。

インタラクティブ・ワークベンチ・セッションでリソースの作業を行うと、それらの変更はそのセッションにのみ適用されます。後続のセッションで継続できるように、作業 (リソース、カテゴリー、パターンなど) を保存したい場合、モデル作成ノードを更新する必要があります。詳しくは、78 ページの『モデル作成ノードの更新および保存』のトピックを参照してください。

変更を、テンプレートの内容がモデル作成ノードにコピーされている元のテンプレートに保存し直して、更新されたこのテンプレートを他のノードに読み込めるようにする場合は、リソースからテンプレートを作成できます。詳しくは、164 ページの『テンプレートの作成および更新』のトピックを参照してください。

エディターのインターフェース

テンプレート・エディターまたは リソース・エディター で実行する操作は、言語リソースの管理および調整を中心に展開しています。これらのリソースは、テンプレートおよびライブラリーの形で保存されています。詳しくは、189 ページの『キーワード辞書』のトピックを参照してください。

「ライブラリー・リソース」タブ

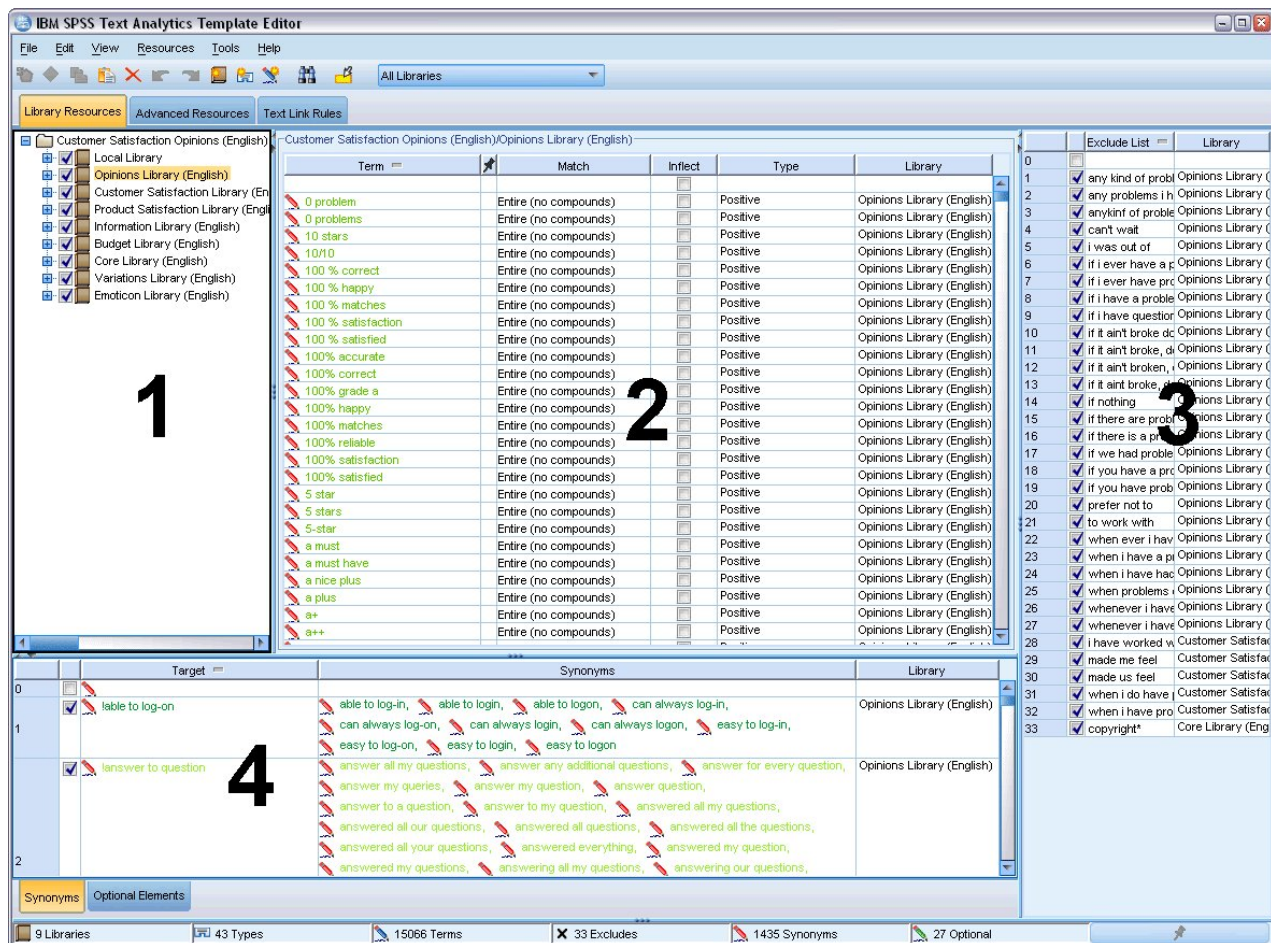


図 36. テキストマイニング・テンプレート・エディター

インターフェースは、次のような 4 つの部分で構成されています。

1. ライブラリー・ツリー・ペイン: 左上のこのパネルにはライブラリーのツリーが表示されます。このツリーでライブラリーを有効化および無効化し、ツリーのライブラリーを選択して、その他のパネルのビューをフィルタリングできます。コンテキスト・メニューを使用して、このツリーで多くの操作を実行できます。ツリーのライブラリーを展開して、含まれるタイプのセットを表示できます。特定のライブラリーのみに焦点を当てたい場合、「表示」メニューからこのリストをフィルタリングすることもできます。
2. 「キーワード辞書」ペインの用語リスト。ライブラリー・ツリーの右側にあるこのパネルには、ツリーで選択されたライブラリーのキーワード辞書のキーワード・リストが表示されます。キーワード辞書は、1つのラベル、またはタイプ、名前に基づいてグループ化されたキーワードの集合です。抽出エンジンがテキスト・データを読み取る場合、テキストの単語を、キーワード辞書のキーワードと比較します。抽出した概念がキーワード辞書でキーワードとして表示されている場合、そのタイプ名が割り当てられます。キーワード辞書を、共通点のあるキーワードの独立した辞書として見なすことができます。例えば、コア・ライブラリーの <Location> タイプには、new orleans、great britain、および new york などの概念が含まれます。これらのキーワードはすべて、地名を示します。ライブラリーには、1つまたは複数のキーワード辞書が含まれます。詳しくは、189 ページの『キーワード辞書』のトピックを参照してください。
3. 不要語辞書ペイン。右側にあるこのパネルには、最終的な抽出結果から除外されるキーワードの集合が表示されます。不要語辞書に表示されているキーワードは、抽出結果ペインには表示されません。不要語キーワードは選択するライブラリーに保存できます。ただし、「不要語辞書」パネルには、ライブラリー・ツ

リーに表示されるすべてのライブラリーの不要語登録されたすべてのキーワードが表示されます。詳しくは、200 ページの『不要語辞書』のトピックを参照してください。

4. 類義語辞書ペイン: 左下にあるこのパネルには、類義語およびオプションの要素がそれぞれのタブに表示されます。類義語およびオプションの要素を使用すると、最終的な抽出結果の代表語に基づいて類似したキーワードをグループ化できます。この辞書には既知の類義語やユーザー定義の類義語および要素、そして一般的なスペルミスと正しいスペルのペアが含まれています。類義語の定義およびオプションの要素は、選択するライブラリーに保存できます。ただし、類義語辞書ペインには、ライブラリー・ツリーに表示されるすべてのライブラリーのすべての内容が表示されます。このパネルにはすべてのライブラリーのすべての類義語またはオプションの要素が表示されますが、ツリーのすべてのライブラリーの類義語は、このパネルでいっしょに表示されます。ライブラリーには、含まれる類義語辞書は 1 つだけです。詳しくは、197 ページの『類義語辞書』のトピックを参照してください。

注:

- 1 つのライブラリーに関する情報のみ表示されるようフィルタリングしたい場合、ツールバーのドロップダウン・リストを使用して、ライブラリー・ビューを変更できます。ここには「すべてのライブラリー」という上位レベルのエントリーおよび各ライブラリーの追加エントリーが含まれます。詳しくは、182 ページの『ライブラリーの表示』のトピックを参照してください。

「詳細リソース」タブ

エディター・ビューの 2 番目のタブで詳細リソースを使用できるようになりました。このタブで詳細リソースを確認および編集することができます。詳しくは、203 ページの『第 17 章 拡張リソースについて』のトピックを参照してください。

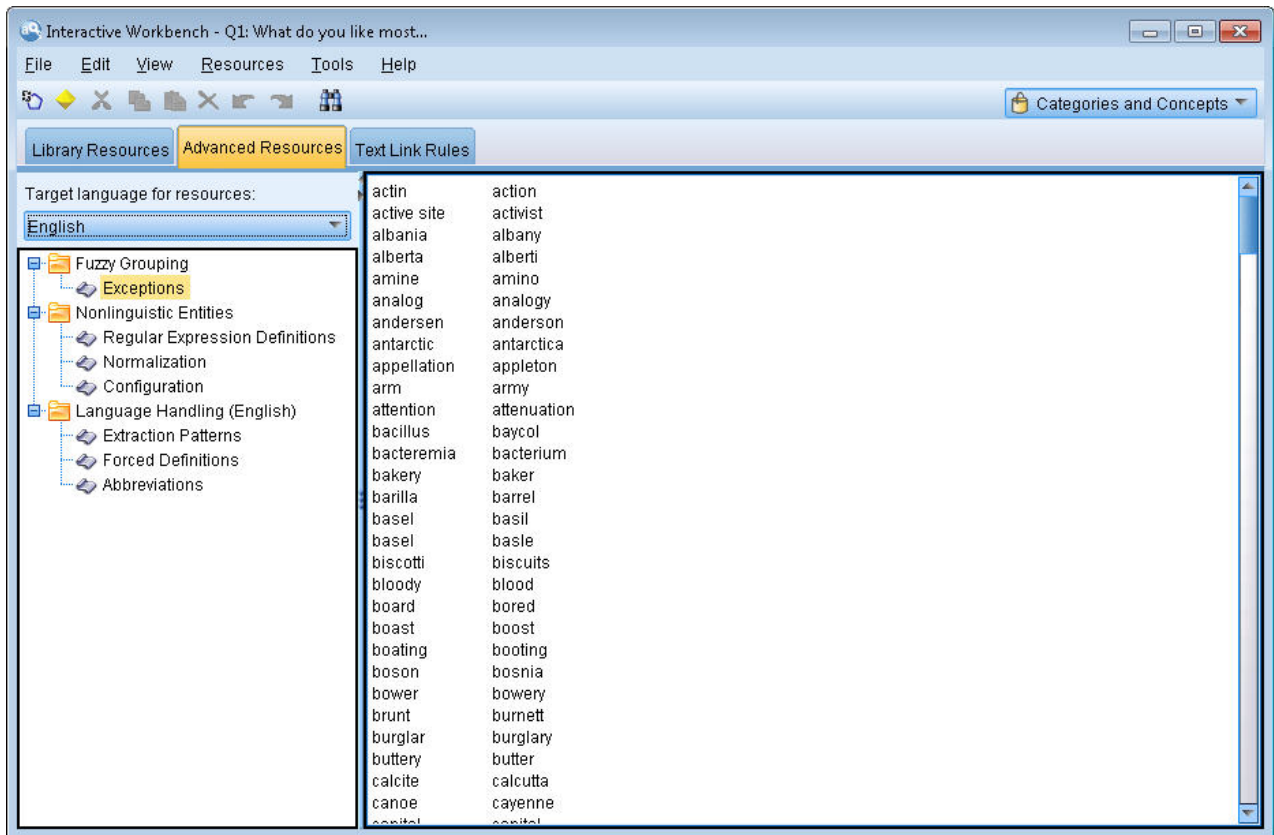


図 37. テキスト マイニング・テンプレート・エディター - 「詳細リソース」タブ

「テキスト リンク規則」タブ

バージョン 14 以降、テキスト リンク分析規則はエディター・ビューの独自のタブで編集できます。条件規則エディターで作業し、独自の規則を作成、シミュレーションを実行して、規則が TLA 結果にどのような影響を与えるかを確認できます。詳しくは、215 ページの『第 18 章 テキスト リンク規則について』のトピックを参照してください。

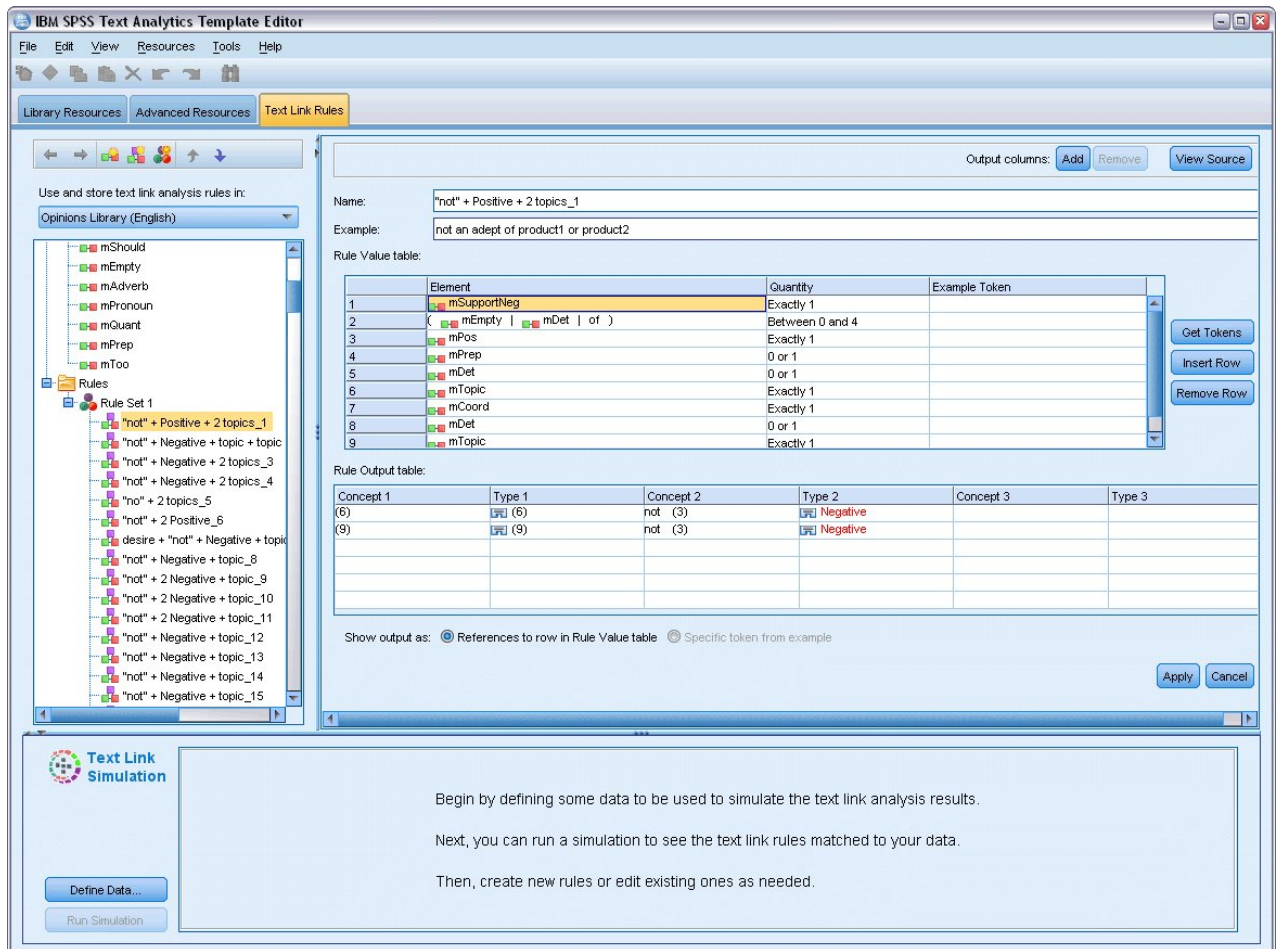


図 38. テキストマイニング・テンプレート・エディター - 「テキストリンク規則」タブ

テンプレートを開く

テンプレート・エディターを起動すると、テンプレートを開くよう要求するメッセージが表示されます。同様に、「ファイル」メニューからテンプレートを開くことができます。いくつかのテキストリンク分析 (TLA) 規則を含むテンプレートが必要な場合、「TLA」列にアイコンのあるテンプレートを選択してください。テンプレートが作成された言語は、「言語」列に表示されます。

テーブルに表示されていないテンプレートをインポートしたい場合、または、テンプレートをエクスポートしたい場合、「テンプレートを開く」ダイアログボックスのボタンを使用します。詳しくは、175 ページの『テンプレートのインポートおよびエクスポート』のトピックを参照してください。

テンプレートを開くには

1. テンプレート・エディターのメニューから「ファイル」>「リソーステンプレートを開く」を選択します。「リソーステンプレートを開く」ダイアログボックスが開きます。
2. テーブルに表示されたテンプレートから、使用したいテンプレートを選択します。
3. 「OK」をクリックすると、このテンプレートが開きます。現在エディターで別のテンプレートが開いている場合、「OK」をクリックすると開いていたテンプレートが中断し、ここで選択したテンプレートが表示されます。リソースに変更を行い、今後使用するためにライブラリーを保存したい場合、別の

テンプレートを開く前にそれらを公開、更新、共有することができます。詳しくは、184 ページの『ライブラリーの共有』のトピックを参照してください。

テンプレートの保存

テンプレート・エディター で、テンプレートに行った変更を保存できます。既存のテンプレート名を使用して保存するのか、新しい名前を付けるのかを選択できます。

以前すでにノードに読み込んだテンプレートに変更を行った場合、最新の変更を取得するには、テンプレートの内容をノードに再読み込みする必要があります。詳しくは、27 ページの『テンプレートおよび TAP からのリソースのコピー』のトピックを参照してください。

または、テキストマイニングノードの「モデル」タブのオプション「保存されたインタラクティブ作業を使用」を使用している場合、つまり以前のインタラクティブ・ワークベンチ・セッションのリソースを使用している場合、インタラクティブ・ワークベンチ・セッションからこのテンプレートのリソースに切り替える必要があります。詳しくは、165 ページの『リソース・テンプレートの切り替え』のトピックを参照してください。

注:ライブラリーを公開して共有することもできます。詳しくは、184 ページの『ライブラリーの共有』のトピックを参照してください。

テンプレートを保存するには

1. テンプレート・エディター のメニューから「ファイル」>「リソース テンプレートを保存」を選択します。「リソース・テンプレートを保存」ダイアログ・ボックスが開きます。
2. このテンプレートを新しいテンプレートとして保存する場合は、「テンプレート名」フィールドに新しい名前を入力します。既存のテンプレートをその時点で読み込まれているリソースで上書きする場合は、テーブルでテンプレートを選択します。
3. 必要に応じて、テーブルにコメントまたは注釈として表示する説明を入力します。
4. 「保存」 をクリックして、テンプレートを保存します。

重要: テンプレートまたは TAP のリソースがノードに読み込まれ/コピーされるため、テンプレートに変更を行い、既存のストリームでこれらの変更を活用する場合、それらを再読み込みすることによってリソースを更新する必要があります。詳しくは、『読み込み後のノード・リソースの更新』のトピックを参照してください。

読み込み後のノード・リソースの更新

デフォルトでは、ノードをストリームに追加する場合、デフォルト テンプレートの一連のリソースが読み込まれ、ノードに組み込まれます。テンプレートを変更または TAP を使用する場合、それらを読み込むとこれらのリソースのコピーがリソースを上書きします。テンプレートおよび TAP がノードに直接リンクしていないため、テンプレートまたは TAP に行った変更は、既存のノードでは自動的に有効になりません。これらの変更を活用するには、そのノードでリソースを更新する必要があります。リソースは、次の 2 通りの方法のいずれかで更新できます。

方法 1: 「モデル」タブでリソースを再読み込みする

新しいまたは更新されたテンプレートまたは TAP を使用してノードのリソースを更新する場合、ノードの「モデル」タブでリソースを再読み込みできます。再読み込みをして、ノードのリソースのコピーを最新の

コピーと置き換えます。元のテンプレート名のほか、更新日時が「モデル」タブに表示されます。詳しくは、27 ページの『テンプレートおよび TAP からのリソースのコピー』のトピックを参照してください。

ただし、テキストマイニングモデル作成ノードでインタラクティブ・セッション・データの作業を行っており、「モデル」タブの「セッション作業を使用 オプションを選択している場合、保存されたセッション作業およびリソースが使用され、「読み込み」ボタンが無効になります。インタラクティブ・ワークベンチ・セッション時に一度、「モデル作成ノードを更新」オプションを選択して、カテゴリ、リソースおよびその他のセッション作業を保持しているため無効になります。これらのリソースを変更または更新したい場合、次の方法でリソース・エディターのリソースを切り替える必要があります。

方法 2:リソース・エディターのリソースの切り替え

インタラクティブ・セッションで異なるリソースを使用したい場合はいつでも、「リソースを切り替え」ダイアログ・ボックスを使用してこれらのリソースを交換することができます。これは、既存のカテゴリ作業を再利用したいがリソースは置き換える場合に特に役立ちます。この場合、テキストマイニングモデル作成ノードの「モデル」タブの「セッション作業を使用」オプションを選択できます。オプションを選択すると、ノードのダイアログ・ボックスを使用してテンプレートを再読み込みする機能が無効になり、代わりにセッション時に行った設定および変更が保持されます。ストリームを実行してインタラクティブ・ワークベンチ・セッションを起動し、リソース・エディターのリソースを切り替えることができます。詳しくは、165 ページの『リソース・テンプレートの切り替え』のトピックを参照してください。

リソースなど後続のセッションにセッション作業を保持するために、リソース（およびその他のデータ）をノードに保存し直すようインタラクティブワークベンチセッション内からモデル作成ノードを更新する必要があります。詳しくは、78 ページの『モデル作成ノードの更新および保存』のトピックを参照してください。

注:インタラクティブ・セッションで別のテンプレートの内容に切り替える場合、ノードに表示されるテンプレートの名前は、最後に読み込まれ、コピーされたテンプレートの名前となります。これらのリソースまたは他のセッション作業を利用するには、セッションを終了する前にモデル作成ノードを更新します。

テンプレートの管理

テンプレート名の変更、テンプレートのインポートおよびエクスポート、または古いテンプレートの削除など、テンプレートを扱う上で時折実行する必要がある基本的な管理作業がいくつかあります。これらの作業は「テンプレートを管理」ダイアログボックスで実行されます。テンプレートをインポートおよびエクスポートすると、テンプレートを他のユーザーと共有できます。詳しくは、175 ページの『テンプレートのインポートおよびエクスポート』のトピックを参照してください。

注:この製品とともにインストールされた（付属の）テンプレートの名前を変更したり、削除することはできません。その代わりとして、名前を変更したい場合には、インストールしたテンプレートを開き、新しい名前で作成します。ユーザー定義のテンプレートは削除することができますが、付属のテンプレートを削除しようとする、このテンプレートがインストールされたときの一番最初のバージョンで置き換えられます。

テンプレートの名前を変更するには

1. メニューの「リソース」>「リソース テンプレートを管理」を選択します。「テンプレートを管理」ダイアログ・ボックスが開きます。
2. 名前を変更したいテンプレートを選択し、「名前を変更」をクリックします。表の名前のセルが編集可能なフィールドとなります。

3. 新しい名前を入力し、Enter キーを押します。確認のダイアログ・ボックスが開きます。
4. 名前を変更する場合は、「はい」 をクリックします。そうでない場合は、「いいえ」 をクリックします。

テンプレートを削除するには

1. メニューの「リソース」>「リソース テンプレートを管理」を選択します。「テンプレートを管理」ダイアログ・ボックスが開きます。
2. 「テンプレートを管理」ダイアログ・ボックスで、削除したいテンプレートを選択します。
3. 「削除」をクリックします。確認のダイアログ・ボックスが開きます。
4. 「はい」 をクリックすると削除され、「いいえ」 をクリックすると操作はキャンセルされます。「はい」 をクリックすると、テンプレートが削除されます。

テンプレートのインポートおよびエクスポート

テンプレートをインポートおよびエクスポートすることによって、テンプレートを他のユーザーまたはマシンと共有できます。テンプレートは内部データベースに格納されますが、ハード・ドライブに *.lrt ファイルとしてエクスポートできます。

テンプレートのインポートまたはエクスポートが必要な状況があるため、これらの機能を提供するダイアログ・ボックスがいくつかあります。

- テンプレート・エディター の「テンプレートを開く」ダイアログ・ボックス
- テキスト マイニング モデル作成ノードおよびテキスト リンク分析ノードの「リソースを読み込む」ダイアログ・ボックス
- テンプレート・エディター および リソース・エディター の「テンプレートを管理」ダイアログ・ボックス

テンプレートをインポートするには

1. ダイアログ・ボックスの「インポート」 をクリックします。「テンプレートをインポート」ダイアログ・ボックスが開きます。
2. インポートするリソース・テンプレート・ファイル (*.lrt) を選択して、「インポート」 をクリックします。インポートしているテンプレートを別の名前前で保存するか、既存のテンプレートを上書きできます。ダイアログ・ボックスが閉じ、表にテンプレートが表示されます。

テンプレートをエクスポートするには

1. ダイアログ・ボックスで、エクスポートしたいテンプレートを選択し、「エクスポート」 をクリックします。「ディレクトリーを選択」ダイアログ・ボックスが開きます。
2. エクスポート先のディレクトリーを選択し、「エクスポート」 をクリックします。ダイアログ・ボックスが閉じ、テンプレートがエクスポートされ、ファイルの拡張子(*.lrt) がつきます。

テンプレート・エディター の終了

テンプレート・エディター の作業を終了したら、作業を保存してエディターを終了できます。

テンプレート・エディター を終了するには

1. メニューから、「ファイル」>「閉じる」を選択します。「保存して閉じる」ダイアログ・ボックスが開きます。

2. エディターを閉じる前に開いているテンプレートを保存するには、「テンプレートへの変更を保存」を選択します。
3. エディターを閉じる前に開いているテンプレートのライブラリーを公開するには、「ライブラリーを公開」を選択します。このオプションを選択すると、公開するライブラリーを選択するよう要求するメッセージが表示されます。詳しくは、186 ページの『ライブラリーの公開』のトピックを参照してください。

リソースのバックアップ

セキュリティ上の観点から、リソースのバックアップが必要な場合があります。

重要: 復元を実行すると、リソースの内容はすべて完全に消去され、この製品ではバックアップ ファイルの内容しか使用できなくなります。これには処理中の作業も含まれます。

注: バックアップと復元ができるのは、ソフトウェアの同じメジャー・バージョンについてのみです。例えば、バージョン 15 からバックアップした場合、バージョン 16 に復元することはできません。

リソースをバックアップするには

1. メニューの「リソース」>「バックアップ ツール」>「リソースをバックアップ」を選択します。「バックアップ」ダイアログ・ボックスが開きます。
2. バックアップ・ファイルの名前を入力して、「保存」をクリックします。ダイアログ・ボックスが閉じ、バックアップ・ファイルが作成されます。

リソースを復元するには

1. メニューの「リソース」>「バックアップ ツール」>「リソースを復元」を選択します。復元するとデータベースの現在の内容が上書きされることの警告が表示されます。
2. 「はい」をクリックして、先に進みます。ダイアログ・ボックスが開きます。
3. 復元したいバックアップ・ファイルを選択し、「開く」をクリックします。ダイアログ・ボックスが閉じ、リソースがアプリケーションに復元されます。

リソース・ファイルのインポート

この製品以外のところでリソースファイルに直接変更を加えた場合、そのライブラリーを選択してインポートの手順を踏むことにより、それらのファイルを選択したライブラリーにインポートできます。ディレクトリーごとインポートする場合、対象となるファイルすべてを、特定の使用中のライブラリーにインポートすることもできます。インポートできるのは、*.txt ファイルのみです。

インポートしようとするファイルには、各行に項目1つを記載し、その内容は、以下の構造となります。

- 語または語句のリスト（各行に1つずつ）。ファイルはキーワード辞書のキーワード・リストとしてインポートされ、キーワード 辞書の名前はファイル名から拡張子を除いたものとなります。
- term1 <TAB> term2 のようなエントリーのリストとして構成されている場合、類義語のリストとしてインポートされます。term1 は基本キーワードで term2 は代表語です。

1 つのリソース・ファイルをインポートするには

1. メニューの「リソース」>「ファイルをインポート」>「単一ファイルをインポート」を選択します。「ファイルをインポート」ダイアログ・ボックスが開きます。
2. インポートしたいファイルを選択し、「インポート」をクリックします。ファイルの内容は内部形式に変換され、ライブラリーに追加されます。

ディレクトリー内のすべてのファイルをインポートするには

1. メニューの「リソース」>「ファイルをインポート」>「フォルダー全体をインポート」を選択します。「フォルダーをインポート」ダイアログ・ボックスが開きます。
2. 「インポート」 リストから、すべてのリソース・ファイルをインポートしたいライブラリーを選択します。「デフォルト」 オプションを選択すると、ディレクトリーの名前で、新しいライブラリーが作成されます。
3. ファイルのインポート元であるディレクトリーを選択します。サブディレクトリーは読み込まれません。
4. 「インポート」 をクリックします。ダイアログ・ボックスが閉じ、インポートされたリソース・ファイルの内容が辞書および拡張リソース・ファイルの形式でエディターに表示されます。

第 15 章 ライブラリーの使用

テキスト・データからキーワードを抽出してグループ化するために抽出エンジンで使用するリソースには、常に 1 つ以上のライブラリーが含まれています。テンプレート・エディター および リソース・エディターの左上部分にあるライブラリー・ツリーに一連のライブラリーが表示されます。ライブラリーは、3 種類の辞書で構成されています：キーワード、類義語、および不要語の3 種類の辞書で構成されています。詳しくは、189 ページの『第 16 章 ライブラリー辞書について』のトピックを参照してください。

リソース・テンプレートまたは選択したTAP のリソースには、複数のライブラリーが含まれており、テキスト・データからすぐにキーワード抽出を開始できるようになっています。しかしユーザーは、独自のライブラリーを作成して、それらを再利用できるよう公開することもできます。詳しくは、186 ページの『ライブラリーの公開』のトピックを参照してください。

例えば、自動車産業に関連するテキスト・データを頻繁に扱っているとします。データを分析した後、カスタム化された言語リソースを作成し、業界特有の用語や隠語を扱えるようにします。テンプレート・エディターを使用して、新しいテンプレートを作成し、テンプレート内にライブラリーを作成して自動車に関するキーワードを抽出し、グループ化します。このライブラリーの情報が再び必要になるため、ライブラリーを中央リポジトリに公開し、「ライブラリーを管理」ダイアログ・ボックスで使用できるようにします。また異なるストリーム・セッションで独立して再利用できるようになります。

また、その業界のさらに各分野(電子機器、エンジン、冷却装置、あるいは場合によっては特定の製造業者や市場)などに特有のキーワードをグループ化する必要性が出てくる場合もあります。グループごとにライブラリーを作成して公開することで、テキスト・データでこれを使用できます。こうすれば、自分のテキスト・データのコンテキストの状況に最も適切なライブラリーを追加できます。

注:追加リソースは、「拡張リソース」タブで設定および管理できます。一部のリソースはすべてのライブラリーに適用され、固有表現、Fuzzy Grouping の例外などを管理します。また、「テキスト リンク規則」タブで、ライブラリー固有のテキスト リンク分析のパターン規則を編集できます。詳しくは、203 ページの『第 17 章 拡張リソースについて』のトピックを参照してください。

付属ライブラリー

デフォルトでは、複数のライブラリーが IBM SPSS Modeler Text Analytics と共にインストールされます。これらの事前に形式設定されたライブラリーを使用して、さまざまなタイプのほか、多くの事前定義されたキーワードや類義語を使用できます。これらの付属ライブラリーは、複数の異なるドメイン向けに調整されます。また、複数の言語で使用できます。

多くのライブラリーがありますが、最も一般的に使用されているのは次のとおりです。

- ローカル・ライブラリー: ユーザー定義の辞書の格納に使用します。デフォルトではすべてのリソースに追加される空のライブラリーです。空白のキーワード辞書も含まれます。カテゴリーとコンセプト・ビュー、クラスター・ビュー、およびテキスト リンク分析ビューのリソースに直接変更または調整を行う(単語をタイプに追加するなど)場合に役立ちます。この場合、これらの変更および調整は、リソース・エディターのライブラリー・ツリーに表示された最初のライブラリーに自動的に保存されます。デフォルトでは、これがローカル・ライブラリーとなります。このライブラリーはセッション データ固有であるため、このライブラリーは公開できません。この内容を公開したい場合は、まずライブラリーの名前を変更する必要があります。

- コア・ライブラリー: 人名、地名、組織名、商品名、そして不明を示す 5 つの基本的なビルトインのタイプで構成されているため、多くのケースで使用されます。キーワード辞書の 1 つに記載されているのは少数のキーワードですが、コア・ライブラリーに記載されているタイプは、テキスト マイニング製品に付属する内部のコンパイル済み辞書の頑健なタイプを補います。これらの内部コンパイル済み辞書には、各タイプの多くのキーワードが含まれています。このため、キーワード辞書のキーワード・リストにキーワードは表示されませんが、コア タイプで抽出およびタイプ指定できます。つまり、*John*のみがコア・ライブラリーの <Person> キーワード辞書で出現する場合に、*George* をタイプ <Person> として抽出できます。同様に、コア・ライブラリーがない場合でも、それらのタイプを含むコンパイル済み辞書が抽出エンジンで使用されるため、抽出結果にそれらのタイプが表示される場合があります。
- 意見ライブラリー: テキスト・データの意見パターンを抽出する場合、最も一般的に使用されます。このライブラリーには、嗜好、識別子、優先順位を示す単語が数多く含まれています。これらは他のキーワードと連携して使用された場合に、主題についての意見を示すものです。このライブラリーには、多くのビルトインのタイプ、類義語および不要語が含まれています。また、テキスト リンク分析に使用されるパターン規則の大きなセットも含まれています。このライブラリーのテキスト リンク分析規則および生成されるパターン結果を利用するには、このライブラリーを「テキスト リンク規則」タブで指定する必要があります。詳しくは、215 ページの『第 18 章 テキスト リンク規則について』のトピックを参照してください。
- 予算ライブラリー: コストを参照するキーワードを抽出するために使用します。このライブラリーには、価格または品質に関する形容詞、識別子、意見を示す多くの単語および句が含まれています。
- バリエーション・ライブラリー: 特定の言語バリエーションが適切にグループ化するために類義語定義が必要なケースを追加するために使用します。ライブラリーには、類義語定義のみが含まれます。

テンプレート外の付属ライブラリーの一部はいくつかのテンプレートの内容に似ていますが、テンプレートは特定のアプリケーション向けに調整され、追加の拡張リソースを含んでいます。一般的なテンプレートに個別のライブラリーを追加するのではなく、処理しているテキスト・データの種類向けに作成されたテンプレートを使用し、これらのリソースに変更を行うことをお勧めします。

コンパイル済み辞書も、IBM SPSS Modeler Text Analytics に付属しています。コンパイル済み辞書は抽出プロセスで常に使用され、デフォルト・ライブラリーのビルトインのキーワード辞書に対する多くの補足的定義が含まれています。これらのリソースはコンパイルされていないため、表示も編集もできません。ただし、これらのコンパイル済み辞書にタイプが割り当てられたキーワードを他の辞書に強制投入することができます。詳しくは、195 ページの『キーワードの強制』のトピックを参照してください。

ライブラリーの作成

ライブラリーはいくつでもできます。新しいライブラリーを作成した後、このライブラリー内で辞書を作成し、キーワード、類義語、不要語を入力できます。

ライブラリーを作成するには

1. メニューの「リソース」>「新規ライブラリー」を選択します。「ライブラリーのプロパティ」ダイアログが開きます。
2. 「名前」テキスト・ボックスにライブラリーの名前を入力します。
3. 必要に応じて、「注釈」テキスト・ボックスにコメントを入力します。
4. ライブラリーに入力する前にこのライブラリーを公開したい場合は、「公開」をクリックします。詳しくは、184 ページの『ライブラリーの共有』のトピックを参照してください。後でいつでも公開することができます。

5. 「OK」 をクリックしてライブラリーを作成します。ダイアログ・ボックスが閉じ、ツリー・ビュー内にライブラリーが表示されます。ツリー内のこのライブラリーを展開すると、空白のキーワード辞書が自動的に含まれていることがわかります。すぐにキーワードを追加することもできます。詳しくは、192 ページの『キーワードの追加』のトピックを参照してください。

パブリック・ライブラリーの追加

別のセッション データからライブラリーを再利用したい場合、パブリック・ライブラリーであれば、ライブラリーを現在のリソースに追加することができます。パブリック・ライブラリーとは、公開されているライブラリーです。詳しくは、186 ページの『ライブラリーの公開』のトピックを参照してください。

パブリック・ライブラリーを追加すると、ローカル コピーがセッション データに埋め込まれます。このライブラリーに変更することはできますが、これらの変更を共有したい場合には、このライブラリーのパブリック・バージョンを再度公開する必要があります。

パブリック・ライブラリーを追加すると、このライブラリーと他のライブラリー間でキーワードのタイプが同じでない場合、「競合を解決してください」ダイアログ・ボックスが表示されます。これらの競合を自分で解決するか、あるいはそこに提示された解決方法を承認して、この操作を完了する必要があります。詳しくは、187 ページの『競合の解決』のトピックを参照してください。

注: インタラクティブ・ワークベンチ・セッションを起動する、またはセッションを閉じて公開するときにライブラリーを常に更新する場合、ライブラリーが同期しない可能性が低くなります。詳しくは、184 ページの『ライブラリーの共有』のトピックを参照してください。

ライブラリーを追加するには

1. メニューの「リソース」>「ライブラリーを追加」を選択します。「ライブラリーを追加」ダイアログ・ボックスが開きます。
2. リストのライブラリーを選択します。
3. 「追加」をクリックします。新しく追加されたライブラリーと既存のライブラリーとの間に競合箇所がある場合には、この競合を解決するか、ライブラリーを変更しないと次に進めません。詳しくは、187 ページの『競合の解決』のトピックを参照してください。

キーワードおよびタイプの検索

エディター内の「検索」機能を使用して、様々な場所を検索できます。エディターのメニューから「編集」>「検索」を選択すると、検索ツールバーが表示されます。このツールバーを使用して、一回に1つの出現ずつ検索できます。もう一度「検索」をクリックすると、次に出現しているキーワードが検索できます。

検索するとき、エディターは検索ツールバーのドロップダウン・リストに表示されている1つあるいは複数のライブラリーのみを検索します。「すべてのライブラリー」が選択されている場合、エディター内のすべてが検索されます。

検索を開始すると、対象となっている部分から検索を始めます。検索はセッションごとに行われ、開始位置であるアクティブなセルに戻るまでループします。矢印を使うことで、検索の順番を逆にできます。検索で大文字と小文字を区別するかどうかを選択することもできます。

ウィンドウ内の文字列を検索するには

1. メニューから「編集」 > 「検索」を選択します。検索ツールバーが表示されます。
2. 検索したい文字列を入力します。
3. 「検索」 ボタンをクリックして検索を開始します。該当するキーワードまたはタイプの出現が強調表示されます。
4. ボタンをもう一度クリックして次の出現しているものに移動します。

キーワードでのアスタリスクの使用

膠着型の言語 (間にスペースを入れずに単語同士を複合させて新しい単語を作る言語) を処理する場合は、キーワードでアスタリスク (*) を使用すると特に便利です。例えば、ドイツ語の単語 *Übernachtungspreis* の構造は *Übernachtung* + *s* + *Preis* です。

例えば、タイプ *Budget* でキーワード *preis** を検索すると、*preiserhöhung* などの抽出されたコンセプトに一致します。同様に、**preis* は *Übernachtung* に一致し、**preis** は *Übernachtungspreiserhöhung* に一致します。

ライブラリーの表示

ある特定のライブラリーまたはすべてのライブラリーの内容を表示できます。これは、ライブラリーがたくさんあったり、あるいは、ある特定のライブラリーの内容を公開前に確認する際に便利です。ビューを変更しても、この「ライブラリー・リソース」タブの表示内容が変わるだけで、これによって抽出のプロセスでライブラリーが使用されなくなる、というものではありません。詳しくは、183 ページの『ローカル・ライブラリーを使用不可に』のトピックを参照してください。

デフォルトのビューは「すべてのライブラリー」で、これはツリー内にすべてのライブラリーを、また他のウィンドウにその内容が表示するものです。ツールバーのドロップダウン・リストまたはメニューの選択 (「表示」 > 「ライブラリー」) によってこの選択範囲を変更できます。1 つのライブラリーが表示されている場合、他のライブラリーのすべての項目がビューから表示されなくなりますが、抽出時に読み取ることができます。

ライブラリー・ビューを変更するには

1. 「ライブラリー・リソース」タブのメニューで「表示」 > 「ライブラリー」を選択します。すべてのローカル・ライブラリーを含んだメニューが開きます。
2. 1つのライブラリーを選択するか、あるいは「すべてのライブラリー」を選択して、その内容を表示させます。ウィンドウの内容はこの選択によって変わってきます。

ローカル・ライブラリーの管理

パブリック・ライブラリーに対し、ローカル・ライブラリーはインタラクティブ・ワークベンチ・セッション内またはテンプレート内のライブラリーです。詳しくは、183 ページの『パブリック・ライブラリーの管理』のトピックを参照してください。ローカルライブラリーの基本的な管理方法としては次のようなものがあります。ローカル・ライブラリーの名前の変更、無効化、削除。

ローカル・ライブラリーの名前の変更

ローカル・ライブラリーの名前は変更できます。ローカルライブラリーの名前を変更した場合、これと同じパブリックバージョンがあった場合、それらの関係性はなくなってしまいます。つまり、それ以降の変更はこのパブリックライブラリーと共有されることはないということです。ローカルライブラリーを新しい名前

のものとして再度公開することはできます。この場合にも、このローカルライブラリーにおける変更は、元の名前のパブリックライブラリーに対して反映されません。

注:パブリック・ライブラリーの名前は変更できません。

1. メニューから「編集」>「ライブラリーのプロパティ」を選択します。「ライブラリーのプロパティ」ダイアログ・ボックスが開きます。

ローカル・ライブラリーの名前を変更するには

1. ツリー・ビュー内で、名前を変更したいライブラリーを選択します。
2. 「名前」テキスト・ボックスにライブラリーの新しい名前を入力します。
3. 「OK」をクリックし、ライブラリーの新しい名前を確定します。ダイアログボックスが閉じ、ツリー・ビュー内にあるライブラリー名が更新されます。

ローカル・ライブラリーを使用不可に

抽出プロセスからライブラリーを一時的に除外したい場合には、ツリー・ビュー内で、このライブラリーの名前の左側にあるチェックボックスをオフにします。これによって、このライブラリーはプロジェクト内に保持されますが、その内容は、競合のチェックならびに抽出のプロセスでは無視されるようになります。

ライブラリーを無効化するには

1. ライブラリー・ツリー・パネルで、使用しないライブラリーを選択します。
2. スペース・バーをクリックします。名前の左側にあるチェック・ボックスがオフになります。

ローカル・ライブラリーの削除

パブリック・バージョンのライブラリーを削除せずにライブラリーを削除することができます。その逆も可能です。ローカル・ライブラリーを削除すると、セッションのみのライブラリーおよびすべての内容が削除されます。ローカル・バージョンのライブラリーを削除しても、他のセッションまたはパブリック・バージョンのライブラリーは削除されません。詳しくは、『パブリック・ライブラリーの管理』のトピックを参照してください。

ローカル・ライブラリーを削除するには

1. ツリー・ビューで、削除したいライブラリーを選択します。
2. ライブラリーを削除するには、メニューから、「編集」>「削除」を選択します。ライブラリーは削除されます。
3. このライブラリーを公開したことがない場合には、このライブラリーを削除するか保存するかを尋ねるメッセージが表示されます。「削除」をクリックして次に進むか、「保持」をクリックし、このライブラリーを保持します。

注:1 つのライブラリーは必ず保持する必要があります。

パブリック・ライブラリーの管理

ローカル・ライブラリーを再利用するために、ローカル・ライブラリーを公開して処理し、「ライブラリーを管理」ダイアログ・ボックスに表示することができます（「リソース」>「ライブラリーを管理」）。詳しくは、184 ページの『ライブラリーの共有』のトピックを参照してください。パブリックライブラリーの基本的な管理方法としては、パブリック・ライブラリーのインポート、エクスポート、または削除があります。パブリック・ライブラリーの名前は変更できません。

パブリック・ライブラリーのインポート

1. 「ライブラリーを管理」ダイアログ・ボックスの「インポート...」をクリックします。「ライブラリーをインポート」ダイアログ・ボックスが開きます。
2. インポートしたいライブラリー (*.lib) を選択し、このライブラリーをローカルに追加したい場合、「現在のプロジェクトにライブラリーを追加」のチェックボックスをオンにします。
3. 「インポート」をクリックします。ダイアログ・ボックスがクローズします。同じ名前のパブリック・ライブラリーがすでに存在する場合、インポートしようとしているライブラリーの名前を変更するか、あるいは現在のパブリックライブラリーを上書きするのかをたずねるメッセージが表示されます。

パブリック・ライブラリーのエクスポート

パブリック・ライブラリーを .lib 形式でエクスポートすると、ライブラリーを共有できるようになります。

1. 「ライブラリーを管理」ダイアログ・ボックスで、リストからエクスポートしたいライブラリーを選択します。
2. 「エクスポート」をクリックします。「ディレクトリーを選択」ダイアログ・ボックスが開きます。
3. エクスポート先のディレクトリーを選択し、「エクスポート」をクリックします。ダイアログ・ボックスが閉じ、ライブラリー・ファイル (*.lib) がエクスポートされます。

パブリック・ライブラリーの削除

パブリック・バージョンのライブラリーを削除せずにローカル・ライブラリーを削除することができます。その逆も可能です。ただし、ライブラリーがこのダイアログ・ボックスから削除されると、ローカル・バージョンがもう一度公開されるまでセッション・リソースに追加できなくなります。

製品とともにインストールされたライブラリーを削除すると、最初にインストールされていたバージョンが復元されます。

1. 「ライブラリーを管理」ダイアログ・ボックスで、削除したいライブラリーを選択します。該当する見出しをクリックして、リストをソートすることができます。
2. 「削除」をクリックしてライブラリーを削除します。IBM SPSS Modeler Text Analytics が、ローカル・バージョンのライブラリーがパブリック・ライブラリーと同じかどうかを検証します。同じであった場合には、警告なくこのライブラリーが削除されます。しかしライブラリーのバージョンが異なっている場合、パブリックバージョンを保持するのか、あるいは削除するのかをたずねる警告が表示されません。

ライブラリーの共有

ライブラリーを使用して、複数のインタラクティブ・ワークベンチ・セッション間で共有しやすい方法でリソースを扱うことができます。ライブラリーには2つの状態、すなわち2つのバージョン（版）があります。エディターで編集可能でインタラクティブ・ワークベンチ・セッションの一部であるライブラリーはローカル・ライブラリーと呼ばれます。インタラクティブ・ワークベンチ・セッションで作業している間、例えば野菜ライブラリーに多くの変更を加えることができます。変更が他のデータでも役立つ場合、この野菜というライブラリーのパブリック・ライブラリー版を作成することで、これらのリソースを他でも使用できるようになります。パブリック・ライブラリーは、その名前のとおり、すべてのインタラクティブ・ワークベンチ・セッションの他のすべてのリソースで使用することができます。






「ライブラリーを管理」ダイアログ・ボックスにパブリック・ライブラリーが表示されます。このようなパブリックバージョンのライブラリーは、他の文脈のリソースに追加できます。こうすることで、ユーザーが作成したカスタムの言語リソースを他でも活用できます。

付属ライブラリー（インストール時に含まれるライブラリー）は、最初はパブリックライブラリーです。これらのライブラリー内のリソースを編集してから、これを新しいパブリックバージョンとすることも可能です。これらの新しいバージョンは、他のインタラクティブ・ワークベンチ・セッションでアクセスすることができます。

自分のライブラリーを使用して作業をし、これに変更を加えていった場合、このライブラリーと他のバージョンのライブラリーが同期しなくなってきました。場合によっては、ローカルバージョンのほうがパブリックバージョンより最新であったり、また反対に、パブリックバージョンのほうがローカルバージョンよりも最新であったりします。他のインタラクティブ・ワークベンチ・セッションの他のプロジェクト内ライブラリーの個々のバージョンの同期がなくなった場合、これらを再度同期させることができます。ライブラリーバージョンの同期は、ローカルライブラリーの再公開や更新によって行います。

インタラクティブ・ワークベンチ・セッションを起動または閉じるたびに、更新や再公開が必要なすべてのライブラリーを同期するためのプロンプトが表示されます。またローカルライブラリーの同期状況は、ツリー・ビュー内のライブラリー名の隣にあるアイコンや、「ライブラリーのプロパティ」ダイアログボックスを表示することによって簡単にわかります。また、同期はメニューからいつでも行うことができます。次の表は、あり得る5つの状態とそれに対応するアイコンです。

表 37. ローカル・ライブラリーの同期の状態：

アイコン	ローカル・ライブラリーの状態の説明
	未公開ローカル・ライブラリーは公開されていません。
	同期ローカル・ライブラリーのバージョンおよびパブリック・ライブラリーのバージョンが同じです。ローカル・ライブラリーにも適用されます。ローカル・ライブラリーはセッション固有のリソースのみを含むとされているため、公開できません。
	古いパブリック・ライブラリー・バージョンの方がローカル・バージョンに比べて新しいものです。ローカルバージョンを更新して、これらの変更を反映させます。
	新しいローカル・ライブラリー・バージョンの方がパブリック・バージョンに比べて新しいものです。このローカルバージョンをパブリックバージョンとして再公開できます。
	非同期ローカル・ライブラリーおよびパブリック・ライブラリーに、一方には含まれていない変更があります。この場合、ローカルライブラリーを更新するのか、あるいはこれを公開するのかを決めなくてはなりません。更新を選んだ場合、最後に更新あるいは公開した以降の変更はすべて失われます。一方、公開を選んだ場合、パブリックバージョンに加えられた変更は上書きされて失われます。

注:インタラクティブ・ワークベンチ・セッションを起動する、またはセッションを閉じて公開するときにライブラリーを更新する場合、ライブラリーが同期しない可能性が低くなります。

ライブラリーに変更するとこのライブラリーを含む他のストリームに役立つと考えた場合はいつでも、ライブラリーを再公開できます。その後、変更が他のストリームに役立つ場合は、それらのストリームのローカル・バージョンを更新することができます。このように、新しいライブラリーを作成または多くのパブリック・ライブラリーをリソースに追加することによってデータに適用される各コンテキストまたはドメインのストリームを作成できます。

あるパブリックライブラリーが共有されている場合、ローカルバージョンとパブリックバージョンの間で差異が出てくる可能性は高くなります。インタラクティブ・ワークベンチ・セッションから起動または終了して公開するたびに、またはテンプレート・エディターからテンプレートを開くまたは閉じるたびに、「ライブラリーの管理」ダイアログ・ボックスでライブラリーと同期されていないバージョンを持つすべてのライブラリーの公開や更新を行うためのメッセージが表示されます。パブリックライブラリーのバージョンがこのローカルバージョンよりも新しい場合、更新するかどうかをたずねるダイアログボックスが開きます。パブリックバージョンで更新する代わりに、現在のローカルバージョンを保持するのか、あるいは更新を現在のローカルライブラリーに反映させるのかを選択できます。

ライブラリーの公開

ライブラリーをこれまで公開していなかった場合、これを公開するとデータベース内に、ローカルライブラリーのパブリックコピーが作成されます。ライブラリーを再度公開した場合、ローカルライブラリーの内容が、既存のパブリックバージョンの内容を置き換えます。再公開した後、他のストリーム・セッションのこのライブラリーを更新し、ローカル・バージョンがパブリック・バージョンと同期するようにできます。ライブラリーを公開できる場合でも、ローカル・バージョンは常にセッションに格納されます。

重要: ローカルライブラリーが変更されており、またこれに対応するパブリックライブラリーも変更されている場合、これは同期していない（非同期）と見なされます。このような場合、まず変更されたパブリックバージョンでローカルバージョンを更新し、次にローカルバージョンを再度公開することで、両方のバージョンを全く同じにすることを推奨します。変更が加えられたローカルバージョンを先に公開してしまうと、パブリックバージョンの変更が上書きされてしまいます。

ローカル・ライブラリーをデータベースに公開するには

1. メニューの「リソース」>「ライブラリーを公開」を選択します。「ライブラリーを公開」ダイアログ・ボックスが開き、デフォルトでは、公開の必要があるライブラリーがすべて選択されています。
2. 公開または再公開したい各ライブラリーの左側にあるチェック・ボックスをオンにします。
3. 「公開」 をクリックし、このライブラリーを「ライブラリーを管理」データベースに公開します。

ライブラリーの更新

インタラクティブ・ワークベンチ・セッションを起動または閉じる場合、パブリック・バージョンと同期しないライブラリーを更新または公開できます。パブリックライブラリーのバージョンがローカルバージョンよりも新しい場合、このライブラリーを更新するかどうかをたずねるダイアログボックスが開きます。パブリックバージョンで更新せずに現在のローカルバージョンを保持するのか、あるいは現在のプロジェクト内のローカルバージョンをパブリックのもので置き換えるのかを選択できます。パブリックバージョンがローカルバージョンよりも新しい場合、ローカルバージョンを更新して、パブリックバージョンの内容と同期させることができます。更新とは、パブリックバージョン内の変更を、ローカルバージョンに適用することです。

注:インタラクティブ・ワークベンチ・セッションを起動する、またはセッションを閉じて公開するときにライブラリーを更新する場合、ライブラリーが同期しない可能性が低くなります。詳しくは、184 ページの『ライブラリーの共有』のトピックを参照してください。

ローカル・ライブラリーを更新するには

1. メニューの「リソース」>「ライブラリーを更新」を選択します。「ライブラリーを更新」ダイアログボックスが開き、デフォルトでは、更新の必要のあるライブラリーがすべて選択されています。
2. 公開または再公開したい各ライブラリーの左側にあるチェック・ボックスをオンにします。
3. 「更新」 をクリックし、ローカル・ライブラリーを更新します。

競合の解決

ローカル・ライブラリーとパブリック・ライブラリーの競合

ストリーム・セッションを起動すると、IBM SPSS Modeler Text Analytics が、ローカル・ライブラリーと「ライブラリーを管理」ダイアログ・ボックスに表示されたライブラリーとの比較を行います。セッションのローカル・ライブラリーがパブリック・バージョンと同期していない場合、「ライブラリーの同期」ダイアログ・ボックスが開きます。ここで使用したいライブラリーのバージョンを選択する際には、次のようないくつかの方法があります。

- すべてをローカル・ライブラリーで: このオプションでは、ローカルライブラリーをすべてそのまま保持します。いつでもこれらを再公開あるいは更新できます。
- すべてをこのマシン上の公開済みライブラリーで: このオプションでは、表示されたローカルライブラリーをデータベース内のバージョンで置き換えます。
- すべてを最新のライブラリーで: このオプションでは、表示されたローカルライブラリーをデータベース内のバージョンで置き換えます。
- その他: このオプションは、使用するバージョンをユーザーが表から選択できるようにします。

強制キーワードの競合

パブリックライブラリーを追加したり、あるいはローカルライブラリーを更新した場合、リソース内で、当該ライブラリーと他のライブラリーとの間で、キーワードやタイプの競合や重複が発見されることがあります。この場合、提示された競合の解決方法から選択します。あるいは競合や重複を変更するための、「強制キーワードを編集」ダイアログボックスが表示されます。詳しくは、195 ページの『キーワードの強制』のトピックを参照してください。

「強制キーワードを編集」ダイアログ・ボックスには、競合するキーワードまたはタイプのペアが表示されます。各ペアを見やすくするために、異なる背景色が使われています。これらの色は、「オプション」ダイアログ・ボックスで変更できます。詳しくは、76 ページの『オプション: 「表示」タブ』のトピックを参照してください。[強制キーワードを編集] ダイアログ・ボックスには、次の 2 つのタブがあります。

- 重複: このタブには、このライブラリー内に含まれる重複するキーワードが表示されます。キーワード部分にある押しピンのアイコンは、キーワードのこの出現が強制されていることを表わします。黒いXのアイコンがある場合、他のところで強制されているため、キーワードのこの出現は抽出時に無視されることを示します。
- ユーザー定義: このタブには、競合ではなく、キーワード辞書のキーワードパネルにおいて手作業で強制されたキーワードのリストが含まれています。

注:ライブラリーを追加・更新すると、「強制キーワードを編集」ダイアログボックスが開きます。このダイアログ・ボックスをキャンセルしても、ライブラリーの更新または追加はキャンセルされません。

競合を解決するには

1. 「強制キーワードを編集」ダイアログ・ボックスで、強制したいキーワードの「使用」列で選択します。
2. 完了したら、「OK」をクリックして強制キーワードを適用し、ダイアログ・ボックスを閉じます。「キャンセル」をクリックすると、このダイアログ・ボックスで行った変更がキャンセルされます。

第 16 章 ライブラリー辞書について

テキスト・データの抽出に使用されるこれらのリソースは、テンプレートおよびライブラリーの形式で保存されています。ライブラリーは、3 つの辞書で構成されています。

- **The** キーワード辞書には、1 つのラベル、またはタイプ名に基づいてグループ化されたキーワードの集合が含まれています。抽出エンジンがテキスト・データを読み取る場合、テキストの単語を、キーワード辞書で定義したキーワードと比較します。抽出時、タイプの用語と類義語の活用形が、コンセプトという代表語にグループ化されます。抽出されたコンセプトは、キーワードとして出現するキーワード辞書に割り当てられます。エディターの左上のウィンドウと中央のパネル (ライブラリー・ツリーとキーワードのパネル) でキーワード辞書を管理できます。詳しくは、『キーワード辞書』のトピックを参照してください。
- 類義語辞書には、最終抽出結果で、類義語、またはコンセプトという 1 つの代表語の下で類似したキーワードをグループ化するために使用する類義語またはオプションの要素として定義される単語の集合が含まれています。「類義語」タブおよび「オプション」タブを使用してエディターの左下のパネルで、類義語辞書を管理できます。詳しくは、197 ページの『類義語辞書』のトピックを参照してください。
- 不要語辞書には、最終抽出結果から削除されるキーワードおよびタイプの集合が含まれています。エディターの一番右側のパネルで不要語辞書を管理できます。詳しくは、200 ページの『不要語辞書』のトピックを参照してください。

詳しくは、179 ページの『第 15 章 ライブラリーの使用』のトピックを参照してください。

キーワード辞書

キーワード辞書は、1 つのタイプ名、またはラベル、およびキーワードのリストで構成されています。キーワード辞書は、エディターの「ライブラリー・リソース」タブの左上のパネルおよび中央のパネルで管理されます。メニューの「ビュー」>「リソース・エディター」を使用して、このビューにアクセスできます。それ以外の場合は、テンプレート・エディターで、特定のテンプレート用の辞書を編集できます。

抽出エンジンがテキスト・データを読み取る場合、テキストの単語を、キーワード辞書で定義したキーワードと比較します。キーワードは言語リソースのキーワード辞書にある語および句です。

単語がキーワードに一致した場合、そのキーワードにタイプ名が割り当てられます。抽出時にリソースが読み込まれた場合、テキスト内で見つかったキーワードはいくつかの処理手順を経て、抽出結果ペインでコンセプトとなります。同じキーワード辞書に含まれる複数のキーワードが抽出エンジンによって類義語と判断される場合、最も頻繁に出現するキーワードに基づいてグループ化され、抽出結果ペインでコンセプトとして表示されます。例えば、キーワード `question` および `query` は、最終的にコンセプト名 `question` で表示されます。

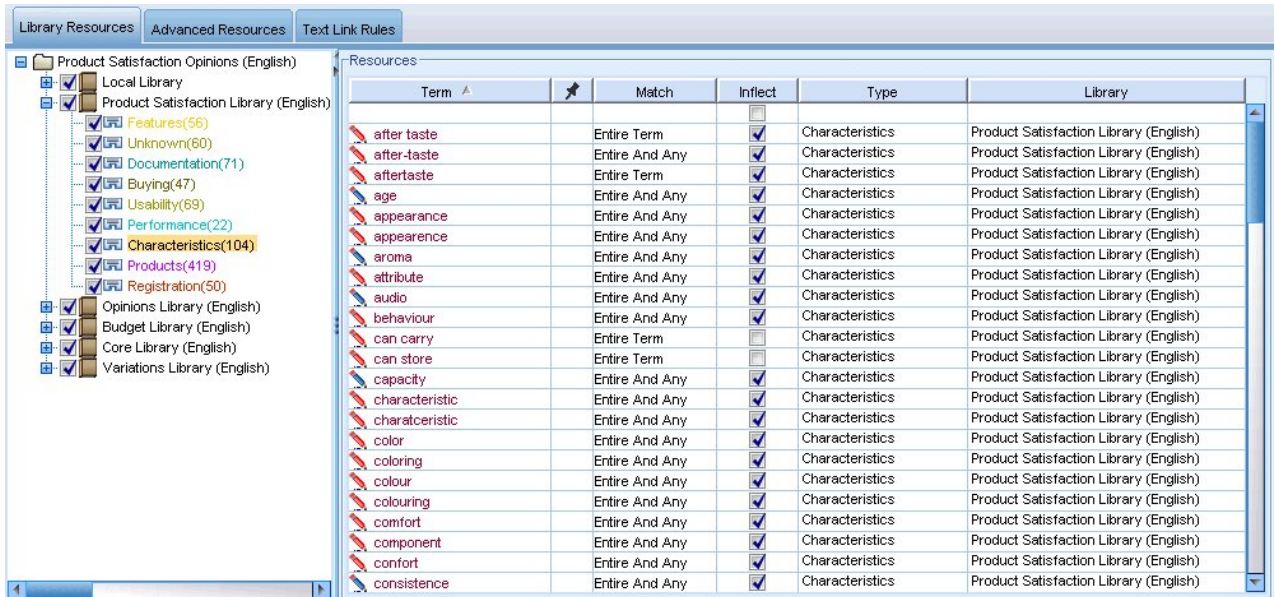


図 39. ライブラリー・ツリーおよびキーワードパネル

キーワード辞書のリストが、左側のライブラリー・ツリー・パネルに表示されます。各キーワード辞書の内容は、中央のパネルに表示されます。キーワード辞書には、キーワードのリスト以上のものが含まれています。テキスト・データの語および語句がキーワード辞書で定義されたキーワードに一致する方法は、定義されたマッチ・オプションによって決まります。マッチ・オプションは、キーワードがテキストデータの候補語または候補句に関してどのように固定されているかを指定します。詳しくは、192 ページの『キーワードの追加』のトピックを参照してください。

また、キーワードの活用形を自動的に生成して辞書に追加するかどうかを指定して、キーワード辞書のキーワードを拡張することができます。活用形を生成して、単数形の複数形、複数形の単数形、および形容詞をキーワードに自動的に追加することができます。詳しくは、192 ページの『キーワードの追加』のトピックを参照してください。

注: 多くの言語の場合、キーワード辞書にはないが、テキストから抽出されたコンセプトは、自動的に <Unknown> のタイプとなります。

キーワードでのアスタリスクの使用

膠着型の言語 (間にスペースを入れずに単語同士を複合させて新しい単語を作る言語) を処理する場合は、キーワードでアスタリスク (*) を使用すると特に便利です。例えば、ドイツ語の単語 *Übernachtungspreis* の構造は *Übernachtung* + *s* + *Preis* です。

例えば、タイプ Budget でキーワード *preis** を検索すると、*preiserhöhung* などの抽出されたコンセプトに一致します。同様に、**preis* は *Übernachtung* に一致し、**preis** は *Übernachtungspreiserhöhung* に一致します。

ビルトインのタイプ

IBM SPSS Modeler Text Analytics は付属ライブラリーおよびコンパイル済み辞書の形式で一連の言語リソースに付属しています。付属ライブラリーには、<地名>、<組織名>、<人名>、および <商品名> を含むビルトインのキーワード辞書が含まれています。

これらのキーワード辞書は、抽出エンジンによって使用され、タイプ <Location> をコンセプト paris に割り当てるように、タイプを抽出したコンセプトに割り当てます。多くのキーワードがビルトインのキーワード辞書で定義されていますが、すべての可能性をカバーしているわけではありません。そのため、辞書に追加するか独自に作成することができます。特定の付属キーワード辞書の内容の詳細は、「タイプのプロパティ」ダイアログ・ボックスの注釈を参照してください。ツリーでタイプを選択し、コンテキスト・メニューから「編集」>「プロパティ」をクリックします。

注:

付属ライブラリーのほか、(抽出エンジンにも使用される) コンパイル済み辞書には、ビルトインのタイプ辞書に対する定義の多くの補足が含まれていますが、その内容は製品では表示されません。ただし、これらのコンパイル済み辞書にタイプが割り当てられたキーワードを他の辞書に強制投入することができます。詳しくは、195 ページの『キーワードの強制』を参照してください。

キーワード辞書の作成

キーワード辞書を作成して、類似したキーワードをグループ化できます。この辞書に出現するキーワードが抽出プロセスで見つかった場合、キーワードはそのタイプ名に割り当てられ、コンセプト名で抽出されます。ライブラリーを作成すると、すぐにキーワードを入力できるような空白のキーワード辞書が含まれます。

食べ物に関するテキストを分析しており、野菜に関連するキーワードをグループ化する場合、独自の <Vegetables> キーワード辞書を作成できます。テキストに出現する重要なキーワードと感じた場合は、carrot、broccoli、および spinach などのキーワードを追加できます。抽出時、これらのキーワードのいずれかが見つかった場合、コンセプトとして抽出され、<Vegetables> タイプに割り当てられます。

キーワードの活用形を生成できるため、単語または表現のすべての形式を定義する必要はありません。このオプションを選択すると、抽出エンジンは、他の形式からキーワードの単数形または複数形を、このタイプに割り当てられているものとして自動的に認識します。このオプションは、動詞または形容詞の活用形が必要な場合が少ないため、タイプには主に名詞が含まれている場合に役立ちます。

「タイプのプロパティ」ダイアログ・ボックスには以下のフィールドがあります。

名前: 作成しているキーワード辞書に指定する名前。複数のタイプ名が同じ単語で始まる場合は特に、タイプ名にスペースを使用しないことをお勧めします。

注: タイプ名や記号の使用にはいくつかの制限があります。例えば、"@ " や "! " は使用できないなどです。
within the name.

デフォルトの合致: デフォルトの合致属性は、抽出エンジンにこのキーワードがテキスト・データにどのように合致するかを指示します。キーワードをキーワード辞書に追加すると、これが自動的キーワードに割り当てられる合致属性となります。キーワード・リストで合致の選択を手動でいつでも選択することができます。オプションの内容: オプションには、「キーワード全体」、「語頭」、「語末」、「任意」、「語頭あるいは語末」、「完全かつ語頭」、「完全かつ語末」、「完全かつ (語頭あるいは語末)」、および「全体 (複合語なし)」があります。詳しくは、192 ページの『キーワードの追加』のトピックを参照してください。

追加先: 新しいキーワード辞書を作成するライブラリーを指定します。

デフォルトで活用形を生成: 抽出エンジンに、文法的形態論を使用して、キーワードの単数形または複数形など、この辞書に追加するキーワードの類似した形式をキャプチャおよびグループ化するように指示します。

このオプションは、タイプにほとんど名詞が含まれている場合に特に役立ちます。このオプションを選択すると、このタイプに追加されたすべての新しいキーワードには自動的にこのオプションが与えられますが、リストで手動で変更できます。

フォントの色: このフィールドを指定すると、インターフェースでこのタイプの結果と他のタイプの結果とを区別できるようになります。「グローバル設定の色を使用」を選択すると、タイプのデフォルト色がこのキーワードに使用されます。このデフォルト色は、「オプション」ダイアログ・ボックスで設定されます。詳しくは、76 ページの『オプション: 「表示」タブ』のトピックを参照してください。「カスタム」を選択すると、ドロップダウン・リストから色を選択できます。

注釈: このフィールドはオプションで、任意のコメントまたは説明に使用できます。

キーワード辞書を作成するには

1. 新しいキーワード辞書を作成するライブラリーを選択します。
2. メニューの「ツール」>「新規タイプ」を選択します。「タイプのプロパティ」ダイアログ・ボックスが開きます。
3. 「名前」テキスト・ボックスにキーワード辞書の名前を入力し、必要なオプションを選択します。
4. 「OK」をクリックしてキーワード辞書を作成します。新しいタイプがライブラリー・ツリー・パネルに表示され、中央パネルに表示されます。すぐにキーワードを追加できます。詳しくは、『キーワードの追加』を参照してください。

注: ここでは、リソース・エディター・ビューまたは テンプレート・エディターでどのように変更を行うかについて説明します。抽出結果ペイン、データ・ペイン、カテゴリー・ペイン、または他のビューの「クラスター定義」ダイアログ・ボックスで直接、この種類の調整を行うこともできます。詳しくは、89 ページの『抽出結果の調整』のトピックを参照してください。

キーワードの追加

ライブラリー・ツリー・パネルにはライブラリーが表示され、ツリーを展開すると、ツリーに含まれているキーワード辞書が表示されます。ツリーの選択によって、中央のパネルのキーワードリストに選択したライブラリーまたはキーワード辞書のキーワードが表示されます。

リソース・エディター では、用語ペインで直接または「新しいキーワードを追加」ダイアログ・ボックスを使用して、キーワードをキーワード辞書に追加できます。追加するキーワードは、単語でも複合語でもかまいません。リストの一番上に空白の行があり、そこに新しいキーワードを追加できます。

注: ここでは、リソース・エディター・ビューまたは テンプレート・エディターでどのように変更を行うかについて説明します。抽出結果ペイン、データ・ペイン、カテゴリー・ペイン、または他のビューの「クラスター定義」ダイアログ・ボックスで直接、この種類の調整を行うこともできます。詳しくは、89 ページの『抽出結果の調整』のトピックを参照してください。

「キーワード」列

この列のセルに、単語または複合語を入力します。キーワードが表示される色は、キーワードが保存または強制投入されるタイプの色によって異なります。「タイプのプロパティ」ダイアログ・ボックスでタイプの色を変更できます。詳しくは、191 ページの『キーワード辞書の作成』のトピックを参照してください。

「強制」列

このセルで押しピンのアイコンを押すと、抽出エンジンは、他のライブラリーのこの同じキーワードの他の出現を無視します。詳しくは、195 ページの『キーワードの強制』のトピックを参照してください。

「一致」列

この列ではマッチ・オプションを選択して、抽出エンジンにこのキーワードがテキスト・データにどのように合致するかを指示します。例については表を参照してください。タイプのプロパティを編集して、デフォルト値を変更できます。詳しくは、191 ページの『キーワード辞書の作成』のトピックを参照してください。メニューから「編集」>「合致方法を変更」を選択します。次に示すのは基本的なマッチ・オプションで、これらを組み合わせて使用することもできます。

- 「語頭」。辞書のキーワードがテキストから抽出したコンセプトの最初の文字に合致する場合、このタイプが割り当てられます。例えば、apple と入力すると、apple tart が合致します。
- 「語末」。辞書のキーワードがテキストから抽出したコンセプトの最後の文字に合致する場合、このタイプが割り当てられます。例えば、apple と入力すると、cider apple が合致します。
- 「任意」。辞書のキーワードがテキストから抽出したコンセプトのいずれかの単語に合致する場合、このタイプが割り当てられます。例えば、apple と入力すると、「任意」のタイプが、apple tart、cider apple、および cider apple tart に割り当てられます。
- キーワード全体:テキストから抽出したコンセプト全体が辞書キーワードにそのまま合致する場合、このタイプが割り当てられます。キーワードを「キーワード全体」として追加すると、「完全かつ語頭」、「完全かつ語末」、「完全かつどこか」、または「全体 (複合語なし)」がキーワードの抽出を強制します。

さらに、<Person> タイプは、*edith piaf* や *mohandas gandhi* のように名前の 2 つの部分抽出するため、姓について記載されていないときに名を抽出する場合、名をこのキーワード辞書に明示的に追加したい場合があります。例えば、*edith* のすべてのインスタンスを名前として取得したい場合、「キーワード全体」または「完全かつ語頭」を使用して、*edith* を <Person> タイプに追加する必要があります。

- 「全体 (複合語なし)」:テキストから抽出したコンセプト全体が辞書キーワードにそのまま合致する場合、このタイプが割り当てられ、キーワードが長い複合語に合致しないよう、抽出が停止します。例えば、apple と入力すると、「全体 (複合語なし)」オプションは、apple にタイプを割り当て、別に強制されていない限り複合語 apple sauce は抽出されません。

次の表では、キーワード *apple* がキーワード辞書にあると想定します。マッチ・オプションによって、この表にはテキスト内で見つかった場合に抽出およびタイプが割り当てられるコンセプトが表示されます。

表 38. 合致の例



マッチ・オプション - キーワード:  apple	抽出コンセプト			
	apple	apple tart	ripe apple	homemade apple tart
「キーワード全体」	✓			

表 38. 合致の例 (続き)

マッチ・オプション - キーワード:  apple	抽出コンセプト			
	apple	apple tart	<i>ripe apple</i>	<i>homemade apple tart</i>
「語頭」		✓		
「語末」			✓	
「語頭あるいは語末」		✓	✓	
「完全かつ語頭」	✓	✓		
「完全かつ語末」	✓		✓	
「完全かつ (語頭あるいは語末)」	✓	✓	✓	
「任意」		✓	✓	✓
「完全かつどこか」	✓	✓	✓	✓
「全体 (複合語なし)」	✓	抽出なし	抽出なし	抽出なし

「活用」列

この列では、抽出エンジンが抽出時にこのキーワードの活用形を生成し、すべてをグループ化するかどうかを選択します。この列のデフォルト値は「タイプのプロパティ」で定義されますが、場合によっては、列で直接このオプションを変更できます。メニューから「編集」>「活用形を変更」を選択します。

「タイプ」列

この列で、ドロップダウン・リストからキーワード辞書を選択します。タイプのリストは、ライブラリー・ツリー・パネルでの選択に応じてフィルタリングされます。リストの最初のタイプは、ライブラリー・ツリー・パネルで選択されたデフォルト・タイプです。メニューから「編集」>「タイプを変更」を選択します。

「ライブラリー」列

キーワードが格納されているライブラリーが表示されます。ライブラリー・ツリー・パネルでキーワードを別のタイプにドラッグ・アンド・ドロップして、そのライブラリーを変更できます。

キーワード辞書に 1 つのキーワードを追加するには

1. ライブラリー・ツリー・パネルで、キーワードを追加したいキーワード辞書を選択します。
2. 中央のパネルのキーワード・リストで、使用できる最初の空白セルにキーワードを入力し、このキーワードに必要なオプションを設定します。

キーワード辞書に複数のキーワードを追加するには

1. ライブラリー・ツリー・パネルで、キーワードを追加したいキーワード辞書を選択します。
2. メニューの「ツール」>「新規キーワード」を選択します。「新しいキーワードを追加」ダイアログ・ボックスが開きます。
3. キーワードを入力するか、キーワードのセットをコピーして貼り付けて、選択したキーワード辞書に追加したいキーワードを入力します。複数のキーワードを入力する場合、「オプション」ダイアログで定義された区切り文字を使用してキーワードを区切るか、各キーワードを新しい行に追加する必要があります。詳しくは、76 ページの『オプションの設定』のトピックを参照してください。
4. 「OK」をクリックすると、キーワードが辞書に追加されます。マッチ・オプションは、このキーワード・ライブラリーのデフォルトのオプションに自動的に設定されます。ダイアログ・ボックスが閉じ、辞書に新しいキーワードが表示されます。

キーワードの強制

キーワードを特定のタイプに割り当てる場合、対応するキーワード辞書に追加することができます。ただし、同じ名前を持つ複数のキーワードがある場合、抽出エンジンはどのタイプを使用するかを認識する必要があります。そのため、使用するタイプを選択するよう要求するメッセージが表示されます。この操作をタイプへのキーワードの強制といいます。このオプションは、コンパイル済み (内部、編集不可能) 辞書からのタイプの割り当てを上書きする場合に役立ちます。通常、キーワードが重複しないようにすることをお勧めします。

強制しても、このキーワードの他の出現を「削除」するわけではありません。抽出エンジンによって無視されます。キーワードを強制または強制を解除することによって、使用する出現を後で変更することができます。また、パブリック・ライブラリーを追加またはパブリック・ライブラリーを更新する場合、キーワードをキーワード辞書に投入することが必要な場合もあります。

キーワード・パネルの 2 番目の列、「強制」列で、度のキーワードが強制されているか、または無視されているかを確認できます。押しピンのアイコンが表示されている場合、キーワードのこの出現が強制されていることを示します。黒い X のアイコンがキーワードの後に表示されている場合、他の場所で強制されているため、キーワードのこの出現は抽出時に無視されることを示します。また、キーワードを強制すると、強制されたタイプの色で表示されます。タイプ 1 および タイプ 2 のキーワードが タイプ 1 に強制投入されると、タイプ 1 に定義されたフォントの色で表示されます。

アイコンをダブルクリックして、状態を変更できます。キーワードが他の場所で表示されている場合、「不一致を解決してください」ダイアログ・ボックスが開き、使用する出現を選択できます。

キーワード辞書の名前変更

「タイプのプロパティ」を編集して、キーワード辞書の名前を変更したり、その他の辞書の設定を変更することができます。

重要: 複数のタイプ名が同じ単語で始まる場合は特に、タイプ名にスペースを使用しないことをお勧めします。コア・ライブラリーまたは意見ライブラリーでタイプの名前を変更したり、デフォルトの合致属性を変更しないことをお勧めします。

キーワード辞書の名前を変更するには

1. ライブラリー・ツリー・パネルで、名前を変更したいキーワード辞書を選択します。
2. マウスを右クリックし、コンテキスト・メニューから「タイプのプロパティ」をクリックします。「タイプのプロパティ」ダイアログ・ボックスが開きます。
3. 「名前」テキスト・ボックスにキーワード辞書の新しい名前を入力します。
4. 「OK」をクリックして、新しい名前を承認します。新しいタイプ名がライブラリー・ツリー・パネルに表示されます。

キーワード辞書の移動

キーワード辞書をライブラリー内の別の場所、またはツリー内の別のライブラリーにドラッグすることができます。

ライブラリー内のキーワード辞書の順序を変更するには

1. ライブラリー・ツリー・パネルで、移動したいキーワード辞書を選択します。
2. メニューで「編集」>「1つ上に移動」を選択すると、ライブラリー・ツリー・パネルでキーワード辞書が1つ上の位置に移動します。「編集」>「1つ下に移動」を選択すると1つ下の位置に移動します。

キーワード辞書を別のライブラリーに移動するには

1. ライブラリー・ツリー・パネルで、移動したいキーワード辞書を選択します。
2. マウスを右クリックし、コンテキスト・メニューから「タイプのプロパティ」をクリックします。「タイプのプロパティ」ダイアログ・ボックスが開きます。(タイプを別のライブラリーにドラッグ・アンド・ドロップすることができます)。
3. 「追加先」リスト・ボックスで、キーワード辞書を移動したいライブラリーを選択します。
4. 「OK」をクリックします。ダイアログ・ボックスが閉じ、タイプが選択したライブラリーに移動します。

キーワード辞書の無効化および削除

キーワード辞書を一時的に削除したい場合、ライブラリー・ツリー・パネルの辞書名の左側にあるチェック・ボックスをオフにして無効化することができます。これは、ライブラリー内に辞書を保存したいが、競合を検証する場合および抽出プロセス時に内容を無視することを示します。

ライブラリーからキーワード辞書を永続的に削除することもできます。

キーワード辞書を無効化するには

1. ライブラリー・ツリー・パネルで、無効化したいキーワード辞書を選択します。
2. スペース・バーをクリックします。タイプ名前の左側にあるチェック・ボックスがオフになります。

キーワード辞書を削除するには

1. ライブラリー・ツリー・パネルで、削除したいキーワード辞書を選択します。
2. メニューから、「編集」>「削除」を選択すると、キーワード辞書が削除されます。

類義語辞書

類義語辞書は、1 つの代表語に基づいて、類似したキーワードをグループ化できるキーワードの集合です。類義語辞書は、「ライブラリー・リソース」タブの一番下のパネルで管理されます。メニューの「ビュー」>「リソース・エディター」を使用して、このビューにアクセスできます。それ以外の場合は、テンプレート・エディターで、特定のテンプレート用の辞書を編集できます。

当辞書では類義語の定義の方法が2つあります。類義語およびオプションの要素。このパネルでタブをクリックして切り替えることができます。

テキスト・データの抽出を実行した後、他のコンセプトの類義語または活用形であるコンセプトをいくつか見つけることができます。オプションの要素および類義語を特定して、抽出エンジンがこれらを1つの代表語にマップされるよう強制できます。

類義語やオプションの要素を使用すると、より頻度の高いドキュメント数のより重要で代表的なコンセプトに結合することによって、抽出結果ペインのコンセプト数を削減することができます。

類義語

類義語とは、同じ意味を持つ複数の語を関連付けたものです。また、類義語はキーワードを略語とグループ化したり、一般的にスペルミスのある単語と正しいスペルの単語とをグループ化したりするために使用できます。「類義語」タブでこれらの類義語を定義できます。

類義語定義は、2つの部分で構成されています。1つ目は代表語で、抽出エンジンがすべての類義語キーワードをグループ化する基準となるキーワードです。この代表語が別の代表語の類義語として使用されていない限り、または不要語として登録されていない限り、抽出結果ペインに表示されるコンセプトとなります。2つめは、代表語の下にグループ化される類義語のリストです。

例えば、`automobile` という単語を `vehicle` という単語に置き換えたい場合、`automobile` が類義語となり、`vehicle` が代表コンセプトとなります。

「類義語」列にはどんな言葉も入力できますが、その語が抽出時に見つからず、キーワードのマッチ・オプションが「キーワード全体」である場合、類義語は出現しません。ただし、類義語をこのキーワードの下にグループ化するために、代表語を抽出する必要がありません。

オプションの要素

オプションの要素は、テキスト内で若干異なったように出現する場合でも類義語を保持するために抽出時に無視できる結合キーワード内の任意の語を示します。オプションの要素は単語で、結合語から削除された場合、別のキーワードと合致を作成できます。これらの単語は、結合語内のどこでも（語頭、語中、語末）出現します。「オプション」タブでオプションの要素を定義できます。

例えば、`ibm` および `ibm corp` のキーワードをグループ化する場合、`corp` をオプションの要素として扱うよう宣言する必要があります。また別の例を示すと、キーワード `access` をオプションの要素として指定し、抽出時に `internet access speed` および `internet speed` が見つかった場合、最も頻繁に出現するキーワードに基づいてグループ化されます。

類義語の定義

「類義語」タブで、テーブルの一番上の空白行に類義語定義を入力できます。まず代表語とその類義語を定義します。この定義を格納するライブラリーも選択できます。抽出時、類義語のすべての出現を、最終的な抽出で代表語に基づいてグループ化します。詳しくは、192 ページの『キーワードの追加』のトピックを参照してください。

例えば、テキストデータにて、多くの電気通信情報が含まれている場合には、下記の単語が出てくる可能性があります：セルラーフォン、ワイアレスフォンおよび モバイルフォン。この場合、cellular および mobile を wireless の類義語として定義する必要があります。これらの類義語を定義すると、cellular phone および mobile phone の抽出されたすべての出現は、wireless phone と同じキーワードとして扱われ、キーワード・リストにいっしょに表示されます。

キーワード辞書を作成している場合、キーワードを入力し、そのキーワードに対して 3 つまたは 4 つの類義語を考えることができます。この場合、類義語辞書にすべてのキーワードと代表語を入力し、類義語をドラッグすることができます。

類義語は、類義語の活用形 (複数形など) にも適用されます。コンテキストに応じて、キーワードをどのように代用するかについて制限が必要な場合があります。特定の文字を使用して、類義語の処理の程度に制限を加えることができます。

- 感嘆符 (!)。!synonym のように類義語の前に直接感嘆符を追加すると、類義語の活用形は代表語によって代用されないことを示します。ただし、!target-term のように代表語の前に感嘆符を追加すると、複合代表語の一部または変異形がさらなる類義語を受け入れないようにすることを示します。
- アスタリスク (*)。synonym* のようにアスタリスクを類義語の直後に置くと、この語を代表語に置き換えることを示します。例えば、manage* を類義語に、management を代表語に定義すると、associate managers が代表語 associate management に置き換えられます。また、internet * のようにごとアスタリスクの間にスペースを追加することもできます (synonym *)。代表語に internet、類義語に internet * * および web * を定義すると、internet access card および web portal は internet に置き換えられます。この辞書では、語または文字列の頭をアスタリスクにすることはできません。
- カレット (^)。^ synonym のように類義語の前にカレットとスペースを追加すると、キーワードが類義語から始まる場合にのみ、類義語のグループ化が適用されることを示します。例えば、^ wage を類義語に、income を代表語に定義して両方が抽出される場合、キーワード income に基づいてグループ化されます。ただし、minimum wage と income が抽出されると、minimum wage が wage で始まっていないため、それらはいっしょにグループ化されません。この記号と類義語の間にスペースを追加する必要があります。
- ドル記号 (\$)。synonym \$ のように類義語の後にスペースドル記号を追加すると、キーワードが類義語で終わる場合にのみ、類義語のグループ化が適用されることを示します。例えば、cash \$ を類義語に、money を代表語に定義して両方が抽出される場合、キーワード money に基づいてグループ化されます。ただし、cash cow と money が抽出されると、cash cow が cash で終わっていないため、それらはいっしょにグループ化されません。この記号と類義語の間にスペースを追加する必要があります。
- カレット (^) およびドル記号 (\$)。^ synonym \$ のように同時に使用すると、完全一致の場合にのみキーワードが類義語と合致します。つまり、類義語のグループ化が行われるよう、抽出されたキーワードの類義語の前後に語があってはいけないことを意味します。例えば、^ van \$ を類義語に、truck を代表語に定義すると、van は truck とグループ化されますが、marie van guerin がそのままになります。また、カレットとドル記号を使用して類義語を定義して、この語がソース・テキスト内のどこにも出現する場合、類義語は自動的に抽出されます。

類義語エントリーを追加するには

1. 類義語パネルを表示し、左下の「類義語」タブをクリックします。
2. テーブルの一番上の空白行で、「代表語」列に代表語を入力します。入力した代表語が色つきで表示されます。この色は、キーワードが表示されるまたは強制されるタイプを示します。キーワードが黒で表示されている場合、キーワード辞書にはこのキーワードがないことを意味します。
3. 代表語の右側の 2 番目のセルをクリックして、類義語のセットを入力します。「オプション」ダイアログ・ボックスで定義したグローバル区切り文字を使用して、各エントリーを区切ります。詳しくは、76 ページの『オプションの設定』のトピックを参照してください。入力したキーワードが色つきで表示されます。この色は、キーワードが出現するタイプを示します。キーワードが黒で表示されている場合、キーワード辞書にはこのキーワードがないことを意味します。
4. 最後のセルをクリックして、この類義語定義を格納するライブラリーを選択します。

注: ここでは、リソース・エディター・ビューまたは テンプレート・エディターでどのように変更を行うかについて説明します。抽出結果ペイン、データ・ペイン、カテゴリ・ペイン、または他のビューの「クラスター定義」ダイアログ・ボックスで直接、この種類の調整を行うこともできます。詳しくは、89 ページの『抽出結果の調整』のトピックを参照してください。

オプションの要素の定義

「オプション」タブで、必要なライブラリーのオプションの要素を定義します。これらのエントリーは、各ライブラリーでグループ化されます。ライブラリーがライブラリー・ツリー・パネルに追加されると、空白のオプションの要素行が「オプション」タブに追加されます。

すべてのエントリーが、自動的に小文字の語に変換されます。抽出エンジンは、エントリーをテキスト内の小文字および大文字の語に合致させます。

注: 「オプション」ダイアログ・ボックスで定義した区切り文字を使用して、キーワードを区切ります。詳しくは、76 ページの『オプションの設定』のトピックを参照してください。入力しているオプションの要素にキーワードの一部と同じ区切り文字が含まれている場合、その前にバックスラッシュを追加する必要があります。

エントリーを追加するには

1. 類義語パネルを表示し、エディターの左下にある「オプション」タブをクリックします。
2. このエントリーを追加するライブラリーの「オプションの要素」列のセルをクリックします。
3. オプションの要素を入力します。「オプション」ダイアログ・ボックスで定義したグローバル区切り文字を使用して、各エントリーを区切ります。詳しくは、76 ページの『オプションの設定』のトピックを参照してください。

類義語の無効化および削除

辞書でエントリーを無効化して、一時的に削除することができます。エントリーを無効化すると、そのエントリーは抽出時に無視されます。

類義語辞書の古いエントリーを削除することもできます。

エントリーを無効化するには

1. 辞書で、無効化したいエントリーを選択します。
2. スペース・バーをクリックします。エントリーの左側にあるチェック・ボックスがオフになります。

注:エントリーの左側にあるチェック・ボックスをオフにして、無効化することもできます。

類義語エントリーを削除するには

1. 辞書で、削除したいエントリーを選択します。
2. メニューから、「編集」>「削除」を選択または **Del** キーを押すと、キーワード辞書が削除されます。エントリーは辞書内から除外されます。

オプションの要素エントリーを削除するには

1. 辞書で、削除したいエントリーをダブルクリックします。
2. キーワードを手動で削除します。
3. **Enter** を押して変更を適用します。

不要語辞書

不要語辞書は、語、句、一部の文字列のリストです。不要語辞書のエントリーに合致またはエントリーを含むキーワードは無視されるか、抽出から除外されます。不要語辞書は、エディターの右側のパネルで管理されます。通常、このリストに追加するコンセプトは、テキストの中で穴埋めとして使われる単語または句で、特に重要な意味を付け加えるようなものではなく、抽出結果を混乱させる場合があります。コンセプトを不要語辞書に追加しておけば、それらが抽出されないようにすることができます。

不要語辞書は、エディターの「ライブラリー・リソース」タブの右上のパネルで管理されます。メニューの「ビュー」>「リソース・エディター」を使用して、このビューにアクセスできます。それ以外の場合は、テンプレート・エディターで、特定のテンプレート用の辞書を編集できます。

不要語辞書で、テーブルの一番上の空白行に語、句、または一部の文字列を入力できます。文字列を不要語辞書に 1 つまたは複数の語として、またはアスタリスクをワイルドカードとして使用し、単語の一部として追加することができます。不要語辞書で宣言されたエントリーを使用して、コンセプトを抽出から除外します。エントリーが、キーワード辞書などインターフェースの別の場所でも宣言されている場合、他の辞書では、現在除外されていることを示す取り消し線が表示されます。この文字列は、テキスト・データに出現する必要はなく、また適用されるキーワード辞書の一部として宣言する必要はありません。

注: 類義語エントリーで代表語としても機能するコンセプトを不要語辞書に追加すると、代表語およびそのすべての類義語も不要語として登録されます。詳しくは、198 ページの『類義語の定義』のトピックを参照してください。

ワイルドカード (*) の使用

アスタリスクのワイルドカードを使用して、不要語のエントリーを文字列の一部として扱うことを示します。抽出エンジンで見つかった、不要語辞書で入力した文字列で始まるまたは終わる語を含むキーワードは、最終的な抽出からは除外されます。ただし、ワイルドカードの使用が認められない場合が 2 つあります。

- *- のようにアスタリスクのワイルドカードの後にダッシュ (-) 文字が追加されている場合
- *'s のようにアスタリスクのワイルドカードの後にアポストロフィ (') が追加されている場合

表 39. 不要語エントリーの例

投入	例	結果
語	<i>next</i>	<i>next</i> という語が含まれている場合、コンセプト (またはそのキーワード) は抽出されません。

表 39. 不要語エントリーの例 (続き)

投入	例	結果
句	<i>for example</i>	<i>for example</i> という句が含まれている場合、コンセプト (またはそのキーワード) は抽出されません。
一部	<i>copyright*</i>	<i>copyrighted</i> 、 <i>copyrighting</i> 、 <i>copyrights</i> 、または <i>copyright 2010</i> のように、 <i>copyright</i> の変異形に合致または変異形を含むコンセプト (またはそのキーワード) は不要語として登録されます。
一部	<i>*ware</i>	<i>freeware</i> 、 <i>shareware</i> 、 <i>software</i> 、 <i>hardware</i> 、 <i>beware</i> 、または <i>silverware</i> のように、 <i>ware</i> の変異形に合致または変異形を含むコンセプト (またはそのキーワード) は不要語として登録されます。

エントリーを追加するには

- テーブルの一番上の空白行で、キーワードを入力します。入力したキーワードが色つきで表示されます。この色は、キーワードが出現するタイプを示します。キーワードが黒で表示されている場合、キーワード辞書にはこのキーワードがないことを意味します。

エントリーを無効化するには

不要語辞書でエントリーを無効化して、エントリーを一時的に削除することができます。エントリーを無効化すると、そのエントリーは抽出時に無視されます。

1. 不要語辞書で、無効化したいエントリーを選択します。
2. スペース・バーをクリックします。エントリーの左側にあるチェック・ボックスがオフになります。

注: エントリーの左側にあるチェック・ボックスをオフにして、無効化することもできます。

エントリーを削除するには

不要語辞書の不要なエントリーを削除することもできます。

1. 不要語辞書で、削除したいエントリーを選択します。
2. メニューから「編集」>「削除」を選択します。エントリーは辞書内から除外されます。

第 17 章 拡張リソースについて

キーワード辞書、不要語辞書および類義語辞書のほか、Fuzzy Grouping 設定や固有表現キーワード辞書など、さまざまな拡張リソースの設定を使用することもできます。これらのリソースは、テンプレート・エディター または リソース・エディター ビューの「高度なリソース」タブで処理することができます。

「拡張リソース」タブに移動すると、次の情報を編集できます。

- リソースの対象言語。 リソースが作成され、調整される言語を選択するために使用されます。詳しくは、205 ページの『リソースの対象言語』のトピックを参照してください。
- あいまいグループ化 (例外)。 リソースが作成され調整される言語の選択に使用されます。詳しくは、205 ページの『Fuzzy Grouping』のトピックを参照してください。
- 固有表現。 抽出時に適用される正規表現や正規化規則のほか、どの固有表現を抽出できるかを有効化および無効化するために使用します。詳しくは、206 ページの『固有表現』のトピックを参照してください。
- 言語処理。 文を構築する (抽出パターンおよび強制定義) そして選択した言語の略語を使用する特別な方法を宣言するために使用します。詳しくは、211 ページの『言語処理』のトピックを参照してください。

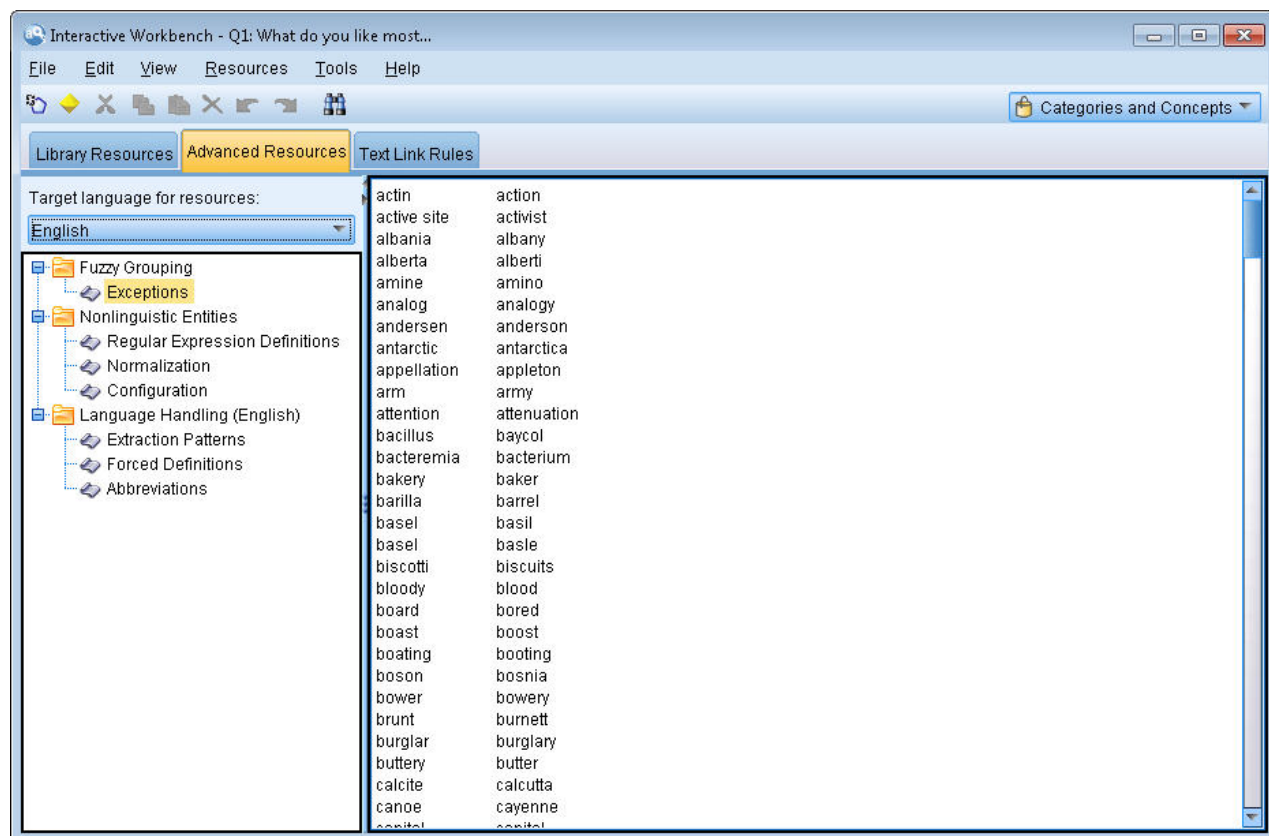


図 40. テキストマイニング・テンプレート・エディター - 「詳細リソース」タブ

注: 検索/置換ツールバーを使用して、情報を迅速に検索したり、一様の変更をセクションに行ったりすることができます。詳しくは、『置換』を参照してください。

アドバンスリソースを編集するには

1. 編集するリソース・セクションを検索および選択します。内容が右側のパネルに表示されます。
2. メニューまたはツールバーのボタンを使用して、必要に応じて内容を切り取り、コピー、または貼り付けることができます。
3. このセクションで、書式設定規則を使用して、変更したいファイルを編集します。編集を行うとすぐに、変更が保存されます。ツールバーの取り消しまたはやり直しの矢印を使用して、以前の変更に戻ります。

検索

特定のセクションで、情報の迅速な検索が必要な場合があります。例えば、テキスト リンク分析を実行する場合、多くのマクロおよびパターン定義がある場合があります。検索機能を使用して、特定の規則をすぐに見つけることができます。セクション内の情報を検索するために、検索ツールバーを使用できます。

検索機能を使用するには

1. 検索するリソース・セクションを検索および選択します。内容がエディターの右側のパネルに表示されます。
2. メニューから「編集」>「検索」を選択します。「拡張リソースを編集」ダイアログ・ボックスの右上に、検索ツールバーが表示されます。
3. テキスト・ボックスに検索したい文字列を入力します。ツールバー・ボタンを使用して、大文字と小文字の区別、部分一致、検索の方向を制御します。
4. 「検索」を選択して、検索を開始します。一致が見つかった場合は、テキストがウィンドウで強調表示されます。
5. 「検索」 もう一度クリックすると、次の一致を検索します。

注: 「テキスト リンク規則」タブ使用時は、「検索」オプションは、ソース・コードの表示中のみ使用できます。

置換

拡張リソースへより広い更新が必要な場合があります。置換機能を使用すると、内容に一様の更新を行うことができます。

置換機能を使用するには

1. 検索および置換するリソース・セクションを検索および選択します。内容がエディターの右側のパネルに表示されます。
2. メニューから「編集」>「置換」を選択します。「置換」ダイアログ・ボックスが開きます。
3. 「検索対象」 テキスト・ボックスで、検索する文字列を入力します。
4. 「以下で置換」 テキスト・ボックスで、見つかったテキストの代わりに使用する文字列を入力します。
5. 完全な語のみを検索または置換する場合、「語全体の一致のみ」を選択します。
6. 大文字と小文字が完全に一致する語のみを検索または置換する場合、「大文字/小文字を区別」を選択します。

7. 「次を検索」 を選択して、一致を検索します。一致が見つかった場合は、テキストがウィンドウで強調表示されます。この一致を置換したくない場合は、置換したい一致が見つかるまで 「次を検索」 をクリックします。
8. 「置換」 を選択して、選択した一致を置換します。
9. 「置換」 を選択して、セクション内のすべての一致を置換します。置換が行われた数を示したメッセージが表示されます。
10. 置換が終了したら、「閉じる」 をクリックします。ダイアログ・ボックスがクローズします。

注: 誤って置換を行った場合、ダイアログ・ボックスを閉じ、メニューの「編集」>「取り消し」を選択して、置換を取り消すことができます。取り消したい変更ごとに 1 回ずつ実行する必要があります。

リソースの対象言語

リソースは、特定のテキスト言語で作成されます。これらのリソースが調整される言語は、「詳細リソース」タブで定義されます。必要に応じて、「リソースの対象言語」コンボボックスでその言語を選択し、別の言語に切り替えることができます。また、ここに示された言語は、これらのリソースで作成するテキスト分析パッケージの言語として表示されます。

重要: リソースの言語を変更する必要は、めったにありません。言語を変更すると、リソースが抽出言語と合致しない場合に問題が発生する場合があります。また、あまり採用されませんが、複数の言語のテキストが期待されるため、抽出時にすべての言語オプションを使用する場合に、言語を変更する場合があります。言語を変更することにより、例えば、関心のある二次言語に関して、抽出パターン、省略形、強制定義などの言語処理リソースにアクセスすることができます。ただし、リソースに対して行った変更を公開または保存したり、その他の抽出を実行する前に、その言語を、抽出対象として関心のある一次言語に設定し戻すように留意してください。

Fuzzy Grouping

テキスト マイニング・ノードおよび抽出設定で、「語幹文字数が次の最小値以上のときにスペルを調整する」を選択している場合、Fuzzy Grouping アルゴリズムが有効化されます。

Fuzzy Grouping を使用すると、抽出語から最初の母音を除くすべての母音および二重および三重の子音を一時的に取り除いて比較し、残りが同じかどうかを確認して、一般的にスペルミスのある語やスペルの近い語をグループ化することができます。抽出プロセス時、Fuzzy Grouping 機能は、抽出キーワードに適用され、結果を比較して、一致があるかどうかを判断します。一致があった場合、元のキーワードは最終的な抽出リストにいっしょにグループ化されます。それらのキーワードは、データ内で最も頻繁に出現するキーワードに基づいてグループ化されます。

注: 比較されている 2 つのキーワードが <Unknown> タイプを除く異なるタイプに割り当てられている場合、このペアに Fuzzy Grouping 手法は適用されません。つまり、この手法を適用するには、キーワードが同じタイプまたは <Unknown> タイプに割り当てられている必要があります。

この機能を有効にして、類似したスペルの 2 つのキーワードが不正にグループ化されていることがわかった場合、それらを Fuzzy Grouping から除外する必要があります。不正に一致したペアを「拡張リソース」タブの「例外」セクションに入力することによって、Fuzzy Grouping から除外することができます。詳しくは、203 ページの『第 17 章 拡張リソースについて』のトピックを参照してください。

次の例は、Fuzzy Grouping がどのように実行されるかについて説明しています。Fuzzy Grouping が有効な場合、これらの語が同じであると示され、以下のように一致します。

```

color -> colr           mountain -> montn
colour -> colr          montana -> montn

modeling -> modlng      furniture -> furntr
modelling -> modlng     furnature -> furntr

```

前者の例では、mountain および montana をグループ化から除外することが考えられます。次のように「例外」セクションに入力します。

```
mountain      montana
```

重要: Fuzzy Grouping の例外では、特定の類義語規則が適用されているため、2 つの単語がグループ化される場合があります。この場合、感嘆符のワイルドカード (!) を使用して類義語を入力し、単語が出力で類義語とならないようにする必要があります。詳しくは、198 ページの『類義語の定義』を参照してください。

Fuzzy Grouping の例外の書式規則

- 1 行につき 1 つの例外ペアのみを定義する。
- 単語または複合語を使用する。
- 語には小文字の実を使用する。大文字は無視されます。
- ペア内の各語を区切るには TAB 文字を使用する。

固有表現

特定の種類のデータを扱っている場合、日付、社会保障番号、またはその他の固有表現を抽出したいと考える場合があります。これらのエンティティは、エンティティを有効化または無効化できる構成ファイルで明示的に宣言されています。詳しくは、210 ページの『構成』のトピックを参照してください。抽出エンジンから出力を最適化するために、非言語処理からの入力を正規化し、事前に定義された形式に従って同様のエンティティをグループ化します。詳しくは、209 ページの『正規化』のトピックを参照してください。

注: 抽出設定で、固有表現の抽出を有効にしたり無効にしたりできます。

使用できる固有表現

次の表の固有表現を抽出できます。カッコ内はタイプ名です。

表 40. 抽出できる固有表現

住所	(<Address>)
アミノ酸	(<Aminoacid>)
通貨	(<Currency>)
日付	(<Date>)
遅延	(<Delay>)
桁	(<Digit>)
電子メール・アドレス	(<email>)
HTTP/URL アドレス	(<url>)
IP アドレス	(<IP>)
組織	(<Organization>)
パーセンテージ	(<Percent>)

表 40. 抽出できる固有表現 (続き)

製品	(<Product>)
プロテイン	(<Gene>)
電話番号	(<PhoneNumber>)
時刻	(<Time>)
米国社会保障番号	(<SocialSecurityNumber>)
計量	(<Weights-Measures>)

処理するテキストの削除

固有表現の抽出を行う前に、入力テキストが削除されます。この手順で、次の一時的な変更が行われ、固有表現が以下のようなものとして特定および抽出されます。

- 複数のスペースの行列は、1 つのスペースに置き換えられます。
- 表形式はスペースに置き換えられます。
- 単一の行末文字または配列文字はスペースに置き換えられ、複数の行末文字の配列はパラグラフの最後としてマークされます。行末は、復帰改行 (CR) および改行 (LF) またはそれらの両方で示されます。
- HTML タグおよび XML タグは一時的に除外および無視されます。

正規表現の定義

固有表現抽出時、正規表現の特定に使用される正規表現定義の編集またはそれへの追加が必要な場合があります。これは「拡張リソース」タブの「正規表現の定義」セクションで行われます。詳しくは、203 ページの『第 17 章 拡張リソースについて』のトピックを参照してください。

ファイルはいくつかのセクションに分割されます。最初のセクションは [macros] です。このセクションのほか、固有表現ごとにセクションが存在する場合があります。このファイルにセクションを追加できます。各セクション内で、規則には番号が付けられています (*regex1*、*regex2*、など)。これらの規則は 1-*n* の順に番号付けされます。番号が欠けていると、このファイルの処理がすべて一時停止します。

エンティティーが言語に依存する場合があります。構成ファイルの言語パラメーターの値が 0 以外の場合、エンティティーは言語に依存すると見なされます。詳しくは、210 ページの『構成』のトピックを参照してください。エンティティーが言語に依存する場合、[english/PhoneNumber] のようにセクション名の前に言語を示す必要があります。PhoneNumber エンティティーの言語の値に 2 が指定されている場合、このセクションには英語の電話番号にのみ適用される規則が含まれます。

重要: エディターでこのファイルや他のファイルに変更を行い、抽出エンジンが必要に応じて機能しない場合、ツールバーの「元のものにリセット」オプションを使用して、ファイルを付属の下の内容に戻します。このファイルは、正規表現に対する特定のレベルの親密度が必要です。この領域においてさらに支援が必要な場合、IBM Corp. にご連絡ください。

特殊文字。[] {} () \ * + ? | ^ \$

以下の特殊文字を除くすべての文字はそれ自身に一致します。これらの特殊文字は、表現内の特別な目的に使用されます。.[{()}*+?|^\$これらの文字をその文字として使用するには、定義で文字の前にバックスラッシュ (\) を追加する必要があります。

例えば、Web アドレスを抽出しようとしている場合、終止符はエンティティーにとって非常に重要な文字であるため、次のようにバックスラッシュを追加する必要があります。

www¥.[a-z]+¥.[a-z]+

繰り返し演算子および識別子: ? + * {}

定義をより柔軟性のあるものにするために、正規表現に標準的なワイルドカードをいくつか使用できます。使用できるのは * ? +

- アスタリスク * は、0 以上の先行文字列があることを示します。例えば、ab*c は、「ac」、「abc」、「abbbc」などに一致します。
- プラス記号 + は、1 つ以上の先行文字列があることを示します。例えば、ab+c は、「abc」、「abbc」、「abbbc」などに一致しますが、「ac」には一致しません。
- 疑問符 ? 0 または 1 つの先行文字列があることを示します。例えば、modell?ing は、「modeling」、および「modeling」のどちらにも一致します。
- 繰り返し制限を示す中カッコ {} は、繰り返しの境界を示します。例えば、

[0-9]{n} は、ちょうど *n* 回繰り返された数値に一致します。例えば、[0-9]{4} は「1998」に一致し、「33」および「19983」には一致しません。

[0-9]{n} は、*n* 回以上繰り返された数値に一致します。例えば、[0-9]{3,} は「199」または「1998」に一致しますが「19」には一致しません。

[0-9]{n,m} は、*n* から *m* 回繰り返された数値に一致します。例えば、[0-9]{3,5} は "199"、"1998"、"19983" に一致しますが、"19" や "199835" とは一致しません。

任意のスペースおよびハイフン

定義内に任意のスペースを追加したい場合があります。例えば、「*uruguayan pesos*」、「*uruguayan peso*」、「*uruguay pesos*」、「*uruguay peso*」、「*pesos*」または「*peso*」などの通貨を抽出したい場合、スペースで区切られた 2 つの単語があるという事実に対処する必要があります。この場合、この定義は (uruguayan |uruguay)?pesos? として記述されます。pesos/peso と共に使用する場合、uruguayan または uruguay の後にはスペースが続くため、任意のスペースは、任意の行列 (uruguayan |uruguay) 内で定義される必要があります。(uruguayan|uruguay)? pesos? のように任意の行列内にスペースがない場合、スペースが必要であるため「pesos」または「peso」とは一致しません。

リスト内にハイフンを含む一連の文字を探している場合、ハイフンを最後に定義する必要があります。例えば、コンマ (,) またはハイフン (-) を検索する場合、[, -] を使用し、[-,] は決して使用しません。

リストおよびマクロの文字列の順序

短い行列の前に最も長い行列を定義する必要があります。短い行列に一致があるため、最も長い行列が読み取られなくなるためです。例えば、billion"または"bill"という文字列を検索している場合、billion" を"bill"の前に定義する必要があります。つまり、(billion|bill) ではなく、(bill|billion) となります。マクロは一連の文字列であるため、これはマクロにも適用されます。

定義セクションの規則の順序

1 行ごとに 1 つの規則を定義します。各セクション内で、規則には番号が付けられています (regex1、regex2、など)。これらの規則は 1-*n* の順に番号付けされます。番号が欠けていると、このファイルの処理がすべて一時停止します。エントリーを無効にするには、正規表現の定義に使用する各行の初めに # 記号を追加します。エントリーを有効にするには、行の前の # 文字を削除します。

各セクションで、確実に処理するために、最も具体的な規則を最も一般的な規則の前に定義する必要があります。例えば、「*month year*」の形式と「*month*」の形式で日付を検索したい場合、「*month year*」の規則を「*month*」の規則の前に定義する必要があります。次に、その定義を示します。

```
#@# January 1932
regexp1=$(MONTH),? [0-9]{4}
```

```
#@# January
regexp2=$(MONTH)
```

そして、次のようには定義してはいけません。

```
#@# January
regexp1=$(MONTH)
```

```
#@# January 1932
regexp2=$(MONTH),? [0-9]{4}
```

規則でのマクロの使用

いくつかの規則で特定の行列を使用している場合、マクロを使用できます。この行列の定義を変更する必要がある場合、該当する箇所を一度だけ変更する必要があります、それを参照するすべての規則を変更する必要はありません。例えば、次のようなマクロがあるとします。

```
MONTH=((january|february|march|april|june|july|august|september|october|
november|december)|(jan|feb|mar|apr|may|jun|jul|aug|sep|oct|nov|dec)(¥.)?)
```

マクロの名前を参照する場合、必ず `$()` で囲みます。例: `regexp1=$(MONTH)`

すべてのマクロは `[macros]` で定義する必要があります。

正規化

固有表現を抽出する場合、出現したエンティティは正規化され、事前定義された形式に従って同様のエンティティとグループ化されます。例えば、語内の通貨記号および同格の文字は同じように扱われます。正規化エントリーは、「拡張リソース」タブの「正規化」セクションで行われます。詳しくは、203 ページの『第 17 章 拡張リソースについて』のトピックを参照してください。ファイルはいくつかのセクションに分割されます。

重要: このファイルはアドバンス・ユーザーのみが使用できます。このファイルの変更が必要な可能性はほとんどありません。この領域においてさらに支援が必要な場合、IBM Corp. にご連絡ください。

正規化の書式規則

- 1 行につき正規化エントリーを 1 つだけ追加する。
- このファイルのセクションを厳密に重視する。新しいセクションを追加することはできません。
- エントリーを無効にするには、その行の初めに # 記号を追加する。エントリーを有効にするには、行の前の # 文字を削除します。

正規化での英語の日付

デフォルトでは、英語版テンプレート内の日付はアメリカ式日付形式 (月、日、年) で認識されます。それを「日、月、年」に変更する必要がある場合、`disable the "format:US"` 行を無効に (行の先頭に # を追加) して `"format:UK"` 行を有効に (行の先頭の # を削除) してください。

構成

固有表現構成ファイルに抽出したい固有表現タイプを有効化および無効化することができます。必要のないエンティティを無効にすることによって、必要な処理時間を短縮することができます。これは「拡張リソース」タブの「構成」セクションで行われます。詳しくは、203 ページの『第 17 章 拡張リソースについて』のトピックを参照してください。固有表現の抽出が有効になると、抽出エンジン抽出時にこの構成ファイルを読み取り、固有表現タイプを抽出するかどうかを判断します。

このファイルのシンタックスは次のようになります。

```
#name<TAB>Language<TAB>Code
```

表 41. 構成ファイルのシンタックス :

列ラベル	説明
#name	固有表現の語は、固有表現抽出の他の 2 つの必須ファイルで参照されます。ここで使用される名前は大文字と小文字が区別されます。
Language	ドキュメントの言語。特定の言語を選択することが最適ですが、「任意」オプションもあります。指定できるオプションは次のとおりです。0 = IP/URL/電子メール・アドレスなど、正規表現が言語特有でなく、さまざまな言語のいくつかのテンプレートで使用できる任意のオプション、1 = フランス語、2 = 英語、4 = ドイツ語、5 = スペイン語、6 = オランダ語、8 = ポルトガル語、10 = イタリア語です。
Code	品詞コード。多くのエンティティは、一部の場合を除いて「s」の値をとります。指定できる値は次のとおりです。s = stopword; a = adjective; n = noun 有効な場合、固有表現が最初に抽出され、より大きなコンテキストでの役割を識別するために、抽出パターンが適用されます。例えば、割合は a の値が与えられます。30% が固有表現として抽出されるとします。固有表現として 30% が抽出されるとします。それは形容詞として特定されます。テキストに「30% salary increase」という文字列が含まれていた場合、「30%」の固有表現は「ann」(形容詞、名詞、名詞)の品詞パターンに適合します。

エンティティ定義の順序

このファイルでエンティティが宣言される順序は重要であり、それらがどのように抽出されるかに影響を与えます。表示された順序に適用されます。順序を変更すると、結果も変わります。最も具体的な固有表現は、最も一般的な固有表現の前に定義する必要があります。

例えば、固有表現「Aminoacid」は次のように定義されます。

```
regex1=($(AA)-?$(NUM))
```

\$(AA) は、特定のアミノ酸に対応する固有の 3 文字の行列である、

「(ala|arg|asn|asp|cys|gln|glu|gly|his|ile|leu|lys|met|phe|pro|ser)」に対応します。

一方、固有表現 Gene はより一般的であり、次のように定義されます。

```
regex1=p[0-9]{2,3}
regex2=[a-z]{2,4}-?[0-9]{1,3}-?[r]
regex3=[a-z]{2,4}-?[0-9]{1,3}-?p?
```

「構成」セクションで「Gene」を「Aminoacid」の前に定義すると、「Gene」の「regex3」が最初に一致するため、「Aminoacid」は一致しなくなります。

構成の書式規則

- 列内の各語を区切るには TAB 文字を使用する。

- 行を削除しない。
- 前述の表示されたシンタックスを重視する。
- エントリーを無効にするには、その行の初めに # 記号を追加する。エンティティを有効にするには、行の前の # 文字を削除します。

言語処理

現在使用される各言語には、キーワードを表現、文を構築、省略形を使用する特別な方法があります。「言語処理」セクションでは、抽出パターンを編集、これらのパターンの定義を強制、「言語」ドロップダウン・リストで選択した言語の省略形を宣言することができます。

- 抽出パターン
- 強制定義
- 省略形

抽出パターン

ドキュメントから情報を抽出する場合、抽出エンジンは、品詞抽出パターンのセットをテキストの単語の「積み重ね」に適用して、抽出の候補のキーワード (単語および句) を特定します。抽出パターンを追加または変更できます。

品詞には、名詞、形容詞、過去分詞、決定詞、前置詞、人の名、イニシャル、助詞など、文法的な要素が含まれます。これら一連の要素が、品詞の抽出パターンを構成しています。IBM Corp. のテキスト マイニング製品では、各品詞が 1 つの文字で表され、パターンが定義しやすくなります。例えば、形容詞は小文字の「a」で表されます。デフォルトでは、サポートされているコードのセットがデフォルトの抽出パターンの上に、パターンのセットと、各パターンの例と共に表示され、使用される各コードについて理解しやすくなります。

抽出パターンの書式規則

- 1 行ごとに 1 つのパターン。
- 行頭に # を使用してパターンを無効化する。

単語の指定された順序は抽出エンジンによって一度だけ読み込まれ、エンジンが合致を検出した最初の抽出パターンに割り当てるため、抽出パターンの表示順が非常に重要になります。

サポートされる品詞コード

英語のコンパイル済み辞書で定義されている、サポートされるすべての品詞コードを以下の表に示します。

特定のテンプレートで使用されているすべての品詞は、「詳細リソース」 > 「抽出パターン」の上部にリストされます。

基本リソース テンプレートと意見テンプレートの主な違いは、基本テンプレートで最小の限定詞 (d) と前置詞 (c) を使用している場合に、意見テンプレートでは範囲の広いもの (e と r) を使用していることです。また、意見テンプレートでは、a と Q の両方の品詞を持つすべての単語を単に Q として処理します。0、1、および 2 は、すべての意見テンプレートで限定的に使用されています。「詳細リソース」 > 「言語処理 (英語)」 > 「強制定義」および「抽出パターン」を参照してください。

それ以外の英語のテンプレートでは、辞書にリストされていない一部の品詞を使用している場合があります (例えば、Market Intelligence テンプレートでは w と W を使用しています)。ただしその場合、それらの品詞は「詳細リソース」 > 「強制定義」で特定の単語に割り当てられています。

表 42. サポートされる品詞コード

コード	意味	例
a	形容詞	abdominal、blue...
A	未使用	未使用
b	副詞	frequently、often、very、...
B	未使用	未使用
c	前置詞	/
C	ミススペルの単語に対する内部コード	
d	限定詞	the
D	未使用	未使用
e	広範囲	限定詞 the、an、my、your...
E	未使用	未使用
f	名	John、Mary...
F	未使用	未使用
g	未使用	未使用
G	国籍の形容詞	french、american...
h	未使用	未使用
H	未使用	未使用
i	未使用	未使用
I	未使用	未使用
j	未使用	未使用
J	未使用	未使用
k	未使用	未使用
K	未使用	未使用
l	未使用	未使用
L	未使用	未使用
m	名詞または不明語	dog、ibm
M	未使用	未使用
n	名詞	dog
N	未使用	未使用
o	接続詞	and、&
O	未使用	未使用
p	過去分詞	abandoned、accessorized...
P	未使用	未使用
q	未使用	未使用
Q	修飾子	expensive、small、good、...
r	広範囲の前置詞	of、among、against、from...
R	未使用	未使用
s	ストップワード	抽出対象外のすべての単語
S	未使用	未使用
t	敬称	mrs.、mrs、captain、brig、...

表 42. サポートされる品詞コード (続き)

コード	意味	例
T	専門用語である形容詞	tumor-restricted... (T はいずれも a でもあります)
u	定義により不明。辞書に存在せず	
U	未使用	未使用
v	動詞	eat、eats、ate、eating、...
V	不定詞の動詞	eat、...
w	未使用	未使用
W	未使用	未使用
x	助動詞	be
X	未使用	未使用
y	小辞	von、di、de、... (人名の抽出に使用します。例: John von Doe)
Y	未使用	未使用
z	未使用	未使用
Z	未使用	未使用
0	意見の副詞	意見のみで使用します。「詳細リソース」 > 「言語処理 (英語)」 > 「強制定義」を参照してください。
1	to (意見の場合)	「詳細リソース」 > 「言語処理 (英語)」 > 「強制定義」を参照してください。
2	特定の識別子	意見のみで使用します。「詳細リソース」 > 「言語処理 (英語)」 > 「強制定義」を参照してください。
3	未使用	未使用
4	未使用	未使用
5	未使用	未使用
6	未使用	未使用
7	未使用	未使用
8	未使用	未使用
9	未使用	未使用

強制定義

ドキュメント から情報を抽出するとき、抽出エンジンはテキストをスキャンし、出現するすべての単語の品詞を特定します。状況によっては、単語がいくつかの異なる役割に適合する場合があります。単語が特定の品詞の役割を取得するよう強制する場合、または処理からその単語を完全に除外する場合、「拡張リソース」タブの「強制定義」セクションで指定できます。詳しくは、203 ページの『第 17 章 拡張リソースについて』のトピックを参照してください。

指定した単語の品詞を強制するには、次のシンタックスを使用して、このセクションに 1 行追加します。

`term:code`

表 43. シンタックスの説明：

投入	説明
term	キーワード名。

表 43. シンタックスの説明 (続き):

投入	説明
code	品詞の役割を示す 1 文字のコード。ユニタームあたり最大 6 つの品詞コードを表示できます。また、 <code>additional:s</code> のように小文字のコード <code>s</code> を使用して、単語が複合語/句に抽出されないようにできます。

強制定義の書式設定規則

- 単語ごとに 1 行。
- キーワードにコロンの使用不可。
- 単語を抽出しないようにするには小文字の `s` を品詞コードに使用する。
- 1 行あたり最大 6 つの品詞コードを使用する。サポートされている品詞コードは、「抽出パターン」セクションに表示されます。詳しくは、211 ページの『抽出パターン』のトピックを参照してください。
- 部分一致の場合、文字列の最後にアスタリスク (*) をワイルドカードとして使用する。例えば、`add*:s` と入力すると、`add`、`additional`、`additionally`、`addendum`、および `additive` のような単語はキーワードとして、あるいは複合語キーワードの一部として抽出されなくなります。ただし、単語の合致がコンパイル済み辞書または強制定義でキーワードとして明示的に宣言されている場合は、抽出されます。例えば、`add*:s` および `addendum:n` を入力すると、`addendum` がテキスト内で見つかった場合に抽出されます。

省略形

抽出エンジンがテキストを処理しているとき、検出されたピリオドは、文の終了を示す指標として認識されます。これは通常、正常な処理ですが、省略形がテキスト内にある場合、ピリオド文字のこうした処理は適用されません。

テキストからキーワードを抽出し、特定の省略形の処理が不適切であったことがわかった場合、このセクションで省略形を明示的に宣言する必要があります。

注: 省略形が類義語定義に出現、またはキーワード辞書でキーワードとして定義されている場合、ここで省略形の投入を追加する必要はありません。

省略形の書式設定規則

- 1 行ごとに 1 つの省略形を定義する。

第 18 章 テキスト リンク規則について

テキスト リンク分析 (TLA) はパターンマッチ手法で、一連の条件規則を使用して、テキスト内の関係性を抽出するのに使用されます。テキスト リンク分析が抽出に対して有効である場合、テキスト・データがこれらの規則に対して比較されます。合致が検出されると、テキスト リンク分析パターンが抽出され、表示されます。これらの条件規則は、「テキスト リンク規則」タブで定義されます。

例えば、組織に関する単純なアイデアを示すコンセプト抽出してもそれが重要でない場合がありますが、TLA を使用して、さまざまな組織または組織に関連する人々の間のつながりについて学習することができます。TLA を使用して、指定された製品または経験についてどう思うかなどのトピックについての意見を抽出することもできます。

TLA を利用するには、テキスト リンク分析 (TLA) 規則を含むリソースが必要です。テンプレートを選択すると TLA 列にアイコンがあるかどうかによって、どのテンプレートに TLA 規則があるかを確認できます。

テキスト・データのテキスト リンク分析パターンは、抽出プロセスのパターン・マッチの段階で検出されます。この段階で、条件規則がテキストデータと比較され、合致が検出されると情報がパターンとして抽出されます。テキスト リンク分析からより多くの情報を抽出したり、どのように合致するかを変更したりすることが必要な場合があります。こうした場合、条件規則が特定のニーズに適応するよう、規則を処理することができます。これは「テキスト リンク規則」タブで実行します。

注: 変数のサポートは、バージョン 13 では継続されていません。代わりにマクロを使用してください。詳しくは、220 ページの『マクロの作業』のトピックを参照してください。

テキスト リンク規則を扱う場所

テンプレート・エディター または リソース・エディター ビューの「テキスト リンク規則」タブで直接、条件規則を編集および作成できます。条件規則がテキストとどのように合致するかを確認するために、このタブでシミュレーションを実行できます。シミュレーション時、抽出はサンプルのシミュレーション・データでのみ実行され、テキスト リンク規則を適用してパターン・マッチがあるかどうかを確認します。テキストに合致する規則は、シミュレーション・パネルに表示されます。合致に基づいて、条件規則およびマクロを編集し、テキストがどのように合致するかを変更することができます。

他の高度なリソースとは異なり、TLA 規則はライブラリー固有の定義です。そのため、一度に 1 つのライブラリーの TLA 規則を使用できます。テンプレート・エディター または リソース・エディター で、「テキスト リンク規則」タブに移動します。使用したい、または編集したい TLA 規則を含むテンプレートのライブラリーを選択します。このため、特別な理由がない限り、すべての規則を 1 つのライブラリーに保存することを強くお勧めします。

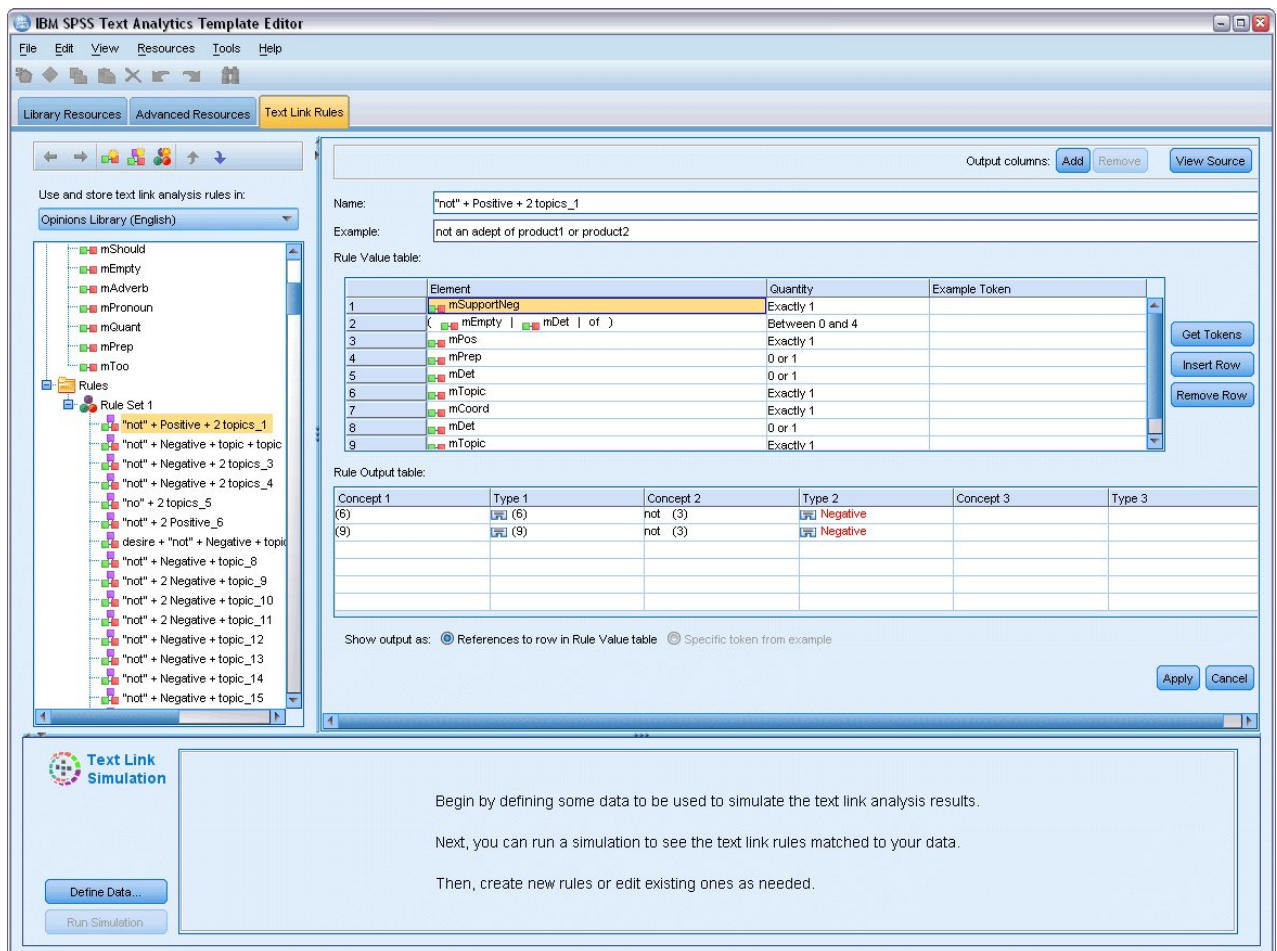


図 41. 「テキスト リンク規則」タブ

作業の開始

「テキスト リンク規則タブ」エディターで作業を開始するにはさまざまな方法があります。

- いくつかのサンプル テキストで結果をシミュレーションし、現在のセットの規則がシミュレーションからパターンをどのように抽出するかに基づいて合致規則を編集または作成する。
- スクラッチから新しい規則を作成するか、既存の規則を編集する。
- ソース・ビューで直接作業する。

規則の編集または作成が必要な場合

各テンプレートに付属するテキスト リンク分析規則は、多くの単純または複雑な関係性をテキストから抽出することに適している場合が多いですが、これらの規則に変更を加えるか、独自の規則をいくつか作成することが必要な場合があります。以下に例を示します。

- 新しい規則またはマクロを作成して既存の規則で抽出されていないアイデアまたは関係性をキャプチャーする。

- リソースに追加したタイプのデフォルトの動作を変更する。通常、mTopic または mNonLingEntities などのマクロを編集することが必要です。詳しくは、223 ページの『特殊マクロ：mTopic、mNonLingEntities、SEP』のトピックを参照してください。
- 既存のテキスト リンク分析規則およびマクロに新しいタイプを追加する。例えば、タイプ<組織名>が広すぎる場合、<医薬品>、<自動車メーカー>、<金融>など、いくつかの企業部門に分かれた組織の新しいタイプを作成することができます。この場合、テキスト リンク分析規則を編集またはマクロを作成して、これらの新しいタイプを考慮に入れ、適宜処理する必要があります。
- 既存のテキスト リンク分析規則およびタイプを追加する。例えば、「john doe called jane doe」というテキストをキャプチャーする規則があり、通話をキャプチャーするこの規則を使用して、電子メールのやり取りもキャプチャーします。電子メールの固有表現タイプを規則に追加し、次のようなテキストもキャプチャーすることができます：johndoe@ibm.com emailed janedoe@ibm.com.
- 規則を新規作成するのではなく、既存の規則を若干変更する。例えば、「xyz is very good」というテキストに合致する規則があり、この規則を使用して、「xyz is very, very good」もキャプチャーします。

テキスト リンク分析結果のシミュレーション

新しいテキスト リンク規則を定義できるようにするために、またはテキスト リンク分析時に特定の文がどのように合致するかを理解できるようにするために、テキストのサンプル部分を抽出してシミュレーションを実行すると役立ちます。シミュレーション時、抽出はサンプルのシミュレーション・データでのみ実行され、テキスト リンク規則を適用してパターン・マッチがあるかどうかを確認します。目標は、シミュレーションの結果を取得し、これらの結果を使用して規則を改善、新しい規則を作成するか、どのように合致が出現するかをより良く理解することです。テキストの各部分 (状況によって文、語、句) に対し、シミュレーションの出力には、トークンの集合およびそのテキストのパターンを明らかにしなかった TLA 規則を示します。トークンは、抽出プロセス時に特定される単語または語句として定義されます。

他の高度なリソースとは異なり、TLA 規則はライブラリー固有の定義です。そのため、一度に 1 つのライブラリーの TLA 規則を使用できます。テンプレート・エディター または リソース・エディター で、「テキスト リンク規則」 タブに移動します。使用したい、または編集したい TLA 規則を含むテンプレートのライブラリーを選択します。このため、特別な理由がない限り、すべての規則を 1 つのライブラリーに保存することを強くお勧めします。

重要: データ・ファイルを使用する場合、含まれるテキストを処理時間短縮のために短くなっていることを確認してください。シミュレーションの目的は、テキストの一部がどのように読み取られるかを確認し、このテキストに条件規則がどのように合致するかを理解することです。この情報を使用して、条件規則を作成および編集できます。テキスト リンク分析ノードを使用、または TLA 抽出を有効にしてインタラクティブ・セッションのあるストリームを実行し、より完全なデータ・セットの結果を取得します。このシミュレーションは、テストおよび条件規則作成のためだけに行われます。

シミュレーションのデータ定義

条件規則がテキストとどのように合致するかを確認するために、サンプルデータを使用してシミュレーションを実行できます。まず、データを定義します。

データの定義

1. 「テキスト リンク規則」 タブの一番下にあるシミュレーション・ペインで「データを定義」をクリックします。または、データがまだ定義されていない場合、メニューから「ツール」 > 「シミュレーションの実行」を選択します。シミュレーション・データ・ウィザードが開きます。
2. 次のいずれかを選択して、データの種別を指定します。

- データを貼り付けまたは直接入力: テキスト ボックスで、クリップボードからテキストを貼り付けるか、処理するテキストを手動で入力します。1 行ごとに 1 文ずつ入力するか、ピリオドやコンマなど、句読点を使用して文を区切ります。テキストを入力すると、「シミュレーションを実行」をクリックして、シミュレーションを開始できます。
- ファイル データ ソースを指定: このオプションでは、テキストを含む、処理対象のファイルを指定します。「次へ」をクリックすると、処理するファイルを定義できるウィザードの手順に進みます。ファイルが選択されると、「シミュレーションを実行」をクリックして、シミュレーションを開始できます。サポートされるファイル タイプは、.txt および .text です。選択したデータ・ファイルは、シミュレーション時にそのまま読み込まれます。ファイル全体は、ファイル リスト ノードをテキスト マイニング・ノードに接続した場合と同じように扱われます。

重要: データ ファイルを使用する場合、処理時間短縮のために、含まれるテキストが短いことを確認してください。シミュレーションの目的は、テキストの一部がどのように読み取られるかを確認し、このテキストに条件規則がどのように合致するかを理解することです。この情報を使用して、条件規則を作成および編集できます。テキスト リンク分析ノードを使用、または TLA 抽出を有効にしてインタラクティブ・セッションのあるストリームを実行し、より完全なデータ・セットの結果を取得します。このシミュレーションは、テストおよび条件規則作成のためだけに行われます。

3. シミュレーション プロセスを開始するには、「シミュレーションを実行」をクリックします。プロセス・ダイアログが表示されます。インタラクティブ・セッションで作業している場合、シミュレーション時に使用する抽出設定は、インタラクティブ・セッションで現在選択している抽出設定となります (コンセプトとカテゴリー・ビューの「ツール」 > 「抽出設定」を参照)。テンプレート・エディターで作業している場合、シミュレーション時に使用する抽出設定は、デフォルトの抽出設定で、テキストリンク分析ノードの「エキスパート」タブに表示されている設定と同じです。詳しくは、『シミュレーション結果の理解』を参照してください。

シミュレーション結果の理解

条件規則がテキストとどのように合致するかを確認するために、サンプルデータを使用してシミュレーションを実行し、結果を確認できます。ここで、データにより正確に適用するよう一連の規則を変更できます。抽出プロセスおよびシミュレーション・プロセスが完了すると、シミュレーションの結果が表示されます。

抽出時に特定された各「文」に対し、正確な「文」を含む情報のいくつかの部分が、この入力テキストの文にあるトークンの分析結果、そしてその文のテキストに合致した規則を含むいくつかの情報が表示されます。「文」は、抽出機能でテキストを読み取り可能な単位にどのように分割するかによって、語、文、句を意味します。

トークンは、抽出プロセス時に特定される単語または語句として定義されます。例えば、文「*My uncle lives in New York*」の場合、*my*、*uncle*、*lives*、*in*、および *new york* というトークンがあります。また、*uncle* はコンセプトとして抽出され、タイプは <不明>、そして *new york* はコンセプトとして抽出され、タイプは <地名> となります。すべてのコンセプトはトークンとなりますが、すべてのトークンがコンセプトとなるわけではありません。トークンは、マクロ、リテラル文字列、そして単語の空所の場合もあります。タイプを指定された単語または語句のみがコンセプトとなります。

インタラクティブ・セッションまたはリソース・エディターで作業している場合、コンセプト・レベルで対応します。TLA 規則はより細かく、抽出されない、またタイプ指定されない場合であっても、文内の各トークンを規則の定義に使用できます。コンセプトでないトークンを使用できると、テキストの複雑な関係性をキャプチャーする場合により柔軟な規則となります。

シミュレーション・データに複数の文がある場合、「次へ」 および 「戻る」 をクリックして、結果を前後に移動することができます。

選択したライブラリー (このタブのツリー上のライブラリー名を参照) の TLA 規則に文が合致しない場合、結果は不一致と見なされ、ボタン「次の不一致」および「前の不一致」が有効になり、合致を検出した規則のないテキストがあることを示し、これらのインスタンスにすばやく移動できます。

新しい規則を作成、規則を編集、またはリソースや抽出設定を変更したと、シミュレーションの再実行が必要な場合があります。シミュレーションを再実行するには、シミュレーション・ペインで「シミュレーションを実行」をクリックします。同じ入力データが再度使用されます。

次のフィールドおよびテーブルがシミュレーション結果に表示されます。

入力テキスト:ウィザードで定義した、シミュレーション・データからの抽出プロセスで特定される実際の「文」。「文」は、抽出機能でテキストを読み取り可能な単位にどのように分割するかによって、語、文、句を意味します。

システム ビュー:抽出プロセスが特定したトークンの集合。

- 入力テキスト・トークン: 入力テキストで検出される各トークン。トークンについては、このトピックの前の項で定義されています。
- 次のタイプを指定: トークンがコンセプトとして特定およびタイプ指定されている場合、(<不明>、<人名>、<地名> など) 関連するタイプ名がこの列に表示されます。
- 合致するマクロ: トークンが既存のマクロに合致する場合、関連するマクロ名がこの列に表示されます。

入力テキストに合致する規則: このテーブルには、入力テキストに対して合致した TLA 規則が表示されます。合致する各規則に対して、「条件規則出力」列に条件規則の名前と、その規則に関連する出力値 (コンセプト + タイプのペア) が表示されます。合致する条件規則の名前をダブルクリックすると、シミュレーション・パネル上のエディター・パネルに規則が表示されます。

「条件規則を生成」ボタン:シミュレーション・パネルでこのボタンをクリックすると、シミュレーション・パネルの上の条件規則エディター・パネルに新しい規則が表示されます。入力テキストを例として採用します。同様に、シミュレーション時にタイプ指定またはマクロに合致したトークンは、「条件規則値テーブル」の「要素」列に自動的に挿入されます。トークンがタイプ指定され、かつマクロに合致した場合、マクロ値は条件規則を単純化するために条件規則内で使用されるマクロ値となります。例えば、Basic English リソースを使用している場合、「I like pizza」という文は、シミュレーション時に <不明> とタイプ指定され、mTopic のマクロに合致します。この場合、生成された規則で mTopic が要素として使用されます。詳しくは、223 ページの『テキスト リンク規則の使用』のトピックを参照してください。

ツリー内の規則およびマクロのナビゲート

抽出時にテキスト リンク分析が実行されると、「テキスト リンク規則」タブで選択されたライブラリーに保存されているテキスト リンク規則が使用されます。

他の高度なリソースとは異なり、TLA 規則はライブラリー固有の定義です。そのため、一度に 1 つのライブラリーの TLA 規則を使用できます。テンプレート・エディター または リソース・エディター で、「テキスト リンク規則」タブに移動します。使用したい、または編集したい TLA 規則を含むテンプレートのライブラリーを選択します。このため、特別な理由がない限り、すべての規則を 1 つのライブラリーに保存することを強くお勧めします。

「テキスト リンク分析規則を使用して保存:」 次のリストで該当するライブラリーを選択して、「テキスト リンク規則」タブで処理するライブラリーを指定できます: このタブのドロップ・ダウン・リスト。抽出時にテキスト リンク分析が実行されると、「テキスト リンク規則」タブで選択されたライブラリーに

保存されているテキスト リンク規則が使用されます。そのため、複数のライブラリーにテキスト リンク規則 (TLA 規則) を定義した場合、TLA 規則がある最初のライブラリーのみがテキスト リンク分析に使用されます。このため、特別な理由がない限り、すべての規則を 1 つのライブラリーに保存することを強くお勧めします。

ツリーでマクロまたは規則を選択する場合、右側のエディター・パネルに内容が表示されます。ツリー内の項目を右クリックすると、次のようなタスクを示すコンテキスト・メニューが表示されます。

- ツリー内に新しいマクロを作成し、右側のエディターで開く。
- ツリー内に新しい規則を作成し、右側のエディターで開く。
- ツリー内に新しいルール・セットを作成する。
- 項目を切り取り、コピーおよび貼り付けて、編集を簡単にする。
- マクロ、規則、ルール・セットを削除して、リソースから除外する。
- マクロ、規則、ルール・セットを無効にして、処理時に無視されるようにする。
- 規則を上下に移動させて処理の順序を変更する。

ツリー上の警告

警告はツリー上で黄色の三角を伴って表示されます。これはどこかに問題があることを知らせるために現れます。マウス・ポインターをエラーのあるマクロやルールに当てると、説明がポップアップ表示されます。大抵の場合は、以下のように表示されます：警告:例がありません。例を入力してください この場合は例を入力します。

例を忘れた場合、または、例がルールにマッチしない場合には、トークン取得機能を使用することができませんので、ルールに沿った例を1つ入力することをお奨めします。

例が黄色でハイライト表示された場合は、TLA エディターに不明なタイプまたはマクロであることを示します。以下のメッセージも同様です：警告:不明なタイプまたはマクロです。これは、ソース・ビューの \$something で定義された項目（例えば、\$myType）が使用するライブラリーのレガシー・タイプではなく、マクロでもないことを知らせるものです。

シンタックス・チェッカーを更新するには、他のルールまたはマクロに切り替える必要があります。コンパイルは何も必要ありません。ですから、例えばルール A が例がないための警告を表す場合には、例を追加する必要があります。上位または下位のルールをクリックしてから、ルール A に戻り今は適正となったことを確認します。

マクロの作業

マクロは、タイプ、その他のマクロおよびリテラル (単語) 文字列を OR 演算子 (|) とグループ化することによって、テキスト リンク分析規則の出現を簡略化することができます。マクロを使用すると、複数のテキスト リンク分析規則でマクロを再利用して簡略化するだけでなく、テキスト リンク分析規則全体で更新する必要はなく、1 つのマクロを更新することができます。付属の TLA 規則の多くには、事前定義されたマクロが含まれています。マクロは、「テキスト リンク規則」タブの左端のパネルの上位ツリーに表示されます。

次のフィールドおよびテーブルがシミュレーション結果に表示されます。

名前: このマクロを示す一意の名前。マクロ名の接頭辞に小文字の m を指定して、規則内でマクロであることをすぐに特定できるようにすることをお勧めします。手動で規則内のマクロを参照する場合 (インライン編集またはソース・ビュー)、抽出プロセスでこの特殊な名前を認識するよう、\$ 文字の接頭辞を使用す

る必要があります。ただし、マクロ名をドラッグ・アンド・ドロップまたはコンテキスト・メニューを使用して追加する場合、その名前はマクロとして自動的に認識され、\$ は追加されません。

マクロ値テーブル:

- このマクロが示すことができるすべての値を示す多くの行。これらの値では大文字と小文字が区別されます。
- 値は、タイプ、リテラル文字列、単語の空所またはマクロの組み合わせがあります。詳しくは、トピック「230 ページの『条件規則およびマクロにサポートされている要素』」を参照してください。
- マクロの要素に値を入力するには、使用する行をダブルクリックします。タイプの参照、マクロの参照、リテラル文字列、または単語の空所を入力できるテキスト・ボックスが表示されます。または、セルを右クリックすると、共通のマクロ、タイプ名、固有表現キーワード名のリストを提供するコンテキスト・メニューが表示されます。タイプまたはマクロを参照するには、マクロ名またはタイプ名の先頭に \$ 文字を追加する必要があります (例: マクロ mTopic の場合 \$mTopic)。引数を結合する場合、カッコ () を使用して引数およびブール型演算子 OR を示す文字 | をグループ化する必要があります。
- 右側のボタンを使用して、マクロ値テーブルの行を追加または削除できます。
- それぞれの行に要素を入力します。例えば、am OR was OR is のような 3 つのリテラル文字列のいずれかを示すマクロを作成する場合、ビューの各行にそれぞれのリテラル文字列を入力すると、マクロ テーブルには3 つの行が作成されます。

マクロの作成および編集

新しいマクロを作成したり、既存のマクロを編集できます。マクロ エディターのガイドラインおよび説明に従います。詳しくは、220 ページの『マクロの作業』のトピックを参照してください。

マクロの新規作成

1. メニューの「ツール」>「新規マクロ」を選択します。または、ツリー・ツールバーの「新規マクロ」アイコンをクリックすると、エディターに新しいマクロが表示されます。
2. 一意の名前を入力して、マクロ値の要素を定義します。
3. エラーの有無を確認したら、「適用」 をクリックします。

マクロの編集

1. ツリー内のマクロ名をクリックします。右側のエディター・パネルにマクロが表示されます。
2. 変更を加えます。
3. エラーの有無を確認したら、「適用」 をクリックします。

マクロの無効化および削除

マクロの無効化

処理中にマクロが無視されるようにするには、マクロを無効にすることができます。マクロを無効にすると、この無効なマクロを参照する規則で警告またはエラーが発生する場合があります。マクロを削除および無効にする場合は注意してください。

1. ツリー内のマクロ名をクリックします。右側のエディター・パネルにマクロが表示されます。
2. 名前を右クリックします。
3. コンテキスト・メニューから、「無効にする」 を選択します。マクロのアイコンがグレーになり、マクロ自体を編集できなくなります。

マクロの削除

マクロを除外する場合、削除することができます。マクロを削除すると、このマクロを参照する規則でエラーが発生する場合があります。マクロを削除および無効にする場合は注意してください。

1. ツリー内のマクロ名をクリックします。右側のエディター・パネルにマクロが表示されます。
2. 名前を右クリックします。
3. コンテキスト・メニューから、「削除」を選択します。リストからマクロが消去されます。

エラーのチェック、保存およびキャンセル

マクロの変更の適用

マクロ エディターの外をクリックまたは「適用」をクリックすると、マクロが自動的にスキャンされ、エラーの有無が確認されます。エラーが見つかると、アプリケーションの別の部分に移る前に修正する必要があります。

ただし、あまり深刻でないエラーが検出された場合は、警告のみが表示されます。例えば、マクロにタイプまたはその他のマクロに対する不完全または参照されない定義が含まれている場合、警告メッセージが表示されます。「適用」をクリックすると、未修正の警告により、左側パネルの規則とマクロ ツリーのマクロ名の左側に警告アイコンが表示されます。

マクロを適用しても、マクロが永続的に保存されるわけではありません。適用すると、検証プロセスで、エラーおよび警告の有無をチェックします。

インタラクティブ・ワークベンチ・セッション内のリソースの保存

1. インタラクティブ・ワークベンチ・セッションでリソースに行った変更を保存し、次回ストリーム実行時にそれらの変更を取得できるようにするには、次の手順を実行する必要があります。
 - モデル作成ノードを更新して、次回ストリーム実行時にこれらの同じリソースを取得できるようにします。詳しくは、78 ページの『モデル作成ノードの更新および保存』のトピックを参照してください。その後ストリームを保存します。ストリームを保存するには、モデル作成ノードを更新した後、IBM SPSS Modeler のメイン・ウィンドウで保存します。
2. インタラクティブ・ワークベンチ・セッションでリソースに行った変更を保存し、他のストリームでそれらを使用できるようにするには、次の手順を実行します。
 - 使用していたテンプレートを更新するか、新しいテンプレートを作成します。詳しくは、164 ページの『テンプレートの作成および更新』のトピックを参照してください。現在のノードの変更を保存するわけではありません (前の手順を参照)。
 - または、使用していた TAP を更新します。詳しくは、138 ページの『テキスト分析パッケージの更新』のトピックを参照してください。

テンプレート・エディター 内のリソースの保存

1. まず、ライブラリーを公開します。詳しくは、186 ページの『ライブラリーの公開』のトピックを参照してください。
2. そして、メニューで「ファイル」>「リソース・テンプレートを保存」を使用してテンプレートを保存します。

マクロの変更のキャンセル

1. 変更を破棄する場合は、「キャンセル」をクリックします。

特殊マクロ : mTopic、 mNonLingEntities、 SEP

基本リソース・テンプレートおよび意見テンプレートは、mTopic および mNonLingEntities という 2 つの特別なマクロに付属しています。

mTopic

デフォルトでは、mTopic マクロは、何らかの意見に結びつけられる可能性のあるテンプレート付属タイプをすべてグループ化します。このような Core ライブラリー・タイプには次のようなものがあります。<Person>、<Organization>、<Location>などです。これらが該当するのは、タイプが意見タイプではない(例えば、<Negative> または <Positive>) か、または [拡張リソース] で固有表現として定義されているタイプです。

意見 (または同様の) テンプレートで新しいタイプを作成する場合、は、このタイプが別のマクロまたは「拡張リソース」タブの固有表現セクションで指定されていない場合、マクロ mTopic で定義された他のタイプと同じ方法で処理されると想定します。

例えば、意見テンプレートからリソースに新しいタイプを作成したとします : <野菜> と <果物>としましょう。変更を行うことなく、新しいタイプは mTopic として扱われ、新しいタイプに関する肯定的、否定的、中立的、文脈上の意見を自動的に明らかにすることができます。抽出時、例えば文「*I enjoy broccoli, but I hate grapefruit*」が次の出力パターンを作成するとします。

broccoli <Vegetables> + like <Positive>

grapefruit <Fruit> + dislike <Negative>

ただし、mTopic の他のタイプと異なる方法でこれらのタイプを処理する場合、タイプ名を mPos のような既存マクロに追加してすべての肯定的意見タイプをグループ化するか、1 つまたは複数の規則で後で参照できる新しいマクロを作成することができます。

重要: <Vegetables> のような新しいタイプを作成した場合、この新しいタイプは mTopic にタイプとして含まれますが、このタイプ名はマクロ定義に明示的に表示されるわけではありません。

mNonLingEntities

同様に、「拡張リソース」タブの「固有表現」セクションに新しい固有表現を追加すると、それらは特に指定のない限り mNonLingEntities として自動的に処理されます。詳しくは、206 ページの『固有表現』のトピックを参照してください。

SEP

定義済みのマクロ SEP も使用できます。これは、ローカル・コンピューターで定義されたグローバル区切り文字、通常コンマ (,) に対応します。

テキスト リンク規則の使用

テキスト リンク分析規則は、文にマッチを実行するために使用するブール型質問です。テキスト リンク分析には、1つまたは複数の引数が含まれます : タイプ、マクロ、リテラル文字列、単語の空所などです。TLA 結果を抽出するには、テキスト リンク分析規則が少なくとも 1 つ必要です。

以下は、規則エディターの「テキスト リンク規則」タブに表示される領域とフィールドです。

「名前」 フィールド: テキスト リンク分析規則の一意の名前。

「例」 フィールド: オプションで、この規則でキャプチャされる例の文または語の連鎖を含むことができます。例の使用をお勧めします。このエディター内でこのテキスト例からトークンを生成し、規則にどのように合致し、どのように出力されるかを確認できます。トークンは、抽出プロセス時に特定される単語または語句として定義されます。例えば、文「*My uncle lives in New York*」の場合、*my*、*uncle*、*lives*、*in*、および *new york* というトークンがあります。また、*uncle* はコンセプトとして抽出され、タイプは <不明>、そして *new york* はコンセプトとして抽出され、タイプは <地名> となります。すべてのコンセプトはトークンとなりますが、すべてのトークンがコンセプトとなるわけではありません。トークンは、マクロ、リテラル文字列、そして単語の空所の場合もあります。タイプを指定された単語または語句のみがコンセプトとなります。

「条件規則値テーブル」: このテーブルには、文に対し規則を合致させるために使用する規則の要素が含まれています。右側のボタンを使用して、テーブルの行を追加または削除できます。テーブルは、次の 3 つの列で構成されています。

- 「要素」 列: 1 つのタイプ、リテラル文字列、単語の空所 (<任意のトークン>)、またはマクロとして値を入力します。詳しくは、トピック「230 ページの『条件規則およびマクロにサポートされている要素』」を参照してください。要素セルをダブルクリックして、情報を直接入力します。または、セルを右クリックすると、共通のマクロ、タイプ名、固有表現キーワード名のリストを提供するコンテキスト・メニューが表示されます。セルに情報を入力する場合、マクロまたはタイプ名の先頭に \$ 文字を追加する必要があります (例: マクロ *mTopic* の場合 *\$mTopic*)。要素の行を作成する順序は、規則がテキストにどのように合致するかにおいて重要です。引数を結合する場合、カッコ () を使用して引数およびブール型演算子 OR を示す文字 | をグループ化する必要があります。値では大文字と小文字が区別されます。
- 「数量」 列: 合致が発生するために要素が検出される必要のある回数の、最大値および最小値を示します。例えば、0 から 3 個の単語のいずれかの箇所の 2 つの要素間に空所、または一連の単語を定義する場合、リストから「範囲: 0 and 3」を選択するか、ダイアログ・ボックスに直接数値を入力することができます。デフォルトは「1 に一致」です。要素をオプションにすることが必要な場合があります。その場合、最小数量が 0 および最大数量が 0 より大きくなります (例: 0 または 1、0 から 2)。規則の最初の要素をオプションにはできません。つまり最初の要素の数量を 0 にすることはできません。
- 「例のトークン」 列: 「トークンを取得」 をクリックすると、「例」 のテキストをトークンに分割し、これらのトークンを使用して、定義した要素と合致するトークンを列に入力します。出力テーブルにこれらのトークンを表示することもできます。

「条件規則出力テーブル」: このテーブルの各行は、TLA パターン出力が結果にどのように表示されるかを定義します。条件規則出力では、それぞれ「スロット」を示す最大 6 つのコンセプト/タイプ列のペアのパターンを生成できます。例えば、タイプ・パターン <地名> + <肯定的> は、2 つのコンセプト/タイプ列のペアで構成される 2 スロットのパターンです。

注: 「条件規則値テーブル」の「エレメント」列、および「条件規則出力テーブル」の「コンセプト」列の用語は、`、#、%、^、*、_、-、:、<、>、/、¥、" のいずれの文字でも始めることはできません。

言語によって、同じ基本的な考えをさまざまな方法で自由に表現できますが、同じ基本的な考えをキャプチャーするために多くの条件規則を定義する場合があります。例えば、「*Paris is a place I love*」と「*I really, really like Paris and Florence*」というテキストは同じ基本の考え「*Paris is liked*」を示しますが、異なる方法で表現され、両方をキャプチャーするには 2 つの異なる条件規則が必要です。ただし、類似した考えがグループ化されている場合、パターン結果を処理するとより簡単です。このため、2 つの異なる条件規則でこれら 2 つの語句をキャプチャーしますが、タイプ・パターン <地名> + <肯定的> のようにいずれのテキストも示すよう、2 つの条件規則に同じ出力を定義することができます。そして、出力が元のテキス

トの構造または語順と必ずしも同じでないことがわかります。さらに、このようなタイプのパターンは他の句と一致する可能性があり、以下のようなコンセプトを生み出す可能性があります：例えば、paris + like および tokyo + like です。

少ないエラーで出力をすばやく定義するには、コンテキスト・メニューを使用して、出力に表示する要素を選択できます。また、条件規則値テーブルから要素を出力にドラッグ・アンド・ドロップすることもできます。例えば、条件規則値テーブルの行 2 に mTopic マクロへの参照を含む規則があり、値を出力に投入する場合、mTopic の要素を条件規則出力テーブルの最初の列のペアにドラッグ/ドロップできます。このようにすると、選択したペアのコンセプトとタイプが自動的に投入されます。または、条件規則値テーブルの 3 番目の要素 (行 3) で定義されたタイプから出力する場合、そのタイプを条件規則値テーブルから出力テーブルの「タイプ 1」セルにドラッグします。テーブルが更新され、(3) に行の参照が表示されます。

また、出力する「コンセプト」列をダブルクリックし、\$ の後に行番号を入力 (例: \$2、条件規則値テーブルの行 2 で定義された要素を参照) して、これらの参照を手動でテーブルに入力することもできます。情報を手動で入力する場合、「タイプ」列を定義し、# の後に行番号を入力 (例: #2、条件規則値テーブルの行 2 で定義された要素を参照) する必要があります。

さらに、これらの方法を組み合わせることもできます。条件規則値テーブルの行 4 にタイプ <肯定的> があるとします。そのタイプを [タイプ 2] 列にドラッグし、[コンセプト 2] 列のセルをクリックしてそのタイプの前に「not」という単語を手動で入力できます。テーブルの出力列は not (4)、あるいは編集モードまたは入力モードの場合は not \$4 となります。「タイプ 1」列を右クリックして、例えば mTopic というマクロを選択します。するとこの出力で、以下のようなコンセプト・パターンが生成されます：car + bad。

多くの条件規則には出力が 1 つしかありませんが、複数の出力を生成できるまたは生成する必要がある場合があります。この場合、条件規則出力テーブルの行ごとに 1 つの出力を定義します。

重要: TLA パターンの間、他の言語処理が実行されます。そのため、出力が t\$3¥t#3 を読み込む場合、パターンは 3 番目の要素の最終コンセプトと 3 番目の最終タイプを、言語処理が適用された後 (類義語およびその他のグループ) に表示します。

- 出力を次の形式で表示: デフォルトでは、「条件規則値テーブルの行への参照」オプションが選択され、「条件規則値」タブで定義されているような行への数値参照を使用して出力が表示されます。前に「トークンを取得」をクリックして、条件規則値テーブルの「例のトークン」列にトークンがある場合、このオプションを選択して、これら特定のトークンの出力を表示できます。

注: 出力テーブルに表示できるコンセプト/タイプの出力ペアが十分でない場合、エディター・ツールバーの「追加」ボタンをクリックして別のペアをクリックできます。現在 3 つのペアが表示され、「追加」をクリックすると、2 つの列 (コンセプト 4 およびタイプ 4) がテーブルに追加されます。つまり、すべての条件規則の出力テーブルに 4 つのペアが表示されます。このライブラリーのルール・セットのルールがそのペアを使用しないかぎり、未使用のペアを削除することもできます。

条件規則の例

リソースに次のテキスト リンク分析規則が含まれ、TLA 結果の抽出を有効にしたとします。

Output columns:

Name:

Example:

Rule Value table:

	Element	Quantity	Example Token
1	mSupportNeg	Exactly 1	isn't
2		0 or 1	
3	(anything [(any a one) thing ?])	Exactly 1	anything
4		Between 0 and 2	that i
5	mNeg	Exactly 1	disliked
6	(about with in)	Exactly 1	about
7		0 or 1	
8	mDet	0 or 1	the

Rule Output table:

Concept 1	Type 1	Concept 2	Type 2	Concept 3	Type 3
product (9)	Products (9)	no dislike (5)	Positive		

Show output as: References to row in Rule Value table Specific token from example

図 42. 「テキスト リンク規則」タブ: 条件規則エディター

抽出すると、抽出エンジンが各文を読み込み、次の順序に合致させようとしています。

表 44. 抽出シーケンスの例

要素 (行)	引数の説明
1	マクロ mPos または mNeg、タイプ <不確定> で示されるいずれかのタイプのコンセプト。
2	マクロ mTopic で示されるいずれかのタイプに指定されたコンセプト。
3	マクロ mBe で示されるいずれかの語。
4	オプションの要素、0 個または 1 つの語で、単語の空所または <任意のトークン> として参照。
5	マクロ mTopic で示されるいずれかのタイプに指定されたコンセプト。

出力テーブルには、この規則の必要なすべてが「条件規則値テーブル」の行 5 で定義された mTopic マクロに対応するコンセプトまたはタイプ + 「条件規則値テーブル」の行 1 で定義された mPos、mNeg、または <不確定> に対応するコンセプトまたはタイプのパターンであることを示します。sausage + like または <不明> + <肯定的> となります。

条件規則の作成および編集

新しい条件規則を作成したり、既存のマクロを編集できます。条件規則エディターのガイドラインおよび説明に従います。詳しくは、223 ページの『テキスト リンク規則の使用』のトピックを参照してください。

条件規則の新規作成

1. メニューの「ツール」>「新規ルール」を選択します。また、ツリー・ツールバーの「新規条件規則」アイコンをクリックすると、エディターに新しい条件規則が表示されます。

2. 一意の名前を入力して、条件規則値の要素を定義します。
3. エラーの有無を確認したら、「適用」 をクリックします。

条件規則の編集

1. ツリー内の条件規則名をクリックします。右側のエディター・パネルに条件規則が表示されます。
2. 変更を加えます。
3. エラーの有無を確認したら、「適用」 をクリックします。

条件規則の無効化および削除

条件規則の無効化

処理中に条件規則が無視されるようにするには、マクロを無効にすることができます。条件規則を削除および無効にする場合は注意してください。

1. ツリー内の条件規則名をクリックします。右側のエディター・パネルに条件規則が表示されます。
2. 名前を右クリックします。
3. コンテキスト・メニューから、「無効にする」 を選択します。条件規則のアイコンがグレーになり、条件規則自体を編集できなくなります。

条件規則の削除

条件規則を除外する場合、削除することができます。条件規則を削除および無効にする場合は注意してください。

1. ツリー内の条件規則名をクリックします。右側のエディター・パネルに条件規則が表示されます。
2. 名前を右クリックします。
3. コンテキスト・メニューから、「削除」 を選択します。リストから条件規則が消去されます。

エラーのチェック、保存およびキャンセル

条件規則の変更の適用

条件規則エディターの外をクリックまたは「適用」 をクリックすると、条件規則が自動的にスキャンされ、エラーの有無が確認されます。エラーが見つかったら、アプリケーションの別の部分に移る前に修正する必要があります。

ただし、あまり深刻でないエラーが検出された場合は、警告のみが表示されます。例えば、条件規則にタイプまたはマクロに対する不完全または参照されない定義が含まれている場合、警告メッセージが表示されず、「適用」 をクリックすると、未修正の警告により、左側パネルのツリーの条件規則名の左側に警告アイコンが表示されます。

条件規則を適用しても、条件規則が永続的に保存されるわけではありません。適用すると、検証プロセスで、エラーおよび警告の有無をチェックします。

インタラクティブ・ワークベンチ・セッション内のリソースの保存

1. インタラクティブ・ワークベンチ・セッションでリソースに行った変更を保存し、次回ストリーム実行時にそれらの変更を取得できるようにするには、次の手順を実行する必要があります。
 - モデル作成ノードを更新して、次回ストリーム実行時にこれらの同じリソースを取得できるようにします。詳しくは、78 ページの『モデル作成ノードの更新および保存』のトピックを参照してください。

ださい。その後ストリームを保存します。ストリームを保存するには、モデル作成ノードを更新した後、IBM SPSS Modeler のメイン・ウィンドウで保存します。

2. インタラクティブ・ワークベンチ・セッションでリソースに行った変更を保存し、他のストリームでそれらを使用できるようにするは、次の手順を実行します。
 - 使用していたテンプレートを更新するか、新しいテンプレートを作成します。詳しくは、164 ページの『テンプレートの作成および更新』のトピックを参照してください。現在のノードの変更を保存するわけではありません (前の手順を参照)。
 - または、使用していた TAP を更新します。詳しくは、138 ページの『テキスト分析パッケージの更新』のトピックを参照してください。

テンプレート・エディター 内のリソースの保存

1. まず、ライブラリーを公開します。詳しくは、186 ページの『ライブラリーの公開』のトピックを参照してください。
2. そして、メニューで「ファイル」>「リソース・テンプレートを保存」を使用してテンプレートを保存します。

条件規則変更のキャンセル

1. 変更を破棄する場合は、エディター・パネルで「キャンセル」をクリックします。

条件規則の処理順序

テキスト リンク分析が抽出時に実行されると、合致が見つかるまで、またはすべての条件規則が終了するまで、「文」(句、語、語句)が各条件規則に対して順番にマッチされます。ツリー内の位置は、条件規則が使用される順序を示します。条件規則の順序は具体的なものから一般的なものへの順序がベスト・プラクティスです。最も具体的な条件規則がツリーの最上位になります。特定の条件規則またはルール・セットを変更するには、条件規則とマクロ ツリーのコンテキスト・メニューで「1つ上に移動」または「1つ下に移動」を選択するか、ツールバーの上方向矢印または下方向矢印を選択します。

ソース・ビュー内で作業している場合、エディター内を移動することによって条件規則の順序を変更することはできません。ソース・ビュー内で上位に表示されている条件規則から処理されます。コピー/貼り付けの問題を回避するために、条件規則の順序の変更はツリー内でのみ行うことを強くお勧めします。

重要: 以前のバージョンの IBM SPSS Modeler Text Analytics では、一意で数値の条件規則 ID が必要でした。バージョン 18.1.1 以降では、条件規則をツリー内の上下に移動して、あるいはソース・ビュー内の位置によってのみ処理順を示します。

例えば、次の 2 つの文を含むテキストがあるとします。

I love anchovies

I love anchovies and green peppers

また、値が次のようになる 2 つのテキスト リンク分析規則があるとします。

A			
	Element	Quantity	Example Token
1	Positive	Exactly 1	
2	mDet	0 or 1	
3	mTopic	Exactly 1	
4			
5			
6			
7			

B			
	Element	Quantity	Example Token
1	Positive	Exactly 1	
2	mDet	0 or 1	
3	mTopic	Exactly 1	
4	(SEP and or)	1 or 2	
5	mDet	0 or 1	
6	mTopic	Exactly 1	
7			

図 43. 2 つの条件規則の例

ソース・ビューで、条件規則値が次のようになります。

A: value = \$Positive \$mDet? \$mTopic

B: value = \$Positive \$mDet? \$mTopic (\$SEP|and|or){1,2} \$mDet? \$mTopic

ルール **A** が **B** よりツリー上で高い位置（頂点近く）にある場合は、**A** が先に処理され、*I love anchovies and green peppers* の分は最初に `$Positive $mDet?$mTopic` とマッチングされ、完全なパターン出力（anchovies + like）を生成します。これは、2つの `$mTopic` 一致を検索しないというルールによるものです。

そのため、テキストの核心をキャプチャーするには、最も具体的な規則（この場合 **B**）がより一般的な規則（この場合 **A**）よりツリー内の上位になければいけません。

ルール・セットの使用 (多段階処理)

ルール・セットとは、多段階処理を実行するために条件規則とマクロ ツリーの関連する一連の条件規則をグループ化する有用な方法です。ルール・セットには、名前以外にそれ自体の定義はありませんが、ルール・セットを使用して、条件規則を意味のあるグループに編成します。一度の通過で処理するにはテキストが多すぎ、また多様である場合があります。例えば、セキュリティーに関する情報データを処理する場合、連絡方法 (*x* が *y* に電話)、家族関係 (*y* の義理の兄弟 *x*)、金銭のやり取り (*x wired \$100 to y*) など明らかに個人間の繋がりがテキストに含まれている場合があります。この場合、連絡先が明らかになる定義、家族構成が明らかになる定義など、特定の種類の関係性に焦点を当てたテキスト リンク分析規則の特殊なセットを作成すると役立ちます。

ルール・セットを作成するには、条件規則とマクロ ツリーのコンテキスト・メニューまたはツールバーで「ルール・セットを作成」を選択します。ツリーのルール・セット・ノードで直接新しい条件規則を作成するか、既存の規則をルール・セットに移動できます。

条件規則がルール・セットにグループ化されるリソースを使用して抽出を実行すると、抽出エンジンは、テキスト全体で複数の通過を行い、通過ごとにさまざまな種類のパターンを合致させます。このように、「文」を各ルール・セットの条件規則に合致させることができます。ルール・セットがない場合は単一の条件規則にのみ合致させることができます。

注: ルール・セットごとに条件規則を、512 個まで定義することができます。

ルール・セットの新規作成

1. メニューの「ツール」>「新規ルール・セット」を選択します。または、ツリーのツールバーで「新規ルール・セット」アイコンをクリックします。条件規則ツリーにルール・セットが表示されます。
2. このルール・セットの新しい条件規則を追加するか、既存の条件規則をセット内に移動します。

ルール・セットの無効化

1. ツリー内のルール・セット名をクリックします。
2. コンテキスト・メニューから、「無効にする」を選択します。ルール・セットのアイコンがグレーになり、そのルール・セットに含まれるすべての条件規則も無効となって、処理時に無視されるようになります。

ルール・セットの削除

1. ツリー内のルール・セット名をクリックします。
2. コンテキスト・メニューから、「削除」を選択します。ルール・セットと、それに含まれるすべての条件規則がリソースから削除されます。

条件規則およびマクロにサポートされている要素

次の引数は、テキスト リンク分析規則およびマクロの値パラメーターに受け入れられます。

マクロ

テキスト リンク分析規則または別のマクロで直接マクロを使用できます。コンテキスト・メニューからマクロ名を選択するのではなく、手動でまたはソース・ビューからマクロ名を入力する場合、`$mTopic` のように、名前の前にドル記号 (\$) を必ず追加してください。マクロ名では大文字と小文字が区別されます。コンテキスト・メニューでマクロを選択する場合、現在の「テキスト リンク規則」タブで定義されたマクロから選択できます。

データ型

テキスト リンク分析規則またはマクロのタイプを直接使用できます。コンテキスト・メニューからタイプ名を選択するのではなく、手動でまたはソース・ビューからタイプ名を入力する場合、`$Person` のように、名前の前にドル記号 (\$) を必ず追加してください。タイプ名では大文字と小文字が区別されます。コンテキスト・メニューを使用する場合、使用されているリソースの現在のセットからタイプを選択できます。

不明のタイプを参照する場合、警告メッセージが表示され、修正されるまで条件規則とマクロ ツリーの条件に警告メッセージが表示されます。

リテラル文字列

抽出されなかった情報を追加するために、抽出エンジンが検索するリテラル文字列を定義できます。抽出されたすべての単語または句はタイプに割り当てられるため、それらをリテラル文字列で使用することはできません。抽出された単語を使用すると、そのタイプが <不明> であっても、その単語は無視されます。

リテラル文字列は、1 単語または複数の単語の場合があります。次の規則は、リテラル文字列のリストを定義する場合に適用されます。

- (his) のように、文字列のリストをカッコで囲む。リテラル文字列を選択する場合、各文字列を (a|an|the) または (his|hers|its) のように OR 演算子で区切る必要があります。

- 単語または複合語を使用する。
- リスト内の単語を、ブール型演算子 OR と同様に機能する | 文字で区切る。
- 単数形および複数形に一致させたい場合は両方を入力する。活用形は自動的に生成されません。
- 小文字の実を使用する。
- リテラル文字列を再利用するには、マクロとしてそれらを定義してから他のマクロおよびテキスト リンク分析規則でそのマクロを使用する。
- 文字列にピリオド (終止符) またはハイフンを使用している場合は、それらを含める必要がある。例えば、テキストの a.k.a を一致させるには、文字 a.k.a と共にピリオドをリテラル文字列として入力してください。

排他演算子




排他演算子として ! を使用し、否定の式が特定のスロットに含まれないようにします。インライン・セル編集 (条件規則値テーブルまたはマクロ値テーブルのセルをダブルクリック) またはソース・ビューから手動でのみ排他演算子を追加できます。例えば、`$mTopic @{0,2} !($Positive) $Budget` をテキスト リンク分析に追加すると、(1) `mTopic` マクロのいずれかのタイプに割り当てられたキーワードを含み、(2) 0 から 2 語の単語の空所を含み、(3) `<Positive>` タイプに割り当てられたキーワードのインスタンスを含まず、(4) `<Budget>` タイプに割り当てられたキーワードを含むテキストを検索します。「*cars have an inflated price tag*」はキャプチャされますが、「*store offers amazing discounts*」は無視されます。

この演算子を使用するには、セルをダブルクリックして、要素のセルに感嘆符のポイントおよびカッコを入力する必要があります。

単語の空所 (<任意のトークン>)

単語の空所 (または <任意のトークン>) は、2 つの要素間に存在するトークンの数値範囲を定義します。単語の空所は、追加の決定詞、前置詞句、形容詞またはその他の単語の有無によってわずかに異なる非常に類似した句と合致させる場合に役立ちます。

表 45. 単語の空所のない条件規則値テーブルの要素の例

#	エレメント
1	 Unknown
2	 mBeHave
3	 Positive



注: ソースビューではこの値は以下のように定義されます: `$Unknown $mBeHave $Positive`

この値は「*the hotel staff was nice*」のような文と合致します。ここで、「*hotel staff*」のタイプは <不明> となり、「*was*」はマクロ `mBeHave`、「*nice*」は <肯定的> となります。ただし、「*the hotel staff was very nice*」とは一致しません。

表 46. <任意のトークン> のある条件規則値テーブルの要素の例

#	エレメント
---	-------

表 46. <任意のトークン> のある条件規則値テーブルの要素の例 (続き)

1	 Unknown
2	 mBeHave
3	
4	 Positive

注: ソースビューではこの値は以下のように定義されます: \$Unknown \$mBeHave @{0,1} \$Positive

単語の空所を条件規則値に追加すると、「the hotel staff was nice」および「the hotel staff was very nice」のいずれの文にも合致します。

ソース・ビューまたはインライン編集の場合、単語の空所のシンタックスは @{#,#} です。ここで、@ は単語の空所を示し、{#,#} は、先行する要素と後続の要素の間に受け入れられる単語数の最小値および最大値を定義します。例えば、@{1,3} は、定義された 2 つの要素の間に少なくとも 1 から 3 個の単語がある場合、この 2 つの要素の間に合致がある可能性があることを示します。@{0,3} は 0、1、2 または 3 この単語がある場合に定義された 2 つの要素の間の合致があると考えられることを示します。

入力モードでの表示および作業

各条件規則およびマクロについて、TLA エディターは TLA 出力を合致および作成する抽出機能で使用する基底のソース・コードを生成します。コード自体を処理する場合、エディター上部の「View Source」ボタンをクリックして、このソースコードを表示したり直接編集したりすることができます。ソース・ビューが、現在選択している条件規則またはマクロにジャンプして、強調表示します。ただし、エディター・パネルを使用してエラーの発生を少なくすることをお勧めします。

ソースの表示または編集を終了するには、「ソースを終了」をクリックします。条件規則に無効なシンタックスを生成した場合、ソース・ビューを終了する前にそのシンタックスを修正する必要があります。

重要: ソース・ビューで編集する場合、一度に 1 つずつ条件規則およびマクロを編集することを強くお勧めします。マクロを編集した後、抽出して結果を検証してください。結果が適切である場合、他の変更を行う前にテンプレートを保存することをお勧めします。結果が適切でない場合、またはエラーが発生した場合、保存したリソースに戻ってください。

ソース・ビューのマクロ

```
[macro]
name = macro_name
value = ([type_name|macro_name|literal_string|word_gap])
```

表 47. マクロ・エントリ

[macro]	各マクロは、「macro」とマークされた行から開始し、マクロの開始を示す必要があります。
name	マクロ定義の名前。それぞれの名前は一意である必要があります。

表 47. マクロ・エンタリー (続き)

value	1 つまたは複数のタイプ、リテラル文字列、単語の空所またはマクロの組み合わせ。詳しくは、230 ページの『条件規則およびマクロにサポートされている要素』のトピックを参照してください。引数を結合する場合、カッコ () を使用して引数およびブール型演算子 OR を示す文字 をグループ化する必要があります。
-------	--

マクロのセクションで説明されているガイドラインおよびシンタックスのほか、ソース・ビューにはエディター・ビューでの作業に不要のガイドラインがいくつかあります。入力モードで作業する場合、マクロは次のことを重点に置く必要があります。

- 各マクロは、[macro] とマークされた行から開始し、マクロの開始を示す必要があります。
- 要素を無効にするには、行の前にコメントを示すインジケータ (#) を追加します。

例: mTopic の値は、次のうちの 1 つに一致するキーワードの存在を示します。mTopic の値は、以下のタイプのいずれかひとつに一致する語の存在を示します。 <Product>、<Person>、<Location>、<Organization>、<Budget>、<Unknown>。

```
[macro]
name=mTopic
value=($Unknown|$Product|$Person|$Location|$Organization|$Budget|$Currency)
```

ソース・ビューの条件規則

```
[pattern(ID)]
name = pattern_name
value = [$type_name|macro_name|word_gaps|literal_strings]
output = $digit[¥t]#digit[¥t]$digit[¥t]#digit[¥t]$digit[¥t]#digit[¥t]
```

表 48. 規則エンタリー

[pattern (<ID>)]	テキスト リンク分析規則の開始を示し、処理の順序を決定する一意の数値 ID を与えます。
name	テキスト リンク分析規則の一意の名前を示します。
value	テキストに一致するシンタックスおよび引数を提供します。詳しくは、230 ページの『条件規則およびマクロにサポートされている要素』のトピックを参照してください。
出力	<p>テキストで見つかった合致パターンの出力形式。出力は、ソース・テキストの要素の正確な元の位置と必ずしも似ているわけではありません。また、各行に出力を配置することで、指定されたテキスト リンク分析規則に複数の出力行を指定することができます。</p> <p>出力のシンタックス:</p> <ul style="list-style-type: none"> • \$1¥t#1¥t\$3¥t#3 のように、出力をタブ・コード ¥t で区切る。 • \$ および数値は、該当する位置の値パラメーターで定義された引数に一致するキーワードを呼び出します。つまり、\$1 は値に定義された最初の引数に一致するキーワードを意味します。 • # および数値は、該当する位置の要素のタイプ名を呼び出します。項目がリテラル文字列のリストである場合、タイプ <不明> が割り当てられます。 • Null¥tNull の値は、出力を作成しません。

条件規則のセクションで説明されているガイドラインおよびシンタックスのほか、ソース・ビューにはエディター・ビューでの作業に不要のガイドラインがいくつかあります。入力モードで作業する場合、条件規則は次のことを重点に置く必要があります。

- 複数の要素が定義されている場合、オプションであるかどうかに関係なく、カッコで囲む必要があります (例: (\$Negative|\$Positive) または (\$mCoord|\$SEP)?)。\$SEP はコンマを意味します。

- テキスト リンク分析規則の最初の要素を、オプションの要素にすることはできません。例えば、value = \$mTopic? または value = @{0,1} から始めることはできません。
- 数量 (またはインスタンス数) をトークンに関連付けることができます。これは、ケースごとに個別の規則を作成するのではなく、すべてのケースを網羅する規則を 1 つだけ作成する場合に役立ちます。例えば、, (コンマ) または and のいずれかに合致させる場合、リテラル文字列 (\$SEP|and) を使用できません。数量を追加してリテラル文字列が (\$SEP|and){1,2} となるよう拡張すると、以下のいずれかに一致するようになります。", " "および" ", and".
- テキスト リンク分析規則 value では、マクロ名と \$ 文字および ? 文字との間にスペースを使用することはできません。
- テキスト リンク分析規則 output にスペースは使用できません。
- 要素を無効にするには、行の前にコメントを示すインジケータ (#) を追加します。

例: リソースに次の TLA テキスト リンク分析規則が含まれ、TLA 結果の抽出を有効にしたとします。

```
## Jean Doe was the former HR director of IBM in France
[pattern(201)]
name= 1_201
value = $Person ($SEP|$mDet|$mSupport|as|then){1,2} @{0,1} $Function
(of|with|for|in|to|at) @{0,1} $Organization @{0,2} $Location
output = $1¥t#1¥t$4¥t#4¥t$7¥t#7¥t$9¥t#9
```

抽出すると、抽出エンジンが各文を読み込み、次の順序に合致させようとしています。

表 49. 抽出シーケンスの例

位置	引数の説明
1	個人の名前 (\$Person)
2	以下のうち1つまたは2つとします: コンマ (\$SEP)、区切り文字 (\$mDet)、助動詞 (\$mSupport)、文字列「then」または「as」
3	0 または 1 つの単語 (@{0,1})
4	役職 (\$Function)
5	以下のいずれかの文字列。 "of"、"with"、"for"、"in"、"to"、"at"、
6	0 または 1 つの単語 (@{0,1})
7	組織の名前 (\$Organization)
8	0, 1 または 2 つの単語 (@{0,2})
9	場所の名前 (\$Location)

このテキスト リンク分析規則のサンプルは、次のような文または句に合致します。

Jean Doe, the HR director of IBM in France

Jean Doe was the former HR director of IBM in France

IBM appointed Jean Doe as the HR director of IBM in France

このテキスト リンク分析規則のサンプルは、次のような出力を作成します。

jean doe <Person> hr director <Function> ibm <Organization> france <Location>

この場合、次のようになります。

- jean doe は、\$1 (テキスト リンク分析規則の最初の要素) に対応するキーワードで、<Person> は jean doe (#1) のタイプです。
- director は、\$4 (テキスト リンク分析規則の 4 番目の要素) に対応するキーワードで、<Function> は hr director (#4) のタイプです。
- ibm は、\$7 (テキスト リンク分析規則の 7 番目の要素) に対応するキーワードで、<Organization> は ibmの以下となります。 (#7),
- france は、\$9 (テキスト リンク分析規則の 9 番目の要素) に対応するキーワードで、<Location> は france (#9) のタイプです。

ソース・ビューのルール・セット

[set(<ID>)]

ここで、[set (<ID>)] ルール・セットの開始を示し、そのセットの処理の順序を決定する一意の数値 ID を与えます。

例: 次の文には、人名、会社内の役職、その会社の合併/買収の活動が含まれています。

Org1 Inc has entered into a definitive merger agreement with Org2 Ltd, said John Doe, CEO of Org2 Ltd.

いくつかの出力で条件規則を作成し、次のような出力をすべて処理できます。

```
## Org1 Inc entered into a definitive merger agreement with Org2 Ltd, said John Doe, CEO of Org2 Ltd.
```

```
[pattern(020)]
name=020
value = $Organization @{0,4} $ActionNouns @{0,6} $mOrg @{1,2}
       $Person @{0,2} $Function @{0,1} $Organization
output = $1¥t#1¥t$3¥t#3¥t$5¥t#5
output = $7¥t#7¥t$9¥t#9¥t$11¥t#11
```

この場合、次の 2 つの出力パターンが出力されます。

- org1 inc<Organization> + merges with <ActiveVerb> + org2 ltd<Organization>
- john doe <Person> + ceo <Function> + org2 ltd<Organization>

重要: TLA パターンの間、他の言語処理が実行されます。この場合、merger は、抽出プロセスの類義語グループ段階で merges with の下にグループ化されます。merges with は <ActiveVerb> タイプに含まれるため、このタイプ名は最終 TLA パターン出力に表示されます。そのため、出力が t\$3¥t#3 を読み込む場合、パターンは 3 番目の要素の最終コンセプトと 3 番目の最終タイプを、言語処理が適用された後 (類義語およびその他のグループ) に表示します。

前述のような複雑な条件規則を作成する代わりに、2 つの条件規則を使用して容易に管理できます。1 つ目は、会社間で合併者/買収者を検出します。

```
[set(1)]
## Org1 Inc has entered into a definitive merger agreement with Org2 Ltd
[pattern(44)]
name=firm + action + firm_0044
value=$mOrg @{0,20} $ActionNouns @{0,6} $mOrg
output(1)=$1¥t#1¥t$3¥t#3¥t$5¥t#5
```

これにより、org1 inc<Organization> + merges with <ActiveVerb> + org2 ltd <Organization> が生成されます。

2 つ目は人名/役職/会社に特化されています。

```
[set(2)]
## said John Doe, CEO of Org2 Ltd
[pattern(52)]
name=individual + role + firm_0007
value=$Person @{0,3} $mFunction (at|of)? ($mOrg|$Media|$Unknown)
output(1)=$1¥t#1¥t$3¥tFunction¥t$5¥t#5
```

これにより、john doe <Person> + ceo <Function> + org2 ltd <Organization> が生成されます。

特記事項

本情報は全世界で提供する製品およびサービスについて作成したものです。

本書に記載の製品、サービス、または機能が日本においては提供されていない場合があります。日本で利用可能な製品、サービス、および機能については、日本 IBM の営業担当員にお尋ねください。本書で IBM 製品、プログラム、またはサービスに言及していても、その IBM 製品、プログラム、またはサービスのみが使用可能であることを意味するものではありません。これらに代えて、IBM の知的所有権を侵害することのない、機能的に同等の製品、プログラム、またはサービスを使用することができます。ただし、IBM 以外の製品とプログラムの操作またはサービスの評価および検証は、お客様の責任で行っていただきます。

IBM は、本書に記載されている内容に関して特許権 (特許出願中のものを含む) を保有している場合があります。本書の提供は、お客様にこれらの特許権について実施権を許諾することを意味するものではありません。実施権についてのお問い合わせは、書面にて下記宛先にお送りください。

〒103-8510

東京都中央区日本橋箱崎町19番21号

日本アイ・ビー・エム株式会社

法務・知的財産

知的財産権ライセンス渉外

IBM およびその直接または間接の子会社は、本書を特定物として現存するままの状態を提供し、商品性の保証、特定目的適合性の保証および法律上の瑕疵担保責任を含むすべての明示もしくは黙示の保証責任を負わないものとします。国または地域によっては、法律の強行規定により、保証責任の制限が禁じられる場合、強行規定の制限を受けるものとします。

この情報には、技術的に不適切な記述や誤植を含む場合があります。本書は定期的に見直され、必要な変更は本書の次版に組み込まれます。IBM は予告なしに、随時、この文書に記載されている製品またはプログラムに対して、改良または変更を行うことがあります。

本書において IBM 以外の Web サイトに言及している場合がありますが、便宜のため記載しただけであり、決してそれらの Web サイトを推奨するものではありません。それらの Web サイトにある資料は、この IBM 製品の資料の一部ではありません。それらの Web サイトは、お客様の責任でご使用ください。

IBM は、お客様が提供するいかなる情報も、お客様に対してなんら義務も負うことのない、自ら適切と信ずる方法で、使用もしくは配布することができるものとします。

本プログラムのライセンス保持者で、(i) 独自に作成したプログラムとその他のプログラム (本プログラムを含む) との間での情報交換、および (ii) 交換された情報の相互利用を可能にすることを目的として、本プログラムに関する情報を必要とする方は、下記に連絡してください。

IBM Director of Licensing

IBM Corporation

North Castle Drive, MD-NC119

Armonk, NY 10504-1785

US

本プログラムに関する上記の情報は、適切な使用条件の下で使用することができますが、有償の場合もあります。

本書で説明されているライセンス・プログラムまたはその他のライセンス資料は、IBM 所定のプログラム契約の契約条項、IBM プログラムのご使用条件、またはそれと同等の条項に基づいて、IBM より提供されます。

記載されている性能データとお客様事例は、例として示す目的でのみ提供されています。実際の結果は特定の構成や稼働条件によって異なります。

IBM 以外の製品に関する情報は、その製品の供給者、出版物、もしくはその他の公に利用可能なソースから入手したものです。IBM は、それらの製品のテストは行っておりません。したがって、他社製品に関する実行性、互換性、またはその他の要求については確認できません。IBM 以外の製品の性能に関する質問は、それらの製品の供給者にお願いします。

IBM の将来の方向性および指針に関する記述は、予告なく変更または撤回される場合があります。これらは目標および目的を提示するものにすぎません。

本書には、日常の業務処理で用いられるデータや報告書の例が含まれています。より具体性を与えるために、それらの例には、個人、企業、ブランド、あるいは製品などの名前が含まれている場合があります。これらの名称はすべて架空のものであり、類似する個人や企業が実在しているとしても、それは偶然にすぎません。

商標

IBM、IBM ロゴおよび `ibm.com` は、世界の多くの国で登録された International Business Machines Corporation の商標です。他の製品名およびサービス名等は、それぞれ IBM または各社の商標である場合があります。現時点での IBM の商標リストについては、<http://www.ibm.com/legal/copytrade.shtml> をご覧ください。

インテル、Intel、Intel ロゴ、Intel Inside、Intel Inside ロゴ、Centrino、Intel Centrino ロゴ、Celeron、Xeon、Intel SpeedStep、Itanium、および Pentium は、Intel Corporation または子会社の米国およびその他の国における商標または登録商標です。

Linux は、Linus Torvalds の米国およびその他の国における登録商標です。

Microsoft、Windows、Windows NT および Windows ロゴは、Microsoft Corporation の米国およびその他の国における商標です。

UNIX は The Open Group の米国およびその他の国における登録商標です。

Java およびすべての Java 関連の商標およびロゴは Oracle やその関連会社の米国およびその他の国における商標または登録商標です。

索引

日本語, 数字, 英字, 特殊文字の順に配列されています。なお, 濁音と半濁音は清音と同等に扱われています。

[ア行]

アスタリスク (*)
 不要語辞書 200
 類義語 198
アップグレード 1
アミノ酸 (固有表現) 206
アンチリンク 110
意見ライブラリー 190
移動
 カテゴリー 140
 キーワード辞書 196
インタラクティブ・ワークベンチ 24, 25, 27, 67, 78
インタラクティブ・ワークベンチにおけるビュー
 カテゴリーとコンセプト 68, 95
 クラスター 70
 テキスト リンク分析 72
 リソース・エディター 74
インタラクティブ・ワークベンチを起動 24
インデント形式 134
エクスポート
 定義済みカテゴリー 135
 テンプレート 175
 パブリック・ライブラリー 183
オプション 76
 サウンド・オプション 77
 セッション・オプション 76
 表示オプション (色) 76
オプションの要素 197
 エントリーの削除 199
 対象 199
 追加 199
 の定義 197

[カ行]

回答およびカテゴリーの関連性 105
外部リンク 143
拡張リソース 203
 エディターの検索と置換 204
活用形 111, 189, 191, 192
活用形の生成 189, 191, 192

カテゴリー 19, 95, 96, 103, 139
 移動 140
 拡張 110, 116
 方略 99
 関連性 105
 記述子 99, 100, 103
 空白カテゴリーを新規作成 119
 結果の調整 139
 結合 141
 削除 141
 作成 98, 106, 108, 110, 115, 116, 120
 手動作成 119
 スコアリング 96
 注釈 104
 テキスト マイニング カテゴリー モデル ナゲット 26
 テキスト分析パッケージ 136, 138
 名前 104
 名前変更 119
 フラット化 140
 プロパティ 104
 への追加 139
 編集 139
 labels 104
カテゴリー モデル ナゲット 19, 40
 出力 41
 生成 78
 設定タブ 42
 ノードで作成 26
 「フィールド」タブ 43
 フィールドまたはレコードとしてのコンセプト 42
 「モデル」タブ 41
 「要約」タブ 43
 例 44
 ワークベンチで作成 25
カテゴリー Web グラフ/テーブル 156, 157
カテゴリー化 7, 95
 共起規則 108, 110, 114
 グループ手法の使用 108
 言語学的手法 106, 116
 出現頻度に基づく手法 115
 手動で 119
 手法の使用 110
 セマンティック・ネットワーク 108, 110, 113
 内包関係のコンセプト 108, 110, 112
 派生関係のコンセプトの語幹 108, 110, 111
 方法 98

カテゴリー規則 121, 127, 129, 130
 共起規則 108, 110, 116
 構文 121
 コンセプト共起から 108, 110, 114, 116
 類義語から 108, 110, 116
 例 127
カテゴリー結合 141
カテゴリーとコンセプト・ビュー 68, 95
 カテゴリー・ペイン 96
 データ・ペイン 104
カテゴリーの拡張 116
カテゴリーの結合 141
カテゴリーの作成 7, 106, 108
 共起規則手法 116
 セマンティック・ネットワーク手法 116
 内包関係のコンセプトの手法 116
 派生関係のコンセプトの語幹による手法 116
 分類リンクの例外 110
カテゴリーのフラット化 140
カテゴリーのラベル 104
カテゴリー棒グラフ 156
カテゴリー名 96
カテゴリー・ペイン 96
カテゴリー・ペインの列の表示 96
画面読み上げソフトウェア 79, 80
カレット記号 (^) 198
感嘆符 (!) 198
管理
 カテゴリー 139
 パブリック・ライブラリー 183
 ローカル・ライブラリー 182
キーボードのショートカットのナビゲート 79
キーボード・ショートカット 79, 80
キーワードおよびタイプの検索 181
キーワード辞書 179
 移動 196
 オプションの要素 189
 キーワードの強制 195
 項の追加 192
 削除 196
 使用不可化 196
 タイプの作成 191
 名前変更 195
 ビルトインのタイプ 190
 類義語 189
キーワードのコンポーネント化 111
記述子 96

記述子 (続き)

カテゴリー 99, 103

カテゴリーの編集 139

クラスター 147

最適の選択 100

規則の演算子 & | !() 129

技法

共起規則 108, 110, 114, 116

セマンティック・ネットワーク 108,
110, 113, 116

ドラッグ・アンド・ドロップ 120

内包関係のコンセプト 108, 110, 112,
116

派生関係のコンセプトの語幹 108, 110,
111, 116

頻度 115

基本キーワード 34

キャッシュ

データおよびセッション抽出結果 25

Web フィード 14

共起規則手法 108, 110, 114, 116

強制

コンセプト抽出 93

用語 195

強制定義 211, 213

区切り文字 76

クラスター 25, 70, 143

概要 143

記述子 147

クラスター Web グラフ 157, 158

検討 146

コンセプト Web グラフ 157

作成 144

類似度リンク値 145

「クラスター」ビュー 70

クラスターのリンク 143

グラフ 158, 159

クラスター Web グラフ 157, 158

コンセプト Web グラフ 157

コンセプト・マップ 86

タイプ Web グラフ 158, 159

探索的分析モード 159

編集 159

TLA コンセプト Web グラフ 158,
159

グローバルな区切り文字 76

計量 (固有表現) 206

結果の調整

カテゴリー 139

コンセプト抽出を強制 93

コンセプトの除外 92

コンセプトのタイプへの追加 91

タイプの作成 91

抽出結果 89

類義語の追加 90

結果のフィルタリング 85, 152

言語学的手法 2

言語処理セクション 203, 211

強制定義 211, 213

省略形 211, 214

抽出パターン 211

言語ノード 11, 17, 60

スクリプトのプロパティ 60

設定タブ 18

言語リソース 48, 179

テキスト分析パッケージ 136, 138

テンプレート 163

リソース テンプレート 167

検索と置換 (拡張リソース) 204

コード・フレーム 130, 131

コア・ライブラリー 190

更新

テンプレート 164, 173

ノード・リソースおよびテンプレート
173

モデル作成ノード 78

ライブラリー 184, 186

固有表現

アミノ酸 206

計量 206

時間 206

住所 206

使用可能化および使用不可化 210

数値 206

正規化、NonLingNorm.ini 209

正規表現、RegExp.ini 207

通貨 206

電子メール・アドレス 206

電話番号 206

パーセンテージ 206

日付 206

日付形式 209

プロテイン 206

米国社会保障番号 206

HTTP アドレス/URL 206

IP アドレス 206

固有表現の無効化 210

固有表現の有効化 210

固有表現有効化 210

コンセプト 19, 32

カテゴリー内 99, 103

カテゴリーに追加 99, 103, 139

クラスター 147

コンセプト・マップ 86

最適な記述子 100

スコアリングするフィールドまたはレ
コードとして 34, 42

タイプの作成 89

タイプへの追加 91

抽出 81

抽出からの除外 92

抽出への強制投入 93

コンセプト (続き)

フィルター操作 85

コンセプト モデル ナゲット 19, 31

スコアリングのコンセプト 32

設定タブ 34

ノードで作成 26

「フィールド」タブ 35

フィールドまたはレコードとしてのコ
ンセプト 34

「モデル」タブ 32

「要約」タブ 36

類義語 34

例 36

コンセプト Web グラフ 157

コンセプトのマッピング 86

コンセプトの無視 92

コンセプト・パターン 151

コンセプト・マップ 86, 89

インデックスの作成 89

コンセプト・マップのインデックス 89

コンセプト・マップ・インデックスの作成
89

コンパクト形式 133

コンポーネント化 111

[サ行]

最小リンク値 108

再利用

データおよびセッション抽出結果 25

Web フィード 14

サウンドのミュート 77

サウンド・オプション 77

削除

オプションの要素 199

カテゴリー 141

カテゴリー規則 130

キーワード辞書 196

除外されたエントリー 200

ライブラリー 183

ライブラリーの無効化 183

リソース テンプレート 174

類義語 199

作成

オプションの要素 199

カテゴリー 2, 7, 26, 98, 106, 108,
110, 111, 112, 113, 114, 115, 116,
119, 120

カテゴリー規則 121, 129

キーワード辞書 191

規則のあるカテゴリー 121

クラスター 144

タイプ 91

テンプレート 173

不要語辞書のエントリー 200

作成 (続き)

モデル作成ノードおよびカテゴリ モデル ナゲット 78
ライブラリー 180
リソースからテンプレート 164
類義語 89, 90, 198
サンプル・ノード
テキスト マイニング 30
視覚化ペイン 155
クラスター Web グラフ 157, 158
コンセプト Web グラフ 157
タイプ Web グラフ 158, 159
テキスト リンク分析ビュー 158, 159
TLA コンセプト Web グラフ 158, 159
時間 (固有表現) 206
辞書 74, 189
タイプ 179, 189
不要語 179, 189, 200
類義語 179, 189, 197
社会保障番号 (固有表現) 206
住所 (固有表現) 206
ショートカット・キー 79, 80
使用不可化
キーワード辞書 196
固有表現 210
不要語辞書 200
ライブラリー 183
類義語辞書 199, 205
省略形 211, 214
除外
カテゴリ・リンク 110
辞書の無効化 196, 199
抽出からコンセプト 92
不要語エントリーの無効化 200
ライブラリーの無効化 183
fuzzy の除外 205
新規カテゴリ 119
数値 (固有表現) 206
「スコア」ボタン 96
スコアリング 96
コンセプト 33
スコアリングするコンセプトの選択 33
すべてのドキュメント 96
スペル ミス 205
正規化 209
生成するカテゴリ数の最大値 108
セッション情報 24, 25, 27
セッションの終了 78
設定 76, 77
セマンティック・ネットワーク手法 108, 110, 113, 116
ソース ノード
ファイル リスト 8, 11
Web フィールド 8, 13

[夕行]

対象言語 205
代表語 198
タイプ 189
エディターの検索 181
コンセプトの追加 89
作成 191
辞書 179
タイプの頻度 115
抽出 81
デフォルトの色 76, 191
ビルトインのタイプ 190
フィルター操作 85, 152
タイプ Web グラフ 158, 159
タイプの頻度 115
タイプ・パターン 151
多段階処理 229
単語の空所 230
探索的分析モード 159
注釈
カテゴリ 104
抽出 1, 2, 5, 49, 81, 82, 179, 189
結果の調整 89
単語の強制 93
抽出結果 81
データのパターン 47
ユニターム 5
TLA パターン 150
抽出の結果 81
結果のフィルタリング 85, 152
抽出パターン 211
追加
オプションの要素 199
カテゴリにコンセプト 139
キーワード辞書へのキーワード 192
キーワードを不要語辞書へ 200
記述子 100
サウンド 76, 77
タイプ 91
パブリック・ライブラリー 181
類義語 90, 198
通貨 (固有表現) 206
データ
カテゴリ化 95, 106, 119
カテゴリの作成 108, 110, 116
クラスタリング 143
結果の調整 89
結果のフィルタリング 85, 152
再構成 51
抽出 81, 82, 150
データ・ペイン 104, 153
テキスト リンク分析 149
テキスト リンク分析パターンの抽出 149
データ区分モード 21

データ・ペイン

カテゴリとコンセプト・ビュー 104
テキスト リンク分析ビュー 153
「表示」ボタン 96
データ・ペインの列の表示 153
定義 99, 103
定義済みカテゴリ 130, 131, 135
インデント形式 134
コンパクト形式 133
フラット・リスト形式 132
テキスト マイニング 2
テキスト マイニング モデル ナゲット 8
TMWBModelApplier のスクリプトの
プロパティ 63
テキスト マイニング モデル作成ノード
8, 19, 20, 59
新しいノードの生成 78
「エキスパート」タブ 28
更新 78
「フィールド」タブ 21
「モデル」タブ 24
例 30
TextMiningWorkbench のスクリプト
のプロパティ 61
テキスト マイニングの
.doc/.docx/.docm ファイル 12
テキスト マイニングの .htm/.html ファ
イル 12
テキスト マイニングの .pdf ファイル 12
テキスト マイニングの .ppt/.pptx/.pptm
ファイル 12
テキスト マイニングの .shtml ファイル
12
テキスト マイニングの .txt/.text ファ
イル 12
テキスト マイニングの .xls/.xlsx/.xslm
ファイル 12
テキスト マイニングの .xml ファイル 12
テキスト リンク分析 (TLA) 47, 72, 149,
151, 215, 216, 217, 218, 219, 223, 226,
227, 228, 232
開始ポイント 216
グラフの表示 158, 159
視覚化ペイン 158, 159
シミュレーション結果 217, 218
条件規則エディター 215
条件規則の処理順 228
条件規則の無効化および削除 227
多段階処理 229
ツリー上の警告 219
データ・ペイン 153
テキスト マイニング モデル作成ノ
ード 25
ナビゲートのルールとマクロ 219
入力モード 232
パターンの検証 149

テキスト リンク分析 (TLA) (続き)
パターンのフィルタリング 152
引数 230
編集が必要な場合 216
マクロ 220
マクロの編集および削除 215
ライブラリーの指定 215, 219
TLA ノード 47
Web グラフ 158, 159
テキスト リンク分析結果のシミュレーション 217, 218
データの定義 217
テキスト リンク分析ノード 8, 47, 48, 49, 51, 52, 64
「エキスパート」タブ 49
出力 51
スクリプトのプロパティ 64
データの再構成 51
「フィールド」タブ 48
例 52
TLA のキャッシュ 51
テキスト区切り文字 76
テキスト分析 2
テキスト分析パッケージ 136, 137, 138
ロード 137
テキスト・マッチ 104
デフォルトのライブラリー 179
電子メール (固有表現) 206
テンプレート 5, 47, 48, 74, 149, 163, 167
インポートとエクスポート 175
更新または名前を付けて保存 164
削除 174
テンプレートの切り替え 165
テンプレートを開く 172
名前変更 174
バックアップ 176
復元 176
保存 173
リソースから作成 164
「リソース・テンプレートを読み込む」ダイアログ・ボックス 27
TLA 165
テンプレートを開く 172
テンプレート・エディター 167, 168, 172, 173, 174, 175
インポートとエクスポート 175
エディターの終了 175
テンプレートの削除 174
テンプレートの名前変更 174
テンプレートの保存 173
テンプレートを開く 172
ノードのリソースの更新 173
リソース・ライブラリー 179
電話番号 (固有表現) 206
ドキュメント フィールド 55
ドキュメント列 96

ドラッグ・アンド・ドロップ 120
ドル記号 (\$) 198

[ナ行]

内部リンク 143
内包関係のコンセプトの手法 108, 110, 112, 116
名前
カテゴリー 104
キーワード辞書 195
ライブラリー 182
名前変更
カテゴリー 119
キーワード辞書 195
ライブラリー 182
リソース・テンプレート 174
ノード
カテゴリー モデル ナゲット 40
言語 17
コンセプト モデル ナゲット 31
テキスト マイニング ビューアー 8, 55
テキスト マイニング モデル ナゲット 8
テキスト マイニング モデル作成ノード 8, 20
テキスト リンク分析 8, 47
ファイル リスト 8, 11
Web フィード 8, 13
ノードおよびモデル ナゲットの生成 78

[ハ行]

パーセント (固有表現) 206
排他演算子 230
派生関係のコンセプトの語幹による手法 108, 110, 111, 116
パターン 25, 47, 81, 149, 151, 215, 219, 223
多段階処理 229
テキスト リンク規則エディター 215
引数 230
パブリッシュ 186
パブリック・ライブラリーの追加 181
ライブラリー 184
日付 (固有表現) 206, 209
日付形式
固有表現 209
ビューアー ノード 8, 55, 56
設定タブ 55
テキスト マイニング 55
例 56
表 80

表示
クラスター 157
テキスト リンク分析 158, 159
文書 55
ライブラリー 182
「表示」ボタン 96
表示設定 76
表題 55
品詞 211, 213
頻度 115
ブール型演算子 129
ファイル リスト ノード 8, 11, 12, 13
拡張子リスト 12
スクリプトのプロパティ 59
設定タブ 12
その他のタブ 13
例 13
ファイル リスト ノードの拡張子リスト 12
フォントの色 191
複数形 191
付属 (デフォルト) ライブラリー 179
不要語辞書 179, 200
フラット・リスト形式 132
プリファレンス 76, 77
プロテイン (固有表現) 206
プロパティ
カテゴリー 104
文書 104, 153
リスト 55
変更
テンプレート 165, 172
編集
カテゴリー 139
カテゴリー規則 130
抽出結果の調整 89
編集モード 159
保存
インタラクティブ・ワークベンチ 78
データおよびセッション抽出結果 25
テンプレート 173
テンプレートとしてのリソース 164
リソース 176
Web フィード 14

[マ行]

マクロ 220, 221, 222
mNonLingEntities 223
mTopic 223
マッチ・オプション 189, 191, 192
未カテゴリー化 96
モデル ナゲット 24
インタラクティブ・ワークベンチからの生成 78

モデル ナゲット (続き)
カテゴリー モデル ナゲット 19, 24,
26, 40, 41
コンセプト モデル ナゲット 19, 24,
26, 31, 32

[ヤ行]

ユーザー定義の色 76
用語
エディターの検索 181
活用形 189
キーワードの強制 195
タイプへの追加 192
不要語辞書に追加 200
マッチ・オプション 189
color 191
予算ライブラリー 190
呼び出し
定義済みカテゴリー 131
テンプレート 175
パブリック・ライブラリー 183

[ラ行]

ライブラリー 74, 179, 189
意見ライブラリー 190
エクスポート 183
共有および公開 184
コア・ライブラリー 190
更新 186
削除 183
作成 180
辞書 179
使用不可化 183
追加 181
同期 184
名前 182
名前変更 182
パブリック・ライブラリー 184
パブリッシュ 186
表示 182
付属のデフォルト・ライブラリー 179
予算ライブラリー 190
呼び出し 183
ライブラリー同期の警告 184
リンク 181
ローカル・ライブラリー 184
ライブラリーの共有 184
更新 186
パブリック・ライブラリーの追加 181
パブリッシュ 186
ライブラリーの同期 184, 186
ライブラリーのフィルタリング 182

リソース
拡張リソースを編集 203
テンプレート・リソースに切り替え
165
バックアップ 176
復元 176
付属のデフォルト・ライブラリー 179
リソース テンプレート 5, 47, 48, 74, 149
リソースからテンプレートを作成 164
リソースのバックアップ 176
リソースをテンプレートに置き換え 165
リソース・エディター 74, 163, 164, 165,
168, 203
テンプレートの更新 164
テンプレートの作成 164
リソースの切り替え 165
リソース・テンプレート 163, 167
リソース・テンプレートの読み込み 27,
48, 173
リテラル文字列 230
リンク値 145
リンクの例外 110
類義語 89, 197
エントリーの削除 199
コンセプト モデル ナゲット 34
代表語 198
追加 90, 198
の定義 197
colors 198
Fuzzy Grouping の例外 205
! ^ * \$ 記号 198
類義語辞書 179, 197, 198, 199
類似度リンク値 145
類似度リンク値の計算 145
レコード 104, 153
列の折り返し 76

[ワ行]

ワークベンチ 24, 25, 27

A

AND 演算子 129

B

Budget キーワード辞書 190

C

Clem 式ビルダー 80
colors
色オプションの設定 76
タイプおよびキーワード 191

colors (続き)
不要語辞書 200
類義語 198

D

delimiter 76

F

filelistnode スクリプトのプロパティ 59
Fuzzy Grouping の例外 203, 205

H

HTTP/URL (固有表現) 206

I

ID フィールド 48
IP アドレス (固有表現) 206

L

label
Web フィールドの再利用 14
language
リソースの対象言語の設定 205
languageidentifier プロパティ 60
Location キーワード辞書 190

M

Microsoft Excel .xls / .xlsx ファイル
定義済みカテゴリーのインポート 130,
131
定義済みカテゴリーのエクスポート
135
mNonLingEntities 223
mTopic 223

N

Negative キーワード辞書 190
NOT 演算子 129

O

OR 演算子 129
Organization キーワード辞書 190

P

Person キーワード辞書 190
Positive キーワード辞書 190
Product キーワード辞書 190

R

rules 226
 共起規則手法 114
 構文 121
 削除 130
 作成 129
 ブール型演算子 129
 編集 130

T

textlinkanalysis プロパティ 64
TextMiningWorkbench のスクリプトのプロ
 パティ 61
TLA 165
TLA コンセプト Web グラフ 158, 159
TMWBModelApplier スクリプトのプロパ
 ティ 63

U

Uncertain キーワード辞書 190
Unknown キーワード辞書 190
URL 14, 15

W

Web グラフ
 クラスター Web グラフ 157, 158
 コンセプト Web グラフ 157
 タイプ Web グラフ 158, 159
 TLA コンセプト Web グラフ 158,
 159
Web フィード ノード 8, 11, 13, 14, 15,
 59
 キャッシュのラベルおよび再利用 14
 「コンテンツ」タブ 16
 スクリプトのプロパティ 59
 「入力」タブ 14
 例 17
 「レコード」タブ 15
Web フィードの HTML 形式 13, 15
Web フィードの RSS 形式 13, 15
webfeednode プロパティ 59

[特殊文字]

! 類義語の ^ * \$ 記号 198
& | !() 規則の演算子 129
*.lib 183
*.tap テキスト分析パッケージ 136, 137,
 138
.テキスト マイニングの rtf ファイル 12



Printed in Japan