Chapter 1 デー	ータマイニングガイド	1-1
1.1.1	データマイニングの概要	1-2
1.2.1	CRISP-DMプロセスモデル	1-3
1.3.1	ビジネスの理解	1-4
1.3.2	データの理解	1-6
1.3.3	データの準備	1-8
1.3.4	モデルの作成	1-10
1.3.5	評価	1-12
1.3.6	導入	1-14
1.4.1	IBM SPSS Modelerの起動	1-16
1.4.2	パレットとノード	1-18
1.4.3	ストリーム	1-19
1.4.4	マネージャウィンドウ	1-20
~ -		
Chapter 2 7-	ータのインポート	2-1
2.1.1	データのインポートの概要	2-2
2.2.1	Microsoft Excel形式のデータファイルの読込み	2-3
2.3.1	CSV形式のデータファイルの読込み	2-14
2.4.1	IBM SPSS Statistics形式のデータファイルの読込み	2-23
2.5.1	データベースの読込み - データソースの定義	2-30
2.5.2	データベースの読込み ーデータのインポートー	2-34
Chapter 3 フィ	ィールドのデタ型	3-1
3.1.1	フィールドのデータ型	3-2
3.1.2	ストリームの確認	3-4
3.2.1	フィールドの尺度の設定	3-5
3.2.2	フィールドのインスタンス化	3-10
3.3.1	フィールドのロールの設定	3-12
3.4.1	欠損値の設定	3-14
3.4.2	ユーザー指定の欠損値(空白値)の設定	3-15
Chapter 4 17	マィールドの要約	4-1
- 4.1.1	フィールドの要約の概要	4-2
4.1.2	要約統計量	4-3
4.1.3	ストリームの確認	
4.2.1	記述統計ノード	4-6

4.2.2	記述統計の結果の解釈	4-10
4.3.1	データ検査の実行	4-12
4.3.2	データ検査の結果の解釈 -連続型フィールド	4-14
4.3.3	データ検査の結果の解釈 ーカテゴリ型フィールドー	4-17
4.4.1	欠損値の概要	4-20
4.4.2	データのインポート	4-21
4.4.3	データ検査ノードによる欠損値のチェック	4-23
4.4.4	欠損レコードを除外する条件抽出ノードの生成	4-28

Chapter 5 グラ	ラフの作成と編集	5-1
5.1.1	グラフ作成ノードの概要	5-2
5.1.2	ストリームの確認	5-3
5.2.1	棒グラフノード	5-4
5.2.2	棒グラフの確認と編集	5-7
5.3.1	棒グラフの作成 ー複数のフラグ型ー	5-14
5.4.1	Webグラフの概要	5-16
5.4.2	Webグラフの作成	5-17
5.4.3	Webグラフの解釈	

Chapter 6 27	マールドの関係性	6-1
6.1.1	2フィールドの関係性の概要	6-2
6.1.2	ストリームの確認	6-3
6.2.1	クロス集計表の概要	6-4
6.2.2	セルの内容 - 度数とパーセンテージー	6-5
6.2.3	仮説検定	6-7
6.2.4	ピアソンのカイ2乗検定	6-8
6.2.5	クロス集計ノードの実行	6-9
6.2.6	クロス集計表の結果の解釈	6-12
6.2.7	カイ2乗検定の結果の解釈	6-14
6.2.8	オーバレイ棒グラフの作成	6-16
6.3.1	相関係数の概要	6-19
6.3.2	正の相関、負の相関、無相関	6-20
6.3.3	ピアソンの相関係数	6-23
6.3.4	相関係数の検定	6-25
6.3.5	相関係数の出力	6-26
6.3.6	相関係数の結果の解釈	

6.3.7	散布図の作成	.6-35
6.4.1	平均値の比較の概要	.6-38
6.4.2	平均値ノードの実行	.6-40
6.4.3	平均値ノードの実行結果の解釈	.6-43
6.4.4	平均値を比較するグラフの作成	.6-45





本書では、IBM SPSS Modeler 18.2Jを使用しています。 IBM およびSPSSは、International Business Machines Corp.の登録商標です。

本書を無断で複写複製(コピー)することは、著作権法上の例外を除き、禁じられています。

2 データのインポート

IBM SPSS Modelerでは、データベースやMicrosoft Excel形式、CSV形式などのさまざまなデータ ソースを読み込んで、必要なデータ加工や前処理を行い、分析やグラフ作成、モデリングを行うこ とができます。ここでは、Microsoft Excel形式とCSV形式のデータファイル、IBM SPSS Statistics形 式とデータベースのテーブルを読み込む手順と注意点を確認します。

Contents



§2.1.1 データのインポートの概要

IBM SPSS Modelerでは、入力ノードを使用することで、フラットファイル(カンマ区切りや タブ区切りなどのデータファイル)、IBM SPSS Statistics(.sav)、Microsoft Excel(.xls、.xlsx)、 およびODBC 準拠のリレーショナルデータベースも含めたさまざまなフォーマットのデータ ソースを、インポートすることができます。また、ユーザー入力ノードを使用して、合成デ ータを生成することもできます。

入力ノードは、入力パレットに含まれています。主な入力ノードの例は以下の通りです。

	Excel	Microsoft Excel (.xls、.xlsx)形式のデータをインポートしま す。ODBC データソースは不要です。
	可変長ファイル	フィールドがカンマやタブなどの区切り文字によって区切られているデータファイルをインボートします。
	Statistics	IBM SPSS Statistics(.sav)形式のデータファイルをインポー トします。
SQL	データベース	Microsoft SQL Server、DB2、Oracle などODBCを使用するデ ータベーステーブルをインポートします。
	IBM Cognos Analytics	IBM Cognos BI データベースからデータをインポートします。
	Figure2.1.	入力パレットに含まれる主な入力ノードの例

POINT

IBM SPSS Modelerに分析用のデータソースをインポートするためには、入力パレット内の 入力ノードを使用します。

TIPS

レコード設定パレットの、レコード結合(フィールドの結合)ノードやレコード追加(レコ ードの結合)ノードを使用することで、異なるデータソースをIBM SPSS Modeler上で統合 することができます。

§2.4.1 IBM SPSS Statistics形式のデータファイルの読込み

IBM SPSS Modelerでは、統計解析ソフトウェアのIBM SPSS Statistics形式のデータファイ ルを直接読み込むことができます。IBM SPSS Statisticsでは、変数の補足説明となる**変数ラベル** や、カテゴリ値の意味付けをする**値ラベル**の設定を行うことができ、IBM SPSS Modelerにも取 り込むことができます。また、欠損値の設定などもIBM SPSS Statisticsで設定されている通りに IBM SPSS Modelerに読み込むことができます。

ここでは、以下のIBM SPSS Statistics形式のデータファイルをIBM SPSS Modelerに読み込み ます。このデータファイルは、顧客の属性や利用サービスに関するデータが記録されています。 フィールド(変数)は17個、8,632行のデータが含まれています。



Figure 2.4.1 analysis 1.savファイル

列にフィールド(変数)、行にレコードが含まれています。また、Excelと同じように1行目 には変数名を入力しておき、分析用のデータ値は2行目以降に入力しておきます。1列1変数、 1行1ケースの形式になっていない場合、IBM SPSS Modelerでは正しくデータを読み込むこと ができません。

1. 入力パレットをクリックします。

▲お気に入り ●入力 ●レコード設定 ●フィールド設定 ▲グラフ作成 ●モデル作成 ■出力 ■エクスボート 🥹 IBM® SPSS® Statistics 💈 Python 😭 Spark) ixar <XMI > Ξ, データベース 可変長ファイル 固定長ファイル Statistics ファイル Data Collection IBM Cognos Analytics TM1インポート TWCインポート lytic Server SASファイル Excel XML



操作手順

2. 可変長ファイルノードを選択し、ストリームキャンバスに挿入します。



POINT

Statisticsファイルノードは、IBM SPSS Statistics形式のデータファイルのインポートに使用します。

3. Statisticsノードをダブルクリックして編集画面を開きます。

Statistics ファイル ×
ぼうしどユー(P) で更新 ぼうしどユー(P) ぼうしどユー(P) ぼうしどユー(P) ぼうしどユー(P) ぼうしどユー(P) ぼうしどユー(P) ぼうしどコー(P) ぼうしどコー(P) ぼうしどコー(P) ぼうしどコー(P) ぼうしどコー(P) ぼうしどコー(P) ぼうしどコー(P) ぼうしどコー(P) ぼうしどコー(P) ぼうしどコー(P) ぼうしどコー(P) ぼうしどコー(P) ぼうしどコー(P) ぼうしどコー(P) ぼうしどコー(P) ぼうしどコー(P) ぼうしどコー(P) ぼうしどコー(P) ぼうしどコー(P) ぼうしどコー(P) ぼうしどコー(P) ぼうしどコー(P) ぼうしどコー(P) ぼうしどコー(P) ぼうしどコー(P) ぼうしどコー(P) ぼうしどコー(P) ぼうしどコー(P) ぼうしどコー(P) ぼうしどコー(P) ぼうしどコー(P) ぼうしどコー(P)
データ フィルター データ型 注釈
インポート ファイル(!):
□ ファイルはパスワードで暗号化(C)
パスワード(<u>W)</u> :
変数名: ● 名前とラベルを読み込み(№) 〇 ラベルを名前として読み込み(앁)
値: ④ データとラベルを読み込み(型) 〇 ラベルをデータとして読み込み(型)
OK キャンセル 適用(A) リセット(R)
Figure2.44 Statisticsファイルノードの編集画面 - タ タブを利用して、読込むファイルの指定を行います。
ンポートファイル ボックスに、読込みが必要なデータファイル名を指定します

変数名では、フィールド(変数)の名前とラベルの両方を読み込むか、ラベルのみを読み込 むかを選択します。また、値では、カテゴリ値と値ラベルの両方を読み込むか、値ラベルの みを読み込むかを選択します。この例では、デフォルトのまま進めます。

- 4. インポートファイルのファイルを探すボタンをクリックします。
- 5. ファイルの一覧からanalysis1.savを選択します。

🛞 Statistics ファイル	×
● プレビュー(P) ご 更新 C:\Training\ModelerA\analysis1.sav	
データ フィルター データ型 注釈	_
インポート ファイル(I): C:\Training\ModelerA\analysis1.sav	
□ ファイルはパスワードで暗号化(<u>C</u>)	
バスワード(<u>W</u>):	
変数名: ● 名前とラベルを読み込み(N) ○ ラベルを名前として読み込みし 値: ● データとラベルを読み込み(D) ○ ラベルをデータとして読み込み(D)	
OK キャンセル 適用(A) リセット(R) Figure2.4.5 データスァイルの選択が完了したStatisticsファイルノードのデータタ] 1ブ
操作手順	

- 6. 変数名の名前とラベルを読み込みを選択します。
- 7. 値のデータとラベルを読み込みを選択します。

TIPS

IBM SPSS Statistics形式のデータファイルでは、変数名を補足するための**変数ラベル**と、カ テゴリ値の意味付けをするための**値ラベル**を設定することができます。

8. **データ型**タブをクリックします。

🛞 Statistics ファイル) Statistics วิราไม่ X						
● プ C:\Trainin データ フィル	レビュー(P) C 更新 gModelerA\analysis1.sav パター データ型 注釈	R			0		
∖ - ⊙ Ø	▶ 値の読み込み	値の	肖去	すべての値の)消去		
フィールドロ	尺度	値	欠捐值	检查	ロール		
	▲ 連続型	12	7 CIALE		$\otimes t $		
	☞ 連続型			なし	× λ.π		
A 性別	▲ 名義型	F.M		なし	▲ 入力		
A 婚姻状況	♣ 名義型	既婚,未婚	「「」」 「「」」 「「」」 「」」 「」」 「」」 「」」 「」」 「」」		▲ 入力		
A 支払方法	山 順序型	カード ボ	なし		▲ 入力		
● 支払金額	☞ 連続型			なし	入 力		
	☞ 連続型			なし	🔪 入力		
	☞ 連続型			なし	🔪 እ.ታ		
	☞ 連続型			なし	🔪 ኢታ		
🕸 サービスA	☞ 連続型			なし	🔨 ኢታ		
🕸 サービスB	☞ 連続型			なし	እ እ. አ.		
🕸 サービスC	☞ 連続型			なし	🔪 ኢታ		
🕸 家の所有	☞ 連続型			なし	🔪 入力		
🕸 車の所有	🔗 連続型			なし	🔪 ኢታ		
🗶 顧客紹介	🔗 連続型 🛛 🚺			なし	🔪 ኢታ		
A 申込方法	🖧 名義型	Web.書面		なし	🔪 ኢታ		
🕸 新規契約	♣ 名義型	0.0,1.0		なし	🔪 入力		
● 現在のフィー	ルドを表示 🔵 未使用の	ウフィールド語	設定を表示		,		
ок	キャンセル		j	適用(<u>A</u>)	リセット(<u>R</u>)		

Figure2.4.6 Statisticsノードのデータ型タブ

データ型タブでは、各フィールドの**尺度や値、欠損値やロール**の設定を行うことができま す。IBM SPSS Statisticsで設定済みの定義は、自動的にセットされます。これらの設定は、集 計やグラフ作成、モデリングに影響するため非常に重要な設定ですが、尺度の設定等は次の 章で確認しますので、ここでは設定しません。

TIPS

IBM SPSS Statisticsでは、変数ビューを利用して、変数の測定の尺度や値ラベル、ユーザ ー欠損値の指定などを行っておくことができ、これらの情報をIBM SPSS Modelerに読み込 むことができます。

9. **OK**ボタンをクリックします。



Figure2.4.7 ファイルを設定したStatisticsファイルノード

ストリームキャンバスに表示されるノード名が、設定したファイル名に変わります。次に、 **テーブル**ノードを利用して、データをIBM SPSS Modelerに読み込みます。



操作手順

- 12. **テーブル**ノードをダブルクリックします。
- 13. ツールバーの選択内容を実行ボタンをクリックします。

	顧客番号	年齢	性別	婚姻状況	支払方法	支払金額	契約 <mark>A</mark>	契約B	
1	47.000	49.000	F	既婚	ポイント	3179.000	0.000	0.000	d.
2	1093.000	51.000	M	未婚	現金	9888.000	0.000	1.000	1
3	3191.000	39.000	F	未婚	現金	9173.000	1.000	1.000	
4	3308.000	58.000	F	既婚	ポイント	4279.000	1.000	1.000	
5	5580.000	38.000	M	未婚	カード	5720.000	0.000	1.000	
6	5974.000	40.000	F	未婚	カード	6164.000	0.000	1.000	
7	9099.000	62.000	F	既婚	ポイント	4365.000	1.000	1.000	
8	9112.000	45.000	F	未婚	現金	14901.000	0.000	0.000	
9	9128.000	39.000	F	未婚	現金	9440.000	1.000	1.000	
10	9388.000	53.000	F	既婚	現金	9049.000	0.000	1.000	
11	9390.000	48.000	M	既婚	ボイント	4319.000	1.000	1.000	
12	12236.000	57.000	F	既婚	現金	6518.000	0.000	1.000	
13	12309.000	35.000	F	未婚	現金	14281.000	1.000	1.000	
14	12343.000	26.000	F	未婚	カード	10639.000	0.000	1.000	
15	138/1.000	37.000	M	未婚	カード	12762.000	1.000	1.000	
16	14211.000	35.000	M	未婚	現金	4644.000	0.000	0.000	
1/	14921.000	44.000	M	未婚	現金	5427.000	0.000	1.000	
18	15068.000	40.000	F	5.大汉语 800.555	ガード	20065.000	0.000	1.000	
19	15167.000	63.000	F	民党刘昚 800455	ホイント 現今	6470.000	1.000	1.000	
20	15252.000	50.000	F	民发现管	現金	11/51.000	1.000	1.000	
								ОК	
	Fi	igure2.4.9	テーブル	に出力さ	いたanalys	sis1.csvのラ	ギータ		

14. テーブルを閉じます。



フィールドは、変数 variables とも呼ばれ、年齢や性別、購入商品やアンケートの設問など、個々 のデータを記録する項目です。個々のフィールドに注目してフィールドの持つ情報を要約すること を、1フィールドの要約と呼びます。

フィールドの要約を行う場合、各フィールドの**尺度**の設定が重要です。1フィールドの要約は、 クロス集計表や相関係数を計算する2フィールドの関係性の分析、ディシジョンツリーやロジステ ィック回帰分析などのようなモデリングの基礎にもなる重要な処理です。

Contents



フィールドの要約/要約統計量/平均値/中央値/最頻値/最小値/最大値/範囲/ 標準偏差/分散/四分位範囲/パーセンテージ/記述統計ノード/データ検査ノード/ ヒストグラム/棒グラフ/欠損値/欠損値選択ノード/

§4.1.1 フィールドの要約の概要

フィールドが持つデータは、さまざまな方法で要約することができます。最も基本となるのは**1フィールドの要約**です。基本的に、連続型フィールドは平均値をはじめとする要約統計量によって要約し、カテゴリ型フィールドはパーセンテージなどによって要約します。



Figure4.1.1 テーブルノードでの生データの表示

データ分析においては、できるだけ生データを確認することが重要です。しかし、フィー ルドやレコードの数が大きくなるほど、生データから何らかの傾向や特徴を見出すことは困 難になります。そこで、さまざまな統計量を用いてデータの要約を図ります。

フィールドの要約の最も基本となるのが、1フィールドの要約です。さらに、2フィールド の関係性、多数のフィールドに基づくモデリングと拡張されていきます。

§4.1.2 要約統計量

要約統計量 summary statistics は、フィールドの持つ情報を要約するための指標です。記述統計量とも呼ばれます。要約統計量は、使用するフィールドの尺度によって、適切に使い分ける必要があります。フィールドの要約では、**中心傾向**と**散らばり**の2つの視点が必要です。

POINT

中心傾向を示す要約統計量には、平均値、中央値、最頻値があります。

平均値 mean は、フィールドのすべてのデータ値を足して、データ数で割ることによって得 られる値です。もっともよく利用される中心傾向の要約統計量ですが、外れ値や分布の歪みの 影響を受けやすい特性を持ちますので、使用する際はデータの分<u>布</u>を確認することが重要です。

中央値 median は、データ値を大きさの順番に並び替えたときに、ちょうど真ん中に位置する値です。中央値より大きい値と小さい値がちょうど同数になります。データ数が偶数の場合は、真ん中の2つの値の平均値が中央値になります。中央値は外れ値の影響を受けません。

最頻値 mode は、もっとも頻度の高い値です。多数決的な要約統計量です。

それぞれの測定の尺度で使用することができる、要約統計量は以下のとおりです。

		名義型	順序型	連続型
平均值	Mean	×	\bigtriangleup	0
中央値	Median	×	0	0
最頻値	Mode	0	0	0

Table4.1.1 中心傾向の要約統計量と測定の尺度

TIPS

順序型は、原則として平均値を使用することは適切ではありません。ただし、アンケート 評価のようなフィールドで、5段階以上の尺度がある場合は、便宜上等間隔とみなして統 計処理を行うケースが多いです。この場合は、データを数値として入力しておき、尺度を 順序型ではなく連続型として設定しておく必要があります。

POINT

変数の散らばりを示す要約統計量には、最小値、最大値、範囲、四分位範囲、分散、標準 偏差があります。

最小値 minimum と最大値 maximum は、それぞれのデータ値のうち最も大きな値と最も小さ な値を意味します。最大値と最小値の差が範囲 range です。

四分位範囲 inter quartile range, IQR は、データを並び替えてサンプルの25%ずつが含まれる ように4等分したときの真ん中の50%の範囲(75パーセンタイル値と25パーセンタイル値の差 分)です。

分散 variance は、平均値に対するデータのばらつきを意味し、分散の平方根が標準偏差 Standard Deviation, SDです。分散はデータ値と平均値の差(偏差)を2乗して計算されるため、元 のデータと単位が異なりますが、標準偏差は分散の平方根となるため、元のデータと同じ単位 で解釈できます。

		名義型	順序型	連続型
最小値	Minimum	×	0	0
最大値	Maximum	×	0	0
範囲	Range	×	0	0
四分位範囲	Inter Quartile Range	×	0	0
分散	Variance	×	\bigtriangleup	0
標準偏差	Standard Deviation	×	\bigtriangleup	0

Table4.1.2 散らばりの要約統計量と測定の尺度

POINT

分散はデータ値と平均値の差(偏差)を2乗して計算されるため、元のデータと単位が異な りますが、標準偏差は分散の平方根となるため、元のデータと同じ単位で解釈できます。

§4.1.3 ストリームの確認

1フィールドの要約の手順を確認するために、サンプルのストリームファイルChapter4.str を読み込みます。この章では、記述統計ノードやデータ検査ノードによる分析を行います。

操作手順

- 1. ファイルメニュー > ストリームを開くを選択します。
- 2. / ModelerA / Chapter4.strを選択します。



Figure4.1.2 Chapter4.strファイルのストリーム

このファイルには、カンマ区切りのデータソース**analysis1.csv**を読み取り、テーブルに出 力する簡単なストリームが含まれています。

入力ノードとテーブルノードの間に、データ型ノードが挿入されており、尺度やロールの設 定が行われています。

§4.3.1 データ検査の実行

データ検査ノードを用いると、各フィールドの要約情報を出力することができます。フィ ールドの尺度に応じたグラフとして、カテゴリ型フィールドでは棒グラフ、連続型フィール ドではヒストグラムが自動的に作成されます。

データ検査ノードは、出力パレットに含まれています。

操作手順

1. 出力パレットを開きます。



Figure4.3.2 データ型ノードからリンクしたデータ検査ノード

データ型ノードから16個のフィールを受け取っていることが分かります。



3. データ検査ノードをダブルクリックして編集画面を開きます。

設定タブを利用して、データ検査を行うフィールドを指定します。デフォルトでは、上流 から受け取ったフィールドが使用されます。ユーザー設定フィールドを使用に切り替えると、 個別にフィールドの指定を行うことができます。

表示では、**グラフ作成と基本統計量**が指定されています。基本統計量では、連続型フィー ルドについて平均値、標準偏差、最小値、最大値、歪度を計算して出力します。高度な統計 を選択すると、合計、範囲、平均値の標準誤差、尖度などが追加で出力されます。中央値と 最頻値はデフォルトでは出力されませんので、使用する場合は選択する必要がありますが、 データの並べ替えが行われるため、パフォーマンスが低速になる可能性があります。

操作手順

4. 実行ボタンをクリックします。

§4.3.2 データ検査の結果の解釈 -連続型フィールド-

出力されるデータ検査の結果を確認します。この例では、16個のフィールドの情報が要約 されており、連続型フィールドは最小値、最大値、平均値、標準偏差などによってまとめら れ、フラグ型や名義型などのカテゴリフィールドはカテゴリ数がまとめられています。



Figure4.3.4 データ検査ノードの実行結果

TIPS

フィールドのロールに対象がセットされていると、対象フィールドの値ごとにオーバーレイされたグラフが出力されます。

フィールドー	サンブル グラフ	尺度	最小値	最大値	平均値	標準 偏差	歪度	カテゴリ数	有効
◇ 年齢		☞ 連続型	18	63	43.496	10.858	-0.071		8632

Figure4.3.5 年齢フィールドのデータ検査結果

年齢は連続型のフィールドです。最小値=18、最大値=63であり、顧客の年齢は18歳から63歳の間です。平均値=43.496であり、顧客の平均年齢は約43歳であることが分かります。また、標準偏差=10.858であり、平均年齢に対して約11歳のばらつきを持つと解釈することができます。これらは、記述統計ノードで確認した結果と同じです。

フィールドが正規分布に近い場合、平均値±1標準偏差内におよそ60~70%のデータが含ま れることが示唆されます。つまり、43±11の計算から、顧客の60~70%は、32歳~54歳の範囲 に含まれそうです。



サンプルグラフをダブルクリックすると、拡大して詳細を確認することができます。

1. 年齢のサンプルグラフをダブルクリックします。



連続型フィールドを要約するグラフとして**ヒストグラム**が表示されます。ヒストグラムの横軸には連続型フィールド、縦軸に度数が表示されます。

POINT

連続型フィールドを要約するグラフはヒストグラムです。

操作手順

2. OKボタンをクリックして、ヒストグラムを閉じます。