

# Spectrum™ Technology Platform

バージョン 2019.1.0

Machine Learning ガイド



# 目次

## 1 - はじめに

---

Machine Learning モジュール	5
Machine Learning ワークフロー	6

## 2 - Binning

---

Binning の概要	8
ビニングのプロパティの定義	8
基本オプションの設定	9
ビニング出力	10

## 3 - K-Means Clustering

---

はじめに	12
モデルのプロパティの定義	12
基本オプションの設定	13
高度なオプションの設定	13
モデル出力	14
出力ポート	15

## 4 - Linear Regression

---

はじめに	17
モデルのプロパティの定義	17
基本オプションの設定	18
高度なオプションの設定	19
モデル出力	22
出力ポート	22

## 5 - Logistic Regression

---

はじめに	25
モデルのプロパティの定義	25
基本オプションの設定	26
高度なオプションの設定	26
モデル出力	30
出力ポート	30

## 6 - 主成分分析

---

はじめに	33
モデルのプロパティの定義	33
基本オプションの設定	34
高度なオプションの設定	34
モデル出力	35
出力ポート	35

## 7 - Random Forest Classification

---

はじめに	38
モデルのプロパティの定義	38
基本オプションの設定	39
高度なオプションの設定	40
モデル出力	43
出力ポート	44

## 8 - Random Forest Regression

---

はじめに	47
------	----

モデルのプロパティの定義	47
基本オプションの設定	48
高度なオプションの設定	49
モデル出力	52
出力ポート	52

## 9 - Machine Learning モデル管理

---

Machine Learning モデル管理へのアクセス	55
モデル評価	56
ビニング管理	64

## 10 - Data Science Demonstration Flows

---

はじめに	66
教師あり学習: 貸付返済不能予測	66
教師なし学習: セグメンテーション	67

# 1 - はじめに

## このセクションの構成

---

Machine Learning モジュール	5
Machine Learning ワークフロー	6

# Machine Learning モジュール

Spectrum™ Technology Platform Machine Learning モジュールを使用すると、数値データをグループ化 (ビンニング) して、教師ありと教師なしの機械学習モデルのデータを、これらのモデルに適合させることができます。

注：Machine Learning モジュールは、Windows と Linux の各オペレーティングシステムでのみサポートされています。

注：Machine Learning モジュールは、K-Means Clustering、Linear Regression、Logistic Regression、主成分分析、Random Forest Classification、および Random Forest Regression のモデリングアルゴリズムに、基盤となる H2O.ai ライブラリを使用します。

## *Binning*

Binning は、目標情報を考慮に入れずに、連続変数のレコードをグループ (ビン) に分類します。均等幅ビンと均等個数ビンという 2 つのいずれかの方法で、教師なしビンニングを実行できます。

## *K-Means Clustering*

K-Means Clustering は、分析クラスタリングに基づくモデルを作成します。このクラスタリングでは、一連のレコードをデータ値に基づく類似レコードのクラスタに分割します。

## *Linear Regression*

Linear Regression では、持続的目標と入力変数を使用するデータセットからモデルを作成して、機械学習を実行することができます。

## *Logistic Regression*

Logistic Regression は、バイナリ目標と入力変数を使用するデータセットからモデルを作成します。

## *主成分分析*

主成分分析は、相関のある可能性がある変数群の観測データの集合を、主成分と呼ばれる線型相関のない変数の値の集合に変換する統計的な処理です。

## *Random Forest Classification*

Random Forest Classification では、持続的目標と入力変数を使用するデータセットからモデルを作成して、機械学習を実行することができます。

## Random Forest Regression

Random Forest Regression では、バイナリ目標と入力変数を使用するデータセットからモデルを作成して、機械学習を実行することができます。

## Machine Learning モデル管理

Machine Learning モデル管理には、Spectrum™ Technology Platform サーバー上のすべての機械学習モデルを管理できるモデル評価と、Spectrum™ Technology Platform サーバー上のすべてのビンニングを管理できるビンニング管理が含まれています。

# Machine Learning ワークフロー

標準的な機械学習ワークフローは、1つ以上のデータフローで行われる以下のステップで構成されます。

1. Data Integration など、Spectrum の他のモジュールを使用して、データにアクセスします。
2. Data Integration、Data Quality、および各種の Core モジュールなど、Spectrum の他のモジュールのステージを使用して、データを準備します。
3. 機械学習モデルを適合し、データフローを実行してから、モデル ステージの [モデル出力] タブを確認します。必要に応じてモデルに微調整を加え、データフローを再実行します。その後、Machine Learning モデル管理ツールのモデル評価出力全体を確認する必要があります。モデルを 1 度に 1 つずつ確認するか、2 つのモデルを比較することができます。
4. オプション: モデルをデータのスコアリングに使用する場合は、モデルを Machine Learning モデル管理ツールでエクスポートします。これにより、そのモデルは Java Model Scoring ステージで使用可能になります。
  - a) 上のステップ 1 (6ページ) ~ 2 (6ページ) によって Spectrum™ Technology Platform データフローを作成し、ステップ 3 を Java Model Scoring ステージに置き換えます。このデータフローをバッチ モードで実行するように設定し、更新されたデータに適用されたモデル スコアを、ファイルに設定します (自然な処理の流れとして、X または入力として使用されたフィールドがステップ 1 (6ページ) ~ 2 (6ページ) で更新されます)。
  - b) あるいは、Spectrum™ Technology Platform の Web サービスを使用してオンデマンドでデータをスコアリングします。例えば、Web サイトにアクセスして顧客 ID とモデル入力を取得し、それらをスコアリングして、顧客向けに Web コンテンツをカスタマイズするプロセスにそのスコアを返します。
5. オプション: モデル スコアは、Data Hub グラフ データベースにエンティティ プロパティとして展開するか、マップ上に展開するか、または CES アプリケーションに展開することもできます。

# 2 - Binning

## このセクションの構成

---

Binning の概要	8
ビンニングのプロパティの定義	8
基本オプションの設定	9
ビンニング出力	10

## Binning の概要

Binning ステージは、目標情報を考慮に入れずに、連続変数をグループ (ビン) に分類する、教師なしビニングとして知られる処理を実行します。取得されるデータには、レンジ、個数、各レンジ内の値の割合などがあります。

ビニングの実行には、次のような利点があります。

- データが欠落しているレコードをモデルに含めることができる。
- 外れ値がモデルに与える影響を制御または緩和することができる。
- 最終モデルの係数の重みを同等にすることによって、特性によって尺度が異なる問題を解決する。

Spectrum™ Technology Platform の教師なしビニングでは、データを同じサイズのビンに分割する均等幅ビン、または、データをほぼ同数のレコードを含むグループに分割する均等個数ビンが使用できます。Binning ステージでは、均等幅ビンは、[Equal Ranges] ビン、均等個数ビンは、[Equal Count] ビンと呼ばれます。

Machine Learning モデル管理の**ビニング管理**ツールを使用すると、より多くのビニング関数を実行できます。

コマンドラインの命令を使用して、ビニングのリストを表示したり、ビニングを削除したりすることもできます。『管理ガイド』の「管理ユーティリティ」セクションにある「Binning」を参照してください。

注：Spectrum™ Technology Platform をバージョン 12.0 SP1 から 12.0 SP2 にアップグレードする場合は、Machine Learning モデル管理のビニング管理ツールですべてのアップグレード対象ビニングを手動で公開解除してから、それらを 12.0 SP2 での再ビニングのために使用する必要があります。従来の Binning ではなく、Binning Lookup でアップグレード対象ビニングを使用する場合、この手順は必要ありません。

## ビニングのプロパティの定義

1. [プライマリ ステージ] > [展開済みステージ] > [Machine Learning] の下で、[Binning] ステージをクリックしてキャンバス上にドラッグし、データフロー内の所望の位置に配置して、他のステージに接続します。



注: 入力ステージは、モデルの目標フィールドと入力変数フィールドの両方を含むデータソースでなければなりません。出力ステージは、**【基本オプション】** タブで**【入力データを記録】** オプションを選択しない限り、不要です。Machine Learning モデル管理ツールとは独立して出力を取得する場合は、出力ステージを接続することもできます。

2. Binning ステージをダブルクリックして、**【ビンニング オプション】** ダイアログ ボックスを表示します。
3. デフォルトの名前を使用しない場合は、**【ビンニング名】** を入力します。
4. **【上書き】** チェックボックスをオンにして、既存モデルを新しいデータで上書きします。
5. モデルの**【説明】** を入力します。
6. データをビンニングに含める各フィールドに対し、**【含める】** をクリックします。このリストには、数値フィールドしか表示されません。
7. **【OK】** をクリックして、設定を保存します。

## 基本オプションの設定

1. レンジ幅均等とレコード数均等のどちらの**【ビンニング スタイル】** を実行するかを選択します。
2. **【NULL 値ビン】** で、空のビンフィールド (データが欠落しているために値が不明であることを表します) の処理方法を選択します。
  - NULL 値を最高ビンに割り当てる場合は **【最高】** を選択します。
  - NULL 値を最低ビンに割り当てる場合は **【最低】** を選択します。

最低ビンは、必ずビン 1 です。

3. **【ターゲット内部ビン】** をクリックして、両端のビンの間のビン数を入力します。  
レンジ幅均等ビンニングを実行する場合は、内部ビン処理を選択しても、**【ビン幅】** を選択してもよいですが、両方を選択することはできません。レコード数均等ビンニングを実行する場合は、内部ビン処理しか実行できません。
4. レンジ幅均等ビンニングを実行し、内部ビン処理ではなくビン幅を選択する場合は、**【ビン幅】** をクリックして、各ビンに含める個数を入力します。
5. データをビンニングに含める各フィールドに対し、**【含める】** をクリックします。

注: このリストには、数値フィールドしか表示されません。

6. **【OK】** をクリックして、設定を保存します。

## ビンニング出力

Binning ステージには 2 つの出力ポートがあります。1 つめのポートは、すべての入力フィールドに加えて、選択された各入力フィールドに対するビンニング済みフィールドを出力します。例えば、入力フィールドに **Name**、**Age**、**Income** のフィールドがあり、**Age** と **Income** に対してビンニングを実行する場合、1 つめのポートからは、以下のフィールドが出力されます。

- 名前
- Age
- Binned\_Age
- Income
- Binned\_Income

2 つめのポートは、選択された各入力フィールドに対する 4 種類の情報を出力します。例えば、**Age** に対してビンニングを実行する場合、2 つめのポートからは、以下のフィールドが出力されません。

- Age\_Bins
- Age\_BinValue
- Age\_Count
- Age\_Percentage

# 3 - K-Means Clustering

## このセクションの構成

---

はじめに	12
モデルのプロパティの定義	12
基本オプションの設定	13
高度なオプションの設定	13
モデル出力	14
出力ポート	15

## はじめに

K-Means Clustering は、分析クラスタリングに基づくモデルを作成します。このクラスタリングでは、一連のレコードをデータ値に基づく類似レコードのクラスタに分割します。

モデルを作成するにはまず、**[モデルのプロパティ]** タブを設定する必要があります。**[基本オプション]** タブと **[高度なオプション]** タブには、ジョブを完了するために十分なデフォルト設定が指定されていますが、ニーズに合わせてその設定を変更することができます。続いてジョブを実行すると、結果として得られたモデル出力詳細情報の限定版が **[モデル出力]** タブに表示されます。モデルは Spectrum™ Technology Platform サーバーに格納され、出力全体は、Machine Learning モデル管理ツールで確認できます。

## モデルのプロパティの定義

1. **[プライマリ ステージ]** > **[展開済みステージ]** > **[Machine Learning]** の下で、**[K-Means Clustering]** ステージをクリックしてキャンバス上にドラッグし、データフロー内の所望の位置に配置して、他のステージに接続します。

注：入力ステージは、モデルの入力変数フィールドを含むデータ ソースでなければなりません。出力ステージは、**[基本オプション]** タブで **[入力データを記録]** オプションを選択しない限り、不要です。Machine Learning モデル管理ツールとは独立して出力を取得する場合は、出力ステージを接続することもできます。

2. K-Means ステージをダブルクリックして、**[K-Means Clustering オプション]** ダイアログボックスを表示します。
3. デフォルトのモデル名を使用しない場合は、**[モデル名]** を入力します。
4. オプション: **[上書き]** チェックボックスをオンにして、既存モデルを新しいデータで上書きします。
5. モデルのクラスタ数をデフォルト値 (5) 以外にする場合は、**[クラスタの数]** を入力します。
6. オプション: モデルの **[説明]** を入力します。
7. モデルにデータを追加するフィールドの **[含める]** をクリックします。
8. **[モデル データ タイプ]** ドロップダウンを使って、入力フィールドを数値、カテゴリ値、日付と時刻のいずれのフィールドとして使うかを指定します。
9. **[OK]** をクリックして、モデルと設定を保存するか、次のタブで操作を続行します。

## 基本オプションの設定

1. 数値列を標準化し、平均が 0 で偏差が 1 になるようにするには、**[標準化]** をオンのままにします。  
標準化を使わない場合、実際の寄与の大きさではなく尺度の違いから、他の属性と比べて偏差が大きくなる変数によって左右される要素が結果に含まれる場合があります。
2. **[クラスタ数を見積もる]** をオンにすると、K-Means アルゴリズムによって、モデルに含めるクラスタ数の判定が試みられます。**[モデルのプロパティ]** タブで所望のクラスタ数を指定した場合でも、データから判断して異なるクラスタ数の方が適切であることが、この処理によって検出される可能性があります。
3. 入力データがトレーニングおよびテストのデータ サンプルにランダムに分割される場合は、1 ~ 100 の値を **[トレーニング データの比率]** に指定します。
4. ステップ 3 (13 ページ) で指定した値を 100 から引いた値を **[テスト データの比率]** に入力します。
5. データフローを何度実行しても、データが必ず同じ方法でテスト データとトレーニング データに分割されるようにするには、**[テスト データ用シード]** に数値を入力します。フローを実行するたびにランダムな分割を行う場合は、このフィールドをオフにします。
6. **[OK]** をクリックして、モデルと設定を保存するか、次のタブで操作を続行します。

## 高度なオプションの設定

1. **[定数フィールドを無視]** をオンにすると、各レコードで値が同じフィールドがスキップされます。
2. **[アルゴリズムのシード]** をオンにしてシード数を入力すると、データフローを何度実行しても、データが必ず同じ方法でテスト データとトレーニング データに分割されます。フローを実行するたびにランダムな分割を行う場合は、このフィールドをオフにします。
3. **[初期化]** ドロップダウンで、適切な初期化モードを選択します。

### 初期化モード 説明

<b>Furthest</b>	最初の中心点はランダムに初期化しますが、2 つめの中心点はそれから最も遠いデータ ポイントになるように初期化します。互いに大きく分散するように、中心点を初期化します。
-----------------	-------------------------------------------------------------------------------------

## 初期化モード 説明

**Plus-Plus (++)** 標準の  $k$ -means の再帰的最適化を行う前に、クラスタの中心を初期化します。  $k$ -means++ の初期化を行うと、アルゴリズムによって、最適な  $k$ -means ソリューションに  $O(\log k)$  近似のソリューションが検出されることが保証されます。

**Random**  $N$  個のオブザベーション集合から  $K$  個のクラスタを、各オブザベーションの選択確率が等しくなるようにランダムに選択します。これはデフォルトの初期化モードです。

4. **[アルゴリズムのシード]** をオンにしてシード数を入力すると、データフローを何度実行しても、データが必ず同じ方法でテスト データとトレーニング データに分割されます。フローを実行するたびにランダムな分割を行う場合は、このフィールドをオフにします。
5. 相互検証を実行する場合は **[N フォールド]** をオンにし、フォールドの数を入力します。
6. 相互検証を実行する場合は、**[フォールド割り当て]** をオンにしてドロップダウン リストから選択します。

## フォールド割り当て 説明

**Auto** オプションの自動選択をアルゴリズムに任せます。現在、**[ランダム]** が選択されます。これがデフォルトです。

**Modulo** データセットをフォールドに等分し、シードを基準としません。

注：このフィールドは、**[N フォールド]** に値が入力済みの場合のみ適用可能です。

7. **[最大反復回数]** をオンにし、実行する必要があるトレーニング反復回数を入力します。
8. **[OK]** をクリックして、モデルと設定を保存するか、次のタブで操作を続行します。

## モデル出力

このタブには、適合されたモデルの評価に使用するメトリクスが表示されます。これらのフィールドは編集できません。**[トレーニング]** 列には、必ずデータが含まれます。**[基本オプション]** タブでテストとトレーニングの分割を選択した場合は、**[テスト]** 列にもデータが設定されます。ただし、**[高度なオプション]** タブで  $N$  フォールド検証を選択した場合を除きます。その場合は、**[N**

**フォールド**] 列にデータが設定されます。出力を再生成するには **[出力]** ボタンをクリックし、出力全体を Machine Learning モデル管理ツールに表示するには **[モデルの詳細]** をクリックします。

## 出力ポート

K-Means Clustering ステージには、1つのオプション出力ポート、モデル メトリクス ポートが含まれています。このポートの機能は、ステージの基本オプションや高度なオプションの設定完了時の選択内容と入力情報によって決まります。例えば、**[高度なオプション]** タブで **[N フォールド]** フィールドをオンにして N フォールド検証の実施を選択した場合、出力メトリクスの **[N フォールド]** 列にデータが設定されます。また、N フォールド検証を実施しないことを選択した場合、**[N フォールド]** 列は空欄になります。

## モデル メトリクス ポート

モデル メトリクス ポートを使用するには、以下の手順に従います。

モデル メトリクス ポートを使用すると、モデル評価メトリクスをデータ ファイルに出力できます。これは、Spectrum™ Technology Platform の内部または外側で生成された多数のモデルを比較したり、メトリクスに関するその他のデータ処理タスクを実行したりするのに役立ちます。

1. K-Means Clustering ステージを使用しているデータフローを開きます。
2. Write to File ステージまたは別のデータ出力ステージを 2 番目の出力ポートに接続します。
3. ジョブを実行します。
4. ステップ **3** (15ページ) の代替: ステップ **2** (15ページ) で追加したシンク ステージに K-Means Clustering ステージを接続しているチャンネルにインスペクション ポイントを追加します。そのためには、チャンネルを右クリックし、**[インスペクション ポイントの追加]** を選択します。その後、Enterprise Designer ツールバーの **[現在のフローのインスペクション]** ボタン (🔍) をクリックします。

インスペクションが実行され、以下に示すような結果が表示されます。

Creation Time	Flow Name	Metrics	Model Name	Model Type	N Fold	Test	Training
10/11/2018 2:37:00 AM	LinearRegressionTest1...	MSE	LinearRegressionTest1...	Linear Regression	0.0050019841090508	0.00321814876650744	1.31531988837612E-10
10/11/2018 2:37:00 AM	LinearRegressionTest1...	RMSE	LinearRegressionTest1...	Linear Regression	0.0707247064967455	0.0567287296394643	1.14687396359675E-05
10/11/2018 2:37:00 AM	LinearRegressionTest1...	Number of observations	LinearRegressionTest1...	Linear Regression	33	17	33
10/11/2018 2:37:00 AM	LinearRegressionTest1...	R2	LinearRegressionTest1...	Linear Regression	-0.366538763835926	0.0268719304981138	0.9999999964065547
10/11/2018 2:37:00 AM	LinearRegressionTest1...	Mean residual deviance	LinearRegressionTest1...	Linear Regression	0.0050019841090508	0.00321814876650744	1.31531988837612E-10

# 4 - Linear Regression

## このセクションの構成

---

はじめに	17
モデルのプロパティの定義	17
基本オプションの設定	18
高度なオプションの設定	19
モデル出力	22
出力ポート	22



## はじめに

Linear Regression では、持続的目標と入力変数を使用するデータセットからモデルを作成して、機械学習を実行することができます。

モデルを作成するにはまず、**[モデルのプロパティ]** タブを設定する必要があります。**[基本オプション]** タブと **[高度なオプション]** タブには、ジョブを完了するために十分なデフォルト設定が指定されていますが、ニーズに合わせてその設定を変更することができます。ジョブを実行すると、最終的なモデルが、限定された形式で **[モデル出力]** タブに表示されます。完全な形式の出力を確認するには、Machine Learning モデル管理ツールを使用します。

## モデルのプロパティの定義

1. **[プライマリ ステージ]** > **[展開済みステージ]** > **[Machine Learning]** にある **[Linear Regression]** ステージをクリックしてキャンバス上にドラッグします。さらに、そのステージを、データフロー内の該当位置に配置して、他のステージに接続します。入力ステージは、モデルの目標フィールドと入力変数フィールドの両方を含むデータ ソースでなければならないことに注意してください。出力ステージは、**[基本オプション]** タブで **[スコア]** 入力データ オプションを選択しない限り、不要です。Machine Learning モデル管理ツールとは独立して出力を取得する場合は、出力ステージを接続することもできます。
2. **[Linear Regression]** ステージをダブルクリックして、**[Logistic Regression オプション]** ダイアログ ボックスを表示します。
3. デフォルトのモデル名を使用しない場合は、**[モデル名]** を入力します。
4. **[上書き]** チェックボックスをオンにして、既存モデルを新しいデータで上書きします。
5. **[目標フィールド]** ドロップダウンをクリックし、数値フィールドを選択します。
6. モデルの **[説明]** を入力します。
7. データをモデルに追加したいそれぞれのフィールドで **[含める]** をクリックします。目標フィールドとして選択したフィールドは必ず含めてください。
8. **[モデル データ タイプ]** ドロップダウンを使用して、各入力フィールドを数値、カテゴリ値、日付と時刻のいずれのフィールドとして使うかを指定します。
9. **[OK]** をクリックして、モデルと設定を保存するか、次のタブで操作を続行します。

## 基本オプションの設定

1. 数値列を標準化し、平均が 0 で偏差が 1 になるようにするには、**【標準化】** をオンのままにします。  
標準化を使わない場合、実際の寄与の大きさではなく尺度の違いから、他の属性と比べて偏差が大きくなる変数によって左右される要素が結果に含まれる場合があります。
2. モデル予測 (スコア) を表す列を入力データに追加するには、**【入力データを記録】** をオンにします。
3. ドロップダウン リストから **【リンク機能】** を選択します。これは、ランダムなコンポーネントと体系的なコンポーネントとの間のリンクを指定するものです。応答で期待される値を説明変数の線型予測因子にどのように関連付けるかを示します。

リンク機能	説明
特定	0 未満または 1 を超えるような無意味な「確率」を予測します。線型確率モデルを得るために二項データで使用されることがあります。 $g(p) = p$
逆変換	実際の見積もり値のためのリンク関数の逆関数を計算します。 $g(\mu) = 1/\mu$
ログ	定められた時間および空間の範囲内での発生回数をカウントします。 $g(\mu) = \log(\mu)$

4. 欠落データの処理方法を指定するには、**【スキップ】** または **【平均値を補完】** をオンにします。後者のオプションを選択すると、欠落データの代わりに平均値が追加されます。
5. 入力データがトレーニングおよびテストのデータ サンプルにランダムに分割される場合は、1 ~ 100 の値を **【トレーニング データの比率】** に指定します。
6. ステップ 5 (18ページ) で指定した値を 100 から引いた値を **【テスト データの比率】** に入力します。

- データフローを何度実行しても、データが必ず同じ方法でテスト データとトレーニング データに分割されるようにするには、**[テスト データ用シード]** に数値を入力します。フローを実行するたびにランダムな分割を行う場合は、このフィールドをオフにします。
- [OK]** をクリックして、モデルと設定を保存するか、次のタブで操作を続行します。

## 高度なオプションの設定

- [定数フィールドを無視]** をオンにすると、各レコードで値が同じフィールドがスキップされます。
- [p 値を計算]** をオンにすると、パラメータを予測するための p 値が計算されます。
- モデルの作成時に共線列を自動的に削除するには、**[共線列を削除]** をオンにします。  
このオプションは、**[p 値を計算]** がオンになっている場合は常にオンにする必要があります。  
これにより、返されるモデルでは係数が 0 になります。
- 定数項 (切片) をモデルに含めるには、**[定数項 (切片) を含める]** をオンにします。  
**[共線列を削除]** がオンの場合は、このオプションを必ずオンにする必要があります。
- ドロップダウン リストから **[ソルバー]** を選択します。

ソルバー	説明
<b>Auto</b>	入力データとパラメータに基づいてソルバーが決定されます。
<b>CoordinateDescent</b>	最も内側のループにおける循環座標降下法の共分散更新バージョンを使う IRLSM。
<b>CoordinateDescentNaive</b>	最も内側のループにおける循環座標降下法のネイティブ更新バージョンを使う IRLSM。
<b>IRLSM</b>	予測因子が少数のときの問題や、L1 ペナルティによるラムダ検索の問題に最適です。

注： **CoordinateDescent** および **CoordinateDescentNaive** は、現時点では実験用です。

- [アルゴリズムのシード]** をオンにしてシード数を入力すると、データフローを何度実行しても、データが必ず同じ方法でテスト データとトレーニング データに分割されます。フローを実行するたびにランダムな分割を行う場合は、このフィールドをオフにします。

7. 相互検証を実行する場合は **[N フォールド]** をオンにし、フォールドの数を入力します。
8. 相互検証を実行する場合は、**[フォールド割り当て]** をクリックしてドロップダウン リストから選択します。このフィールドは、**[N フォールド]** に値が入力済みで、**[フォールド フィールド]** が指定されていない場合にのみ適用可能です。

オプション	説明
<b>Auto</b>	オプションの自動選択をアルゴリズムに任せます。現在、 <b>[ランダム]</b> が選択されます。
<b>Modulo</b>	データセットをフォールドに等分し、シードを基準としません。
<b>Random</b>	データを n フォールドのサブセットにランダムに分割します。大きなデータセットに最適です。

9. 相互検証を実行する場合は、**[フォールド フィールド]** をオンにして、相互検証フォールド インデックス割り当てを含むフィールドをドロップダウン リストから選択します。  
このフィールドは、**[N フォールド]** と **[フォールド割り当て]** に値が入力されていない場合のみ適用可能です。
10. **[最大反復回数]** をオンにし、実行する必要があるトレーニング反復回数を入力します。
11. **[目標イプシロン]** をオンにして、収束のしきい値を入力します。この値は 0 ~ 1 の間でなければなりません。  
目標値がこのしきい値に満たない場合、モデルは収束します。
12. **[ベータ イプシロン]** をオンにして、収束のしきい値を入力します。この値は 0 ~ 1 の間でなければなりません。  
目標値がこのしきい値に満たない場合、モデルは収束します。現在のベータ変化の L1 正則化がこのしきい値に満たない場合、収束の使用を検討してください。
13. 使用する正則化タイプを選択します。

正則化タイプ	説明
<b>LASSO (Least Absolute Shrinkage and Selection Operator)</b>	十分に重要と見なされる大きなラムダの値によって変数の小さなサブセットを選択します。相関する予測子変数がある場合は、適切な実行ができないことがあります。相関があるグループの変数は 1 つが選択され、それ以外はすべて除去されるためです。次元の高さによる制限もあります。モデルに含まれている変数がレコードよりも多い場合、LASSO では選択できる変数の数に制限が生じます。リッジ回帰 (Ridge Regression) にはこの制限がありません。モデルに含まれる変数の数が多い場合や、解が疎であることがわかっている場合は、LASSO が推奨されます。

## 正則化タイプ 説明

### Ridge Regression

すべての予測子変数を保持し、その係数を均等に縮退します。相関のある予測子変数が存在する場合、リッジ回帰 (Ridge Regression) では相関がある変数のグループ全体の係数を互いに均等になるように縮退します。相関がある予測子変数をモデルから除去したくない場合は、リッジ回帰を使用します。

### Elastic Net

LASSO とリッジ回帰 (Ridge Regression) を組み合わせたものであり、変数選択の機能を果たしながらも相関がある変数のグループ化の効果を保持 (同時に相関がある変数の係数を縮退) します。Elastic Net は、高次元による制限がなく、モデルに含まれている変数がレコードより多い場合でもすべての変数を評価できます。

予測モデリングにおける一般的な懸念事項は過剰適合です。これは、分析モデルが特定のデータセットと非常によく (または完全に) 一致しているため、追加のデータや将来の観測データへの適用時にうまく機能しないというものです。過剰適合の緩和に使用する方法の1つが正規化です。

14. **[アルファの値]** をオンにし、デフォルト値 **5** を使用しない場合は値を変更します。

アルファパラメータは、 $\lambda_1$  ペナルティと  $\lambda_2$  ペナルティの配分を制御します。有効な値の範囲は  $0 \sim 1$  です。値  $1.0$  は LASSO を表し、値  $0.0$  はリッジ回帰になります。以下の表に、アルファとラムダが正規化に及ぼす影響を示します。

lambda value	alpha value	Result
lambda == 0	alpha = any value	No regularization. alpha is ignored.
lambda > 0	alpha == 0	Ridge Regression
lambda > 0	alpha == 1	LASSO
lambda > 0	$0 < \alpha < 1$	Elastic Net Penalty

注：単独の等号は "とする" を意味する代入演算子であり、2つの等号を並べたものは "等しい" を意味する等価演算子です。

15. **[ラムダの値]** をオンにし、Linear Regression でラムダ値の計算にデフォルトの方法 (トレーニング データに基づく発見的方法) を使用しない場合は、値を指定します。

ラムダパラメータは適用される正則化の度合いを制御します。例えば、ラムダが  $0.0$  の場合は、正則化が適用されず、アルファパラメータは無視されます。

16. Linear Regression で完全な正則化の手順でモデルを計算するには、**[ラムダの最適値を探索]** をオンにします。

その場合は、ラムダが最大(意味のあるラムダの最大値、すなわちすべての係数を0にする最小の値)の状態を開始し、対数スケールでラムダを最小まで減少させ、ステップごとに正則化の度合いを小さくしていきます。返されるモデルの係数は、トレーニング中に決定されたラムダの最適値に対応したものになります。

17. トレーニングまたはバリデーションセットでそれ以上の改善がない場合に処理を終了するには、**[早期停止]** をオンにします。
18. **[探索する最大のラムダ]** をオンにし、ラムダ探索の処理で使用するラムダの最大数を入力します。
19. **[最大アクティブ予測子数]** をオンにし、計算時に使用する予測子の最大数を入力します。  
この値は、多数の予測子による高コストのモデル構築を防ぐために使用されます。
20. **[OK]** をクリックして、モデルと設定を保存するか、次のタブで操作を続行します。

## モデル出力

このタブには、適合されたモデルの評価に使用するメトリクスが表示されます。これらのフィールドは編集できません。**[トレーニング]**列には、必ずデータが含まれます。**[基本オプション]**タブでテストとトレーニングの分割を選択した場合は、**[テスト]**列にもデータが設定されます。ただし、**[高度なオプション]**タブで**N**フォールド検証を選択した場合を除きます。その場合は、**[Nフォールド]**列にデータが設定されます。

ジョブを実行すると、結果として得られたモデルが Spectrum™ Technology Platform サーバーに格納されます。出力を再生成するには **[出力]** ボタンをクリックし、出力全体を Machine Learning モデル管理ツールに表示するには **[モデルの詳細]** をクリックします。

## 出力ポート

Linear Regression ステージには、2つのオプション出力ポート、モデル スコア ポートとモデルメトリクス ポートが含まれています。これらのポートの機能は、ステージの基本オプションや高度なオプションの設定完了時の選択内容と入力情報によって決まります。例えば、**[高度なオプション]**タブで**[Nフォールド]**フィールドをオンにして**N**フォールド検証の実施を選択した場合、モデルメトリクスポートによって生成された出力メトリクスの**[Nフォールド]**列にデータが設定されます。また、**N**フォールド検証を実施しないことを選択した場合、**[Nフォールド]**列は空欄になります。同様に、**[基本オプション]**タブで**[入力データを記録]**フィールドをオンにすると、**[モデル スコア ポート]**がアクティブになります。

## モデル スコア ポート

[基本オプション] タブの [入力データを記録] チェック ボックスをオンにした場合、モデルの作成時に Linear Regression に予測値の計算が指示され、出力データでそのスコアに対して [Predicted\_Value] 列が追加されます。このポートには、どんな種類のシンクでも接続できます。例えば、Write to File ステージや Write to Null ステージなどです。

## モデル メトリクス ポート

モデル メトリクス ポートを使用するには、以下の手順に従います。

モデル メトリクス ポートを使用すると、モデル評価メトリクスをデータ ファイルに出力できます。これは、Spectrum™ Technology Platform の内部または外側で生成された多数のモデルを比較したり、メトリクスに関するその他のデータ処理タスクを実行したりするのに役立ちます。

1. Linear Regression ステージを使用しているデータフローを開きます。
2. Write to File ステージまたは別のデータ出力ステージを 2 番目の出力ポートに接続します。
3. ジョブを実行します。
4. ステップ 3 (23ページ) の代替: ステップ 2 で追加したシンク ステージに Linear Regression ステージを接続しているチャンネルにインスペクションポイントを追加します。そのためには、チャンネルを右クリックし、[インスペクション ポイントの追加] を選択します。その後、Enterprise Designer ツールバーの [現在のフローのインスペクション] ボタン (🔍) をクリックします。インスペクションが実行され、以下に示すような結果が表示されます。

Creation Time	Flow Name	Metrics	Model Name	Model Type	NFold	Test	Training
10/11/2018 2:37:00 AM	LinearRegressionTest1...	MSE	LinearRegressionTest1...	Linear Regression	0.0050019841090508	0.00321814876650744	1.31531988837612E-10
10/11/2018 2:37:00 AM	LinearRegressionTest1...	RMSE	LinearRegressionTest1...	Linear Regression	0.0707247064967455	0.0567287296394643	1.14687396359675E-05
10/11/2018 2:37:00 AM	LinearRegressionTest1...	Number of observations	LinearRegressionTest1...	Linear Regression	33	17	33
10/11/2018 2:37:00 AM	LinearRegressionTest1...	R2	LinearRegressionTest1...	Linear Regression	-0.366538763835926	0.0268719304981138	0.999999964065547
10/11/2018 2:37:00 AM	LinearRegressionTest1...	Mean residual deviance	LinearRegressionTest1...	Linear Regression	0.0050019841090508	0.00321814876650744	1.31531988837612E-10

# 5 - Logistic Regression

## このセクションの構成

---

はじめに	25
モデルのプロパティの定義	25
基本オプションの設定	26
高度なオプションの設定	26
モデル出力	30
出力ポート	30



## はじめに

Logistic Regression では、バイナリ目標と入力変数を使用するデータセットからモデルを作成して、機械学習を実行することができます。

モデルを作成するにはまず、**【モデルのプロパティ】** タブを設定する必要があります。**【基本オプション】** タブと **【高度なオプション】** タブには、ジョブを完了するために十分なデフォルト設定が指定されていますが、ニーズに合わせてその設定を変更することができます。続いてジョブを実行すると、結果として得られたモデルの限定版が **【モデル出力】** タブに表示されます。出力全体は、Machine Learning モデル管理ツールで確認できます。

## モデルのプロパティの定義

1. **【プライマリ ステージ】** > **【展開済みステージ】** > **【Machine Learning】** にある **【Logistic Regression】** ステージをクリックしてキャンバス上にドラッグします。さらに、そのステージを、データフロー内の該当位置に配置して、他のステージに接続します。

注: 入力ステージは、モデルの目標フィールドと入力変数フィールドの両方を含むデータソースでなければなりません。出力ステージは、**【基本オプション】** タブで**【入力データを記録】** オプションを選択しない限り不要です。Machine Learning モデル管理ツールとは独立して出力を取得する場合は、出力ステージを接続することもできます。

2. Logistic Regression ステージをダブルクリックして、**【Logistic Regression オプション】** ダイアログ ボックスを表示します。
3. デフォルトのモデル名を使用しない場合は、**【モデル名】** を入力します。
4. **【上書き】** チェックボックスをオンにして、既存モデルを新しいデータで上書きします。
5. **【目標フィールド】** ドロップダウンをクリックして "カテゴリ" を選択します。
6. モデルの **【説明】** を入力します。
7. モデルにデータを追加するフィールドの **【含める】** をクリックします。  
目標フィールドとして選択したフィールドは必ず含めてください。
8. **【モデル データ タイプ】** ドロップダウンを使用して、各入力フィールドを数値、カテゴリ値、日付と時刻のいずれのフィールドとして使うかを指定します。
9. **【OK】** をクリックして、モデルと設定を保存するか、次のタブで操作を続行します。

## 基本オプションの設定

1. 数値列を標準化し、平均が 0 で偏差が 1 になるようにするには、**【標準化】** をオンのままにします。  
標準化を使わない場合、実際の寄与の大きさではなく尺度の違いから、他の属性と比べて偏差が大きくなる変数によって左右される要素が結果に含まれる場合があります。
2. モデル予測 (スコア) を表す列を入力データに追加するには、**【入力データを記録】** をオンにします。
3. データがサンプル済みで、応答の平均が実態を反映していない場合は、**【プライア】** をオンにし、 $p(y=1)$  の事前確率をテキスト フィールドに入力します。
4. 欠落データの処理方法を指定するには、**【スキップ】** または **【平均値を補完】** をオンにします。後者のオプションを選択すると、欠落データの代わりに平均値が追加されます。
5. 入力データがトレーニングおよびテストのデータ サンプルにランダムに分割される場合は、1 ~ 100 の値を **【トレーニング データの比率】** に指定します。
6. ステップ 5 で指定した値を 100 から引いた値を **【テスト データの比率】** に入力します。
7. データフローを何度実行しても、データが必ず同じ方法でテスト データとトレーニング データに分割されるようにするには、**【テスト データ用シード】** に数値を入力します。フローを実行するたびにランダムな分割を行う場合は、このフィールドをオフにします。
8. **【OK】** をクリックして、モデルと設定を保存するか、次のタブで操作を続行します。

## 高度なオプションの設定

1. **【定数フィールドを無視】** をオンにすると、各レコードで値が同じフィールドがスキップされます。
2. **【p 値を計算】** をオンにすると、パラメータを予測するための p 値が計算されます。
3. モデルの作成時に共線列を自動的に削除するには、**【共線列を削除】** をオンのままにします。このオプションは、**【p 値を計算】** がオンになっている場合は常にオンにする必要があります。これにより、返されるモデルでは係数が 0 になります。
4. 定数項 (切片) をモデルに含めるには、**【定数項 (切片) を含める】** をオンにします。**【共線列を削除】** がオンの場合は、このオプションを必ずオンにする必要があります。
5. ドロップダウン リストから **【ソルバー】** を選択します。

ソルバー	説明
<b>Auto</b>	入力データとパラメータに基づいてソルバーが決定されます。
<b>CoordinateDescentNaive</b>	最も内側のループにおける循環座標降下法の共分散更新バージョンを使う IRLSM。
<b>CoordinateDescentNaive</b>	最も内側のループにおける循環座標降下法のネイティブ更新バージョンを使う IRLSM。
<b>IRLSM</b>	予測因子が少数のときの問題や、L1ペナルティによるラムダ検索の問題に最適です。
<b>L_BFGS</b>	多数の列が含まれるデータセットに最適です。

注: **CoordinateDescentNaive** および **CoordinateDescentNaive** は現時点で実験用です。

6. **[アルゴリズムのシード]** をオンにしてシード数を入力すると、データフローを何度実行しても、データが必ず同じ方法でテスト データとトレーニング データに分割されます。フローを実行するたびにランダムな分割を行う場合は、このフィールドをオフにします。
7. 相互検証を実行する場合は **[N フォールド]** をオンにし、フォールドの数を入力します。
8. 相互検証を実行する場合は、**[フォールド割り当て]** をオンにしてドロップダウン リストから選択します。

#### フォールド割り当て 説明

<b>Auto</b>	オプションの自動選択をアルゴリズムに任せます。現在、 <b>[ランダム]</b> が選択されます。
<b>Modulo</b>	データセットをフォールドに等分し、シードを基準としません。
<b>Random</b>	データを n フォールドのサブセットにランダムに分割します。大きなデータセットに最適です。
<b>Stratified</b>	分類問題の応答変数に基づいて、フォールドを層化します。データセットをトレーニング データとテスト データに分割する際に、観測値を複数

## フォールド割り当て 説明

このクラスからすべてのセットに均等に分散します。これは、クラスの数が多く、データセットが比較的小さい場合に便利です。

このフィールドは、**[N フォールド]** に値が入力済みで、**[フォールド フィールド]** が指定されていない場合にのみ適用可能です。

9. 相互検証を実行する場合は、**[フォールド フィールド]** をオンにして、相互検証フォールド インデックス割り当てを含むフィールドをドロップダウン リストから選択します。

このフィールドは、**[N フォールド]** と **[フォールド割り当て]** に値が入力されていない場合のみ適用可能です。

10. **[最大反復回数]** をオンにし、実行する必要があるトレーニング反復回数を入力します。
11. **[目標イプシロン]** をオンにして、収束のしきい値を入力します。この値は 0 ~ 1 の間でなければなりません。  
目標値がこのしきい値に満たない場合、モデルは収束します。
12. **[ベータ イプシロン]** をオンにして、収束のしきい値を入力します。この値は 0 ~ 1 の間でなければなりません。  
目標値がこのしきい値に満たない場合、モデルは収束します。現在のベータ変化の L1 正則化がこのしきい値に満たない場合、収束の使用を検討してください。
13. 使用する正則化タイプを選択します。

### 正則化タイプ 説明

#### LASSO (Least Absolute Shrinkage and Selection Operator)

十分に重要と見なされる大きなラムダの値によって変数の小さなサブセットを選択します。相関する予測子変数がある場合は、適切な実行ができないことがあります。相関があるグループの変数は 1 つが選択され、それ以外はすべて除去されるためです。次元の高さによる制限もあります。モデルに含まれている変数がレコードよりも多い場合、LASSO では選択できる変数の数に制限が生じます。リッジ回帰 (Ridge Regression) にはこの制限がありません。モデルに含まれる変数の数が多い場合や、解が疎であることがわかっている場合は、LASSO が推奨されます。

#### Ridge Regression

すべての予測子変数を保持し、その係数を均等に縮退します。相関のある予測子変数が存在する場合、リッジ回帰 (Ridge Regression) では相関がある変数のグループ全体の係数を互いに均等になるように縮退します。相関がある予測子変数をモデルから除去したくない場合は、リッジ回帰を使用します。

## 正則化タイプ 説明

**Elastic Net** LASSO とリッジ回帰 (Ridge Regression) を組み合わせたものであり、変数選択の機能を果たしながらも相関がある変数のグループ化の効果を保持 (同時に相関がある変数の係数を縮退) します。Elastic Net は、高次元による制限がなく、モデルに含まれている変数がレコードより多い場合でもすべての変数を評価できます。

予測モデリングにおける一般的な懸念事項は過剰適合です。これは、分析モデルが特定のデータセットと非常によく (または完全に) 一致しているため、追加のデータや将来の観測データへの適用時にうまく機能しないというものです。過剰適合を緩和するために使用される方法の1つが正則化です。

14. **[アルファの値]** をオンにし、デフォルト値 **5** を使用しない場合は値を変更します。

アルファパラメータは、 $\ell_1$  ペナルティと  $\ell_2$  ペナルティの配分を制御します。有効な値の範囲は  $0 \sim 1$  です。値  $1.0$  は LASSO を表し、値  $0.0$  はリッジ回帰になります。以下の表に、アルファとラムダが正規化に及ぼす影響を示します。

lambda value	alpha value	Result
lambda == 0	alpha = any value	No regularization. alpha is ignored.
lambda > 0	alpha == 0	Ridge Regression
lambda > 0	alpha == 1	LASSO
lambda > 0	$0 < \alpha < 1$	Elastic Net Penalty

注：単独の等号は "とする" を意味する代入演算子であり、2つの等号を並べたものは "等しい" を意味する等価演算子です。

15. **[ラムダの値]** をオンにし、Logistic Regression でラムダ値の計算にデフォルトの方法 (トレーニング データに基づく発見的方法) を使用しない場合は、値を指定します。

ラムダパラメータは適用される正則化の度合いを制御します。例えば、ラムダが  $0.0$  の場合は、正則化が適用されず、アルファパラメータは無視されます。

16. Logistic Regression で完全な正則化の手順でモデルを計算するには、**[ラムダの最適値を探索]** をオンにします。

その場合は、ラムダが最大 (意味のあるラムダの最大値、すなわちすべての係数を  $0$  にする最小の値) の状態で開始し、対数スケールでラムダを最小まで減少させ、ステップごとに正則化の度合いを小さくしていきます。

返されるモデルの係数は、トレーニング中に決定されたラムダの最適値に対応したものになります。

17. トレーニングまたはバリデーションセットでそれ以上の改善がない場合に処理を終了するには、**[早期停止]** をオンにします。
18. **[探索する最大のラムダ]** をオンにし、ラムダ探索の処理で使用するラムダの最大数を入力します。
19. **[最大アクティブ予測子数]** をオンにし、計算時に使用する予測子の最大数を入力します。この値は、多数の予測子による高コストのモデル構築を防ぐために使用されます。
20. **[OK]** をクリックして、モデルと設定を保存するか、次のタブで操作を続行します。

## モデル出力

このタブには、適合されたモデルの評価に使用するメトリクスが表示されます。これらのフィールドは編集できません。[トレーニング]列には、必ずデータが含まれます。[基本オプション]タブでテストとトレーニングの分割を選択した場合は、[テスト]列にもデータが設定されます。ただし、[高度なオプション]タブでNフォールド検証を選択した場合を除きます。その場合は、[Nフォールド]列にデータが設定されます。

ジョブを実行すると、結果として得られたモデルが Spectrum™ Technology Platform サーバーに格納されます。出力を再生成するには **[出力]** ボタンをクリックし、出力全体を Machine Learning モデル管理ツールに表示するには **[モデルの詳細]** をクリックします。

## 出力ポート

Logistic Regression ステージには、2つのオプション出力ポート、モデルスコアポートとモデルメトリクスポートが含まれています。これらのポートの機能は、ステージの基本オプションや高度なオプションの設定完了時の選択内容と入力情報によって決まります。例えば、**[高度なオプション]** タブで **[Nフォールド]** フィールドをオンにしてNフォールド検証の実施を選択した場合、モデルメトリクスポートによって生成された出力メトリクスの **[Nフォールド]** 列にデータが設定されます。また、Nフォールド検証を実施しないことを選択した場合、**[Nフォールド]** 列は空欄になります。同様に、**[基本オプション]** タブで **[入力データを記録]** フィールドをオンにすると、**[モデルスコアポート]** がアクティブになります。

## モデル スコア ポート

[基本オプション] タブの [入力データを記録] チェック ボックスをオンにした場合、モデルの作成時に Logistic Regression に予測値の計算が指示され、出力データでそのスコアに対して [Predicted\_Value]、[Probability\_of\_class\_A]、[Probability\_of\_class\_B] の各列が追加されます。このポートには、どんな種類のシンクでも接続できます。例えば、Write to File ステージや Write to Null ステージなどです。

## モデル メトリクス ポート

モデル メトリクス ポートを使用するには、以下の手順に従います。

モデル メトリクス ポートを使用すると、モデル評価メトリクスをデータ ファイルに出力できます。これは、Spectrum™ Technology Platform の内部または外側で生成された多数のモデルを比較したり、メトリクスに関するその他のデータ処理タスクを実行したりするのに役立ちます。

1. Logistic Regression ステージを使用しているデータフローを開きます。
2. Write to File ステージまたは別のデータ出力ステージを 2 番目の出力ポートに接続します。
3. ジョブを実行します。
4. ステップ 3 (31ページ) の代替: ステップ 2 (31ページ) で追加したシンク ステージに Logistic Regression ステージを接続しているチャンネルにインスペクション ポイントを追加します。そのためには、チャンネルを右クリックし、[インスペクション ポイントの追加] を選択します。その後、Enterprise Designer ツールバーの [現在のフローのインスペクション] ボタン (🔍) をクリックします。インスペクションが実行され、以下に示すような結果が表示されます。

Creation Time	Flow Name	Metrics	Model Name	Model Type	N Fold	Test	Training
10/11/2018 2:37:00 AM	LinearRegressionTest1...	MSE	LinearRegressionTest1...	Linear Regression	0.0050019841090508	0.00321814876650744	1.31531988837612E-10
10/11/2018 2:37:00 AM	LinearRegressionTest1...	RMSE	LinearRegressionTest1...	Linear Regression	0.0707247064967455	0.0567287296394643	1.14687396359675E-05
10/11/2018 2:37:00 AM	LinearRegressionTest1...	Number of observations	LinearRegressionTest1...	Linear Regression	33	17	33
10/11/2018 2:37:00 AM	LinearRegressionTest1...	R2	LinearRegressionTest1...	Linear Regression	-0.366538763835926	0.0268719304981138	0.999999964065547
10/11/2018 2:37:00 AM	LinearRegressionTest1...	Mean residual deviance	LinearRegressionTest1...	Linear Regression	0.0050019841090508	0.00321814876650744	1.31531988837612E-10

# 6 - 主成分分析

## このセクションの構成

---

はじめに	33
モデルのプロパティの定義	33
基本オプションの設定	34
高度なオプションの設定	34
モデル出力	35
出力ポート	35



## はじめに

主成分分析 (PCA) は、相関のある可能性のある変数群の観測データの集合を、主成分と呼ばれる線型相関のない変数の値の集合に変換する統計的な処理です。

モデルを作成するにはまず、**[モデルのプロパティ]** タブを設定する必要があります。**[基本オプション]** タブと **[高度なオプション]** タブには、ジョブを完了するために十分なデフォルト設定が指定されていますが、ニーズに合わせてその設定を変更することができます。ジョブを実行すると、最終的なモデルが、限定された形式で **[モデル出力]** タブに表示されます。完全な形式の出力を確認するには、Machine Learning モデル管理ツールを使用します。モデルの出力に問題がなければ、そのモデルを公開してスコアリング データフローで使用することができます。

## モデルのプロパティの定義

1. **[プライマリ ステージ]** > **[展開済みステージ]** > **[Machine Learning]** にある **[PCA Options]** ステージをクリックしてキャンバス上にドラッグします。さらに、そのステージを、データフロー内の該当位置に配置して、他のステージに接続します。

注：入力ステージは、モデルの主成分を含むデータ ソースでなければなりません。出力ステージは必要ありませんが、Machine Learning モデル管理ツールとは独立して出力を取得する場合は、出力ステージを接続することもできます。

2. **[PCA Options]** ステージをダブルクリックして、**[PCA オプション]** ダイアログ ボックスを表示します。
3. デフォルトのモデル名を使用しない場合は、**[モデル名]** を入力します。
4. オプション: **[上書き]** チェックボックスをオンにして、既存モデルを新しいデータで上書きします。
5. モデルに含める **[主要コンポーネント]** の数を入力します。
6. オプション: モデルの **[説明]** を入力します。
7. **[入力]** テーブルで、モデルにデータを追加するフィールドの **[含める]** をクリックします。
8. **[モデル データ タイプ]** ドロップダウンを使って、入力フィールドをカテゴリ値、日付と時刻、数値、文字列、ユニーク ID のいずれのフィールドとして使うかを指定します。
9. **[OK]** をクリックして、モデルと設定を保存するか、次のタブで操作を続行します。

## 基本オプションの設定

1. **[すべての因子レベルを使用]**を設定します。
  - 第1主成分をスキップするには、このオプションをオフのままにしておきます。この場合、データ内の分散が最大になります。
  - このチェックボックスをオンにすると、第1主成分が保持されます。
2. トレーニング データに対する適切な **[変換]** を選択します。

変形	説明
平均除去	各列の平均値を減算します。
スケール除去	各列の標準偏差による除算を行います。
なし	変換を行いません。
標準化	各列に対して、平均の減算とその範囲 (最大値と最小値の差) による除算を行います。
正規化	ゼロ平均と単位分散を使用します。これがデフォルトの変換です。

3. **[欠落データ]**の処理方法を指定するには、**[スキップ]**または**[平均値を補完]**をオンにします。後者では、欠落データの代わりに平均値が追加されます。
4. **[OK]** をクリックして、モデルと設定を保存するか、次のタブで操作を続行します。

## 高度なオプションの設定

1. **[定数フィールドを無視]** をオンにすると、各レコードで値が同じフィールドがスキップされます。
2. ドロップダウン リストから **[PCA メソッド]** を選択します。GLRM および Power は現時点で実験用であることに注意してください。

PCA メソッド	説明
<b>GLRM</b>	一般化された低ランク モデルのフィッティングを L2 損失関数によって正則化なしに行います。局所行列代数を使用して SVD を求めます。このオプションは、[基本オプション] タブで <b>[すべての因子レベルを使用]</b> をオンにしている場合にのみ有効になります。
<b>GramSVD</b>	Gram 行列の分散型計算を使用した後、JAMA パッケージを使用した局所 SVD を実行します。
<b>Power</b>	べき乗による反復法を使用して SVD を計算します。
<b>Randomized</b>	ランダム化された部分空間反復法を使用します。

3. トレーニングの反復回数を制限しない場合は、**[最大反復回数]** はオフのままにしておきます (デフォルトの状態)。トレーニングの反復回数を制限するには、このボックスをオンにして数値を入力します。
4. **[OK]** をクリックして、モデルと設定を保存するか、次のタブで操作を続行します。

## モデル出力

このタブには、適合されたモデルの評価に使用するメトリクスが表示されます。これらのフィールドは編集できません。

ジョブを実行すると、結果として得られたモデルが Spectrum™ Technology Platform サーバーに格納されます。出力を再生成するには **[出力]** ボタンをクリックし、出力全体を Machine Learning モデル管理ツールに表示するには **[モデルの詳細]** をクリックします。

## 出力ポート

主成分分析ステージには、1つのオプション出力ポート、モデル メトリクス ポートが含まれています。このポートの機能は、ステージの基本オプションや高度なオプションの設定完了時の選択内容と入力情報によって決まります。

## モデル メトリクス ポート

モデル メトリクス ポートを使用するには、以下の手順に従います。

モデル メトリクス ポートを使用すると、モデル評価メトリクスをデータ ファイルに出力できます。これは、Spectrum™ Technology Platform の内部または外側で生成された多数のモデルを比較したり、メトリクスに関するその他のデータ処理タスクを実行したりするのに役立ちます。

1. PCA ステージを使用しているデータフローを開きます。
2. Write to File ステージまたは別のデータ出力ステージを 2 番目の出力ポートに接続します。
3. ジョブを実行します。
4. ステップ 3 (36ページ) の代替: ステップ 2 (36ページ) で追加したシンク ステージに PCA ステージを接続しているチャンネルにインスペクションポイントを追加します。そのためには、チャンネルを右クリックし、[インスペクション ポイントの追加] を選択します。その後、Enterprise Designer ツールバーの [現在のフローのインスペクション] ボタン (🔍) をクリックします。インスペクションが実行され、以下に示すような結果が表示されます。

Creation Time	Cumulative Proportion	Flow Name	Model Name	Model Type	Principal Component	Proportion of Variance	Standard Deviation
10/11/2018 8:36:00 PM	0.120990707073471	PCA_MODEL_ASSESSMENT	PCA Model Assessment	PCA	PC1	0.120990707073471	1.73278529942543
10/11/2018 8:36:00 PM	0.163608702477588	PCA_MODEL_ASSESSMENT	PCA Model Assessment	PCA	PC2	0.0426179954041175	1.02840755055022
10/11/2018 8:36:00 PM	0.20114715020656	PCA_MODEL_ASSESSMENT	PCA Model Assessment	PCA	PC3	0.0375384477289716	0.965176862701583
10/11/2018 8:36:00 PM	0.236650281720107	PCA_MODEL_ASSESSMENT	PCA Model Assessment	PCA	PC4	0.0355031315135469	0.93864652798561
10/11/2018 8:36:00 PM	0.269754397170741	PCA_MODEL_ASSESSMENT	PCA Model Assessment	PCA	PC5	0.0331041154506347	0.90637880520786

# 7 - Random Forest Classification

## このセクションの構成

---

はじめに	38
モデルのプロパティの定義	38
基本オプションの設定	39
高度なオプションの設定	40
モデル出力	43
出力ポート	44

## はじめに

Random Forest Classification では、バイナリ目標と入力変数を使用するデータセットからモデルを作成して、機械学習を実行することができます。

モデルを作成するにはまず、[モデルのプロパティ] タブを設定する必要があります。[基本オプション] タブと [高度なオプション] タブには、ジョブを完了するために十分なデフォルト設定が指定されていますが、ニーズに合わせてその設定を変更することができます。続いてジョブを実行すると、結果として得られたモデルの限定版が [モデル出力] タブに表示されます。出力全体は、Machine Learning モデル管理ツールで確認できます。

注：詳細については、この [Distributed Random Forest \(DRF\)](#) の記事を参照してください。Random Forest Classification とそのオプションに関するその他の情報が記載されています。

## モデルのプロパティの定義

1. [プライマリ ステージ] > [展開済みステージ] > [Machine Learning] にある [Random Forest Classification] ステージをクリックしてキャンバス上にドラッグします。さらに、そのステージを、データフロー内の該当位置に配置して、他のステージに接続します。

注：入力ステージは、モデルの目標フィールドと入力変数フィールドの両方を含むデータソースでなければなりません。出力ステージは、[基本オプション] タブで [入力データを記録] オプションを選択しない限り不要です。Machine Learning モデル管理ツールとは独立して出力を取得する場合は、出力ステージを接続することもできます。

2. [Random Forest Classification] ステージをダブルクリックして、[Random Forest Classification オプション] ダイアログ ボックスを表示します。
3. デフォルトのモデル名を使用しない場合は、[モデル名] を入力します。
4. オプション: [上書き] チェックボックスをオンにして、既存モデルを新しいデータで上書きします。
5. [目標フィールド] ドロップダウンをクリックし、数値フィールドを選択します。
6. [多項レベル] をクリックし、目標フィールドをグループ化して分けることができるカテゴリの最大数を入力します。このオプションをオンにすると、[基本オプション] タブの [入力データを記録] オプションが無効になることに注意してください。

7. オプション: モデルの **[説明]** を入力します。
8. モデルにデータを追加するフィールドの **[含める]** をクリックします。  
目標フィールドとして選択したフィールドは必ず含めてください。
9. **[モデル データ タイプ]** ドロップダウンを使用して、各入力フィールドを数値、カテゴリ値、日付と時刻のいずれのフィールドとして使うかを指定します。
10. **[OK]** をクリックして、モデルと設定を保存するか、次のタブで操作を続行します。

## 基本オプションの設定

1. **[ツリーの数]** に、お使いのモデルでのツリー数の最大値を入力します。
2. **[最大深度]** を入力します。  
これは、モデルに含めるレベルの最大数を示します。
3. **[最小行数]** を入力します。  
これは、モデルに含める行 (またはレコード) の最小数を示します。
4. **[ビンの数]** を入力します。  
これは、ヒストグラムを構築したうえで最良のポイントで分割するビンの数を示します。
5. **[ビンの数 (最上位レベル)]** を入力します。  
これは、ルート レベルに必要なビンの最小数を示します。
6. **[ビンの数 (カテゴリ別)]** を入力します。  
これは、ヒストグラムを構築したうえで最良のポイントで分割するビンの最大数を示します。
7. **[サンプルレート]** をオンにし、各ツリーでサンプルとして使用される行の比率を入力します。  
0.0 ~ 1.0 の値を使用できます。
8. **[各ツリーの列サンプル レート]** をオンにし、各ツリーの列に対するサンプリング率を入力します。  
0.0 ~ 1.0 の値を使用できます。
9. **[各レベルの列数]** をオンにし、すべてのレベルでの列のサンプリングに対する相対変化量を入力します。  
有効な値の範囲は、1.0 から、選択した入力予測因子の数値までです。デフォルトは 1.0 です。
10. モデル予測 (スコア) を表す列を入力データに追加するには、**[入力データを記録]** をオンにします。

注：[モデルのプロパティ] タブで [多項レベル] をオンにした場合はこのオプションが無効になっています。

11. 入力データがトレーニングおよびテストのデータ サンプルにランダムに分割される場合は、1 ~ 100 の値を [トレーニング データの比率] に指定します。
12. ステップ 11 (40ページ) で指定した値を 100 から引いた値を [テスト データの比率] に入力します。
13. [テスト データ用シード] により、データフローを何度実行してもデータが必ず同じ方法でテスト データとトレーニング データに分割されるようになります。フローを実行するたびにランダムな分割を行う場合は、このフィールドをオフにします。
14. [OK] をクリックして、モデルと設定を保存するか、次のタブで操作を続行します。

## 高度なオプションの設定

1. [定数フィールドを無視] をオンにすると、各レコードで値が同じフィールドがスキップされます。
2. [クラスのバランスをとる] をオンにすると、クラス分布のバランスをとるために大多数のクラスでアンダーサンプリングが行われるか、少数のクラスでオーバーサンプリングが行われます。
3. [ヒストグラム タイプ] を選択します。

### ヒストグラム タイプ 説明

<b>Auto</b>	バケットが最小値から最大値まで (最大値 - 最小値)/N の刻み幅でビンングされます。このオプションで、最適な分割ポイントを見つけるために使用するヒストグラムのタイプを指定します。
<b>QuantilesGlobal</b>	各バケットに含める個体数を均等にします。個々の数値列 (二値以外) の nbins 個の分位を計算した後、2 つの分位に挟まれた各バケットに含める内容を均等に (残余はランダムに) 取捨選択して合計 nbins_top_level 個のビンを生成します。
<b>Random</b>	最小値から最大値までの N-1 個のポイントをサンプリングし、それらのポイントをソートしたリストから最適な分割ポイントを見つけます。



## ヒストグラム タイプ 説明

<b>RoundRobin</b>	すべてのヒストグラム タイプを (ツリーごとに 1 つずつ) 順に繰り返し使用します。
<b>UniformAdaptive</b>	個々のフィーチャーをビンニングして刻み幅 (個体数ではない) が均等のバケットを生成します。これは最速の方法ですが、分布に大きな偏りがあると分割が正確でなくなる可能性があります。

## 4. [カテゴリ別エンコーディング] を選択します。

### カテゴリ別エンコーディング 説明

<b>Auto</b>	自動的に 列挙型 エンコーディングを実行します。
<b>Binary</b>	<p>カテゴリを整数に変換してから 2 進数に変換し、その各桁を別々の列に割り当てます。次元数を減らしてデータをエンコードします (距離に歪みが生じます)。</p> <p>注: カテゴリ別のフィーチャーの列の数は 32 以下でなければなりません。</p>
<b>Eigen</b>	カテゴリ別のフィーチャーの $k$ 個の列についてのみ、ワンホット (one-hot) エンコーディング マトリックスを $k$ 次元固有空間に投影し続けます。
<b>列挙</b>	すべてのヒストグラム タイプを (ツリーごとに 1 つずつ) 順に繰り返し使用します。
<b>OneHotExplicit</b>	カテゴリごとに 1 つの列を生成し、列の各セルの値 "1" または "0" でその列のカテゴリが行に含まれているかどうかを表します。

- [アルゴリズムと N フォールドのシード] をオンにしてシード数を入力すると、データフローを何度実行しても、データが必ず同じ方法でテスト データとトレーニング データに分割されます。フローを実行するたびにランダムな分割を行う場合は、このフィールドをオフにします。
- 相互検証を実行する場合は、[N フォールド] をオンにして、フォールドの数を入力します。

7. 相互検証を実行する場合は、**[フォールド割り当て]** をオンにしてドロップダウン リストから選択します。

#### フォールド割り当て 説明

<b>Auto</b>	オプションの自動選択をアルゴリズムに任せます。現在、 <b>[ランダム]</b> が選択されます。
<b>Modulo</b>	データセットをフォールドに等分し、シードを基準としません。
<b>Random</b>	データを $n$ フォールドのサブセットにランダムに分割します。大きなデータセットに最適です。
<b>Stratified</b>	分類問題の応答変数に基づいて、フォールドを層化します。データセットをトレーニング データとテスト データに分割する際に、観測値を複数のクラスからすべてのセットに均等に分散します。これは、クラスの数が多く、データセットが比較的小さい場合に便利です。

このフィールドは、**[N フォールド]** に値が入力済みで、**[フォールド フィールド]** が指定されていない場合にのみ適用可能です。

8. 相互検証を実行する場合は、**[フォールド フィールド]** をオンにして、相互検証フォールド インデックス割り当てを含むフィールドをドロップダウン リストから選択します。

このフィールドは、**[N フォールド]** と **[フォールド割り当て]** に値が入力されていない場合のみ適用可能です。

9. **[停止の基準回数]** をオンにすると、指定した回数のトレーニングで **Stopping\_metric** オプションの改善が見られないとき、トレーニングの停止前に失敗したトレーニングの回数が入力されます。この機能を無効にするには、**0** を指定します。

この指標は **Validation** データに基づいて計算されます (提供されている場合)。そうでなければ、トレーニング データが使われます。

10. **[停止指標]** を選択して、新しいツリーの生成を終了するタイミングを決定します。

#### 停止指標 説明

**AUC**  
ROC 曲線下面積。

注: 二項モデルにのみ適用できます。

停止指標	説明
<b>Auto</b>	デフォルトは deviance です。
<b>Lifftopgroup</b>	上位 1%。
<b>Logloss</b>	対数損失
<b>Meanperclasserror</b>	平均誤分類率。
<b>Misclassification</b>	$(1 - (\text{正しい予測数} / \text{合計予測数})) * 100$ の値。
<b>MSE</b>	平均 2 乗誤差。予測変数の分散とバイアスを包含する誤差です。
<b>RMSE</b>	2 乗平均平方根誤差。モデルや評価関数によって予測された値 (サンプルや母集団の値) と実際に観測した値との差異を表します。MSE の平方根でもあります。

11. **[停止の基準許容値]** をオンにし、指標に基づく停止の相対許容誤差を指定する値を入力すると、改善がこの値未満の場合にトレーニングが終了します。  
このフィールドは、**[停止の基準回数]** をオンにしている場合にのみ有効になります。
12. **[最小分割改善]** をオンにし、2 乗誤差が低減したときに分割が行われるように最小の相対的な改善を指定する値を入力します。  
このオプションは、適切に実行すれば、過剰適合を減らす効果があります。最適な値は  $1e-10 \dots 1e-3$  の範囲でしょう。このフィールドは、**[停止の基準回数]** をオンにしている場合にのみ有効になります。
13. **[OK]** をクリックして、モデルと設定を保存するか、次のタブで操作を続行します。

## モデル出力

このタブには、適合されたモデルの評価に使用するメトリクスが表示されます。これらのフィールドは編集できません。**[トレーニング]** 列には、必ずデータが含まれます。**[基本オプション]** タブでテストとトレーニングの分割を選択した場合は、**[テスト]** 列にもデータが設定されます。ただし、**[高度なオプション]** タブで N フォールド検証を選択した場合を除きます。その場合は、**[N フォールド]** 列にデータが設定されます。

ジョブを実行すると、結果として得られたモデルが Spectrum™ Technology Platform サーバーに格納されます。出力を再生成するには **[出力]** ボタンをクリックし、出力全体を Machine Learning モデル管理ツールに表示するには **[モデルの詳細]** をクリックします。

## 出力ポート

Random Forest Classification ステージには、2つのオプション出力ポート、モデルスコアポートとモデルメトリクスポートが含まれます。これらのポートの機能は、ステージの基本オプションや高度なオプションの設定完了時の選択内容と入力情報によって決まります。例えば、**[高度なオプション]** タブで **[N フォールド]** フィールドをオンにして N フォールド検証の実施を選択した場合、モデルメトリクスポートによって生成された出力メトリクスの **[N フォールド]** 列にデータが設定されます。また、N フォールド検証を実施しないことを選択した場合、**[N フォールド]** 列は空欄になります。同様に、**[基本オプション]** タブで **[入力データを記録]** フィールドをオンにすると、**[モデルスコアポート]** がアクティブになります。

## モデルスコアポート

**[基本オプション]** タブの **[入力データを記録]** チェックボックスをオンにした場合、モデルの作成時に Random Forest Classification に予測値の計算が指示され、出力データでそのスコアに対して **[Predicted\_Value]**、**[Probability\_of\_class\_A]**、**[Probability\_of\_class\_B]** の各列が追加されます。このポートには、どんな種類のシンクでも接続できます。例えば、Write to File ステージや Write to Null ステージなどです。

注：このポートは、Random Forest Classification の多項モデルでは機能しません。

## モデルメトリクスポート

モデルメトリクスポートを使用するには、以下の手順に従います。

モデルメトリクスポートを使用すると、モデル評価メトリクスをデータファイルに出力できます。これは、Spectrum™ Technology Platform の内部または外側で生成された多数のモデルを比較したり、メトリクスに関するその他のデータ処理タスクを実行したりするのに役立ちます。

1. Random Forest Classification ステージを使用しているデータフローを開きます。
2. Write to File ステージまたは別のデータ出力ステージを 2 番目の出力ポートに接続します。

3. ジョブを実行します。
4. ステップ **3** (45ページ) の代替: ステップ **2** (44ページ) で追加したシンク ステージに **Random Forest Classification** ステージを接続しているチャンネルにインスペクション ポイントを追加します。そのためには、チャンネルを右クリックし、[インスペクションポイントの追加] を選択します。その後、Enterprise Designer ツールバーの [現在のフローのインスペクション] ボタン (🔍) をクリックします。インスペクションが実行され、以下に示すような結果が表示されます。

Creation Time	Flow Name	Metrics	Model Name	Model Type	N Fold	Test	Training
10/11/2018 2:37:00 AM	LinearRegressionTest1...	MSE	LinearRegressionTest1...	Linear Regression	0.0050019841090508	0.00321814876650744	1.31531988837612E-10
10/11/2018 2:37:00 AM	LinearRegressionTest1...	RMSE	LinearRegressionTest1...	Linear Regression	0.0707247064967455	0.0567287296394643	1.14687396359675E-05
10/11/2018 2:37:00 AM	LinearRegressionTest1...	Number of observations	LinearRegressionTest1...	Linear Regression	33	17	33
10/11/2018 2:37:00 AM	LinearRegressionTest1...	R2	LinearRegressionTest1...	Linear Regression	-0.366538763835926	0.0268719304981138	0.9999999964065547
10/11/2018 2:37:00 AM	LinearRegressionTest1...	Mean residual deviance	LinearRegressionTest1...	Linear Regression	0.0050019841090508	0.00321814876650744	1.31531988837612E-10

# 8 - Random Forest Regression

## このセクションの構成

---

はじめに	47
モデルのプロパティの定義	47
基本オプションの設定	48
高度なオプションの設定	49
モデル出力	52
出力ポート	52

## はじめに

Random Forest Regression では、バイナリ目標と入力変数を使用するデータセットからモデルを作成して、機械学習を実行することができます。

モデルを作成するにはまず、**[モデルのプロパティ]** タブを設定する必要があります。**[基本オプション]** タブと **[高度なオプション]** タブには、ジョブを完了するために十分なデフォルト設定が指定されていますが、ニーズに合わせてその設定を変更することができます。ジョブを実行すると、最終的なモデルが、限定された形式で **[モデル出力]** タブに表示されます。完全な形式の出力を確認するには、Machine Learning モデル管理ツールを使用します。

注：Random Forest Regression とそのオプションの詳細については、[Distributed Random Forest \(DRF\)](#) を参照してください。

## モデルのプロパティの定義

1. **[プライマリ ステージ]** / **[展開済みステージ]** / **[Machine Learning]** の下で、**[Random Forest Regression]** ステージをクリックしてキャンバス上にドラッグし、データフロー内の所望の位置に配置して、他のステージに接続します。

注：入力ステージは、モデルの目標フィールドと入力変数フィールドの両方を含むデータソースでなければなりません。出力ステージは、**[基本オプション]** タブで **[入力データを記録]** オプションを選択しない限り不要です。Machine Learning モデル管理ツールとは独立して出力を取得する場合は、出力ステージを接続することもできます。

2. **[Random Forest Regression]** ステージをダブルクリックして、**[ランダム フォレスト回帰オプション]** ダイアログ ボックスを表示します。
3. デフォルトのモデル名を使用しない場合は、**[モデル名]** を入力します。
4. オプション: **[上書き]** チェックボックスをオンにして、既存モデルを新しいデータで上書きします。
5. **[目標フィールド]** ドロップダウンをクリックし、数値フィールドを選択します。
6. オプション: モデルの **[説明]** を入力します。
7. データをモデルに追加したいそれぞれのフィールドで **[含める]** をクリックします。目標フィールドとして選択したフィールドは必ず含めてください。

8. **[モデル データ タイプ]** ドロップダウンを使用して、各入力フィールドを数値、カテゴリ値、日付と時刻のいずれのフィールドとして使うかを指定します。
9. **[OK]** をクリックして、モデルと設定を保存するか、次のタブで操作を続行します。

## 基本オプションの設定

1. **[ツリーの数]** に、お使いのモデルでのツリー数の最大値を入力します。デフォルトは 50 です。
2. **[最大深度]** を入力します。  
これは、モデルに含めるレベルの最大数を示します。デフォルトは 5 です。
3. **[最小行数]** を入力します。  
これは、モデルに含める行 (またはレコード) の最小数を示します。デフォルトは 10 です。
4. **[ビンの数]** を入力します。  
これは、ヒストグラムを構築したうえで最良のポイントで分割するビンの数を示します。デフォルトは 20 です。
5. **[ビンの数 (最上位レベル)]** を入力します。  
これは、ルート レベルで必要なビンの最小数を示します。デフォルトは 1024 です。
6. **[ビンの数 (カテゴリ別)]** を入力します。  
これは、ヒストグラムを構築したうえで最良のポイントで分割するビンの最大数を示します。デフォルトは 1024 です。
7. **[サンプルレート]** をオンにし、各ツリーでサンプルとして使用される行の比率を入力します。0.0 ~ 1.0 の値を使用できます。
8. **[各ツリーの列サンプル レート]** をオンにし、各ツリーの列に対するサンプリング率を入力します。  
0.0 ~ 1.0 の値を使用できます。
9. **[各レベルの列数]** をオンにし、すべてのレベルでの列のサンプリングに対する相対変化量を入力します。  
このオプションはデフォルトで 1.0 に設定されており、0.0 ~ 2.0 の値を使用できます。
10. モデル予測 (スコア) を表す列を入力データに追加するには、**[入力データを記録]** をオンにします。
11. 入力データがトレーニングおよびテストのデータ サンプルにランダムに分割される場合は、1 ~ 100 の値を **[トレーニング データの比率]** に指定します。



12. ステップ **11** (48ページ) で指定した値を 100 から引いた値を **[テスト データの比率]** に入力します。
13. **[テスト データ用シード]** により、データフローを何度実行してもデータが必ず同じ方法でテスト データとトレーニング データに分割されるようになります。フローを実行するたびにランダムな分割を行う場合は、このフィールドをオフにします。
14. **[OK]** をクリックして、モデルと設定を保存するか、次のタブで操作を続行します。

## 高度なオプションの設定

1. **[定数フィールドを無視]** をオンにすると、各レコードで値が同じフィールドがスキップされます。
2. **[ヒストグラム タイプ]** を選択します。

### ヒストグラム タイプ 説明

<b>Auto</b>	バケットが最小値から最大値まで (最大値 - 最小値)/N の刻み幅でビニングされます。このオプションで、最適な分割ポイントを見つけるために使用するヒストグラムのタイプを指定します。
<b>QuantilesGlobal</b>	各バケットに含める個体数を均等にします。個々の数値列 (二値以外) の nbins 個の分位を計算した後、2つの分位に挟まれた各バケットに含める内容を均等に (残余はランダムに) 取捨選択して合計 nbins_top_level 個のビンを生成します。
<b>Random</b>	最小値から最大値までの N-1 個のポイントをサンプリングし、それらのポイントをソートしたリストから最適な分割ポイントを見つけます。
<b>RoundRobin</b>	すべてのヒストグラム タイプを (ツリーごとに 1つずつ) 順に繰り返し使用します。
<b>UniformAdaptive</b>	個々のフィーチャーをビニングして刻み幅 (個体数ではない) が均等のバケットを生成します。これは最速の方法ですが、分布に大きな偏りがあると分割が正確でなくなる可能性があります。

3. **[カテゴリ別エンコーディング]** を選択します。

## カテゴリ別エンコーディング

### 説明

#### Auto

自動的に 列挙型 エンコーディングを実行します。

#### Binary

カテゴリを整数に変換してから 2 進数に変換し、その各桁を別々の列に割り当てます。次元数を減らしてデータをエンコードします (距離に歪みが生じます)。

注: カテゴリ別のフィーチャーの列の数は 32 以下でなければなりません。

#### Eigen

カテゴリ別のフィーチャーの  $k$  個の列についてのみ、ワンホット (one-hot) エンコーディング マトリックスを  $k$  次元固有空間に投影し続けます。

#### 列挙

すべてのヒストグラム タイプを (ツリーごとに 1 つずつ) 順に繰り返し使用します。

#### OneHotExplicit

カテゴリごとに 1 つの列を生成し、列の各セルの値 "1" または "0" でその列のカテゴリが行に含まれているかどうかを表します。

4. [アルゴリズムと N フォールドのシード] をオンにしてシード数を入力すると、データフローを何度実行しても、データが必ず同じ方法でテスト データとトレーニング データに分割されます。フローを実行するたびにランダムな分割を行う場合は、このフィールドをオフにします。
5. 相互検証を実行する場合は [N フォールド] をオンにし、フォールドの数を入力します。
6. 相互検証を実行する場合は、[フォールド割り当て] をオンにしてドロップダウン リストから選択します。

## フォールド割り当て

### 説明

#### Auto

オプションの自動選択をアルゴリズムに任せます。現在、[ランダム] が選択されます。

#### Modulo

データセットをフォールドに等分し、シードを基準としません。

#### Random

データを  $n$  フォールドのサブセットにランダムに分割します。大きなデータセットに最適です。

このフィールドは、**[N フォールド]** に値が入力済みで、**[フォールド フィールド]** が指定されていない場合にのみ適用可能です。

- 相互検証を実行する場合は、**[フォールド フィールド]** をオンにして、相互検証フォールド インデックス割り当てを含むフィールドをドロップダウン リストから選択します。

このフィールドは、**[N フォールド]** と **[フォールド割り当て]** に値が入力されていない場合のみ適用可能です。

- [停止の基準回数]** をオンにすると、指定した回数のトレーニングで **Stopping\_metric** オプションの改善が見られないとき、トレーニングの停止前に失敗したトレーニングの回数が入力されます。この機能を無効にするには、**0** を指定します。

この指標は **Validation** データに基づいて計算されます (提供されている場合)。そうでなければ、**トレーニング データ** が使われます。

- [停止指標]** を選択して、新しいツリーの生成を終了するタイミングを決定します。

停止指標	説明
<b>Auto</b>	デフォルトは <b>deviance</b> です。
<b>deviance</b>	平均残差逸脱度 (MSE)。
<b>MAE</b>	平均絶対誤差。2 つの連続変数の間の差異です。
<b>MSE</b>	平均 2 乗誤差。予測変数の分散とバイアスを包含する誤差です。
<b>RMSE</b>	2 乗平均平方根誤差。モデルや評価関数によって予測された値 (サンプルや母集団の値) と実際に観測した値との差異を表します。MSE の平方根でもあります。
<b>RMSLE</b>	2 乗対数平均平方根誤差。予測値と実測値の比率を表します。

- [停止の基準許容値]** をオンにし、指標に基づく停止の相対許容誤差を指定する値を入力すると、改善がこの値未満の場合にトレーニングが終了します。

- [最小分割改善]** をオンにし、2 乗誤差が低減したときに分割が行われるように最小の相対的な改善を指定する値を入力します。

このオプションは、適切に実行すれば、過剰適合を減らす効果があります。最適な値は **1e-10...1e-3** の範囲でしょう。このフィールドは、**[停止の基準回数]** をオンにしている場合のみ有効になります。

- [OK]** をクリックして、モデルと設定を保存するか、次のタブで操作を続行します。

## モデル出力

このタブには、適合されたモデルの評価に使用するメトリクスが表示されます。これらのフィールドは編集できません。[トレーニング]列には、必ずデータが含まれます。[基本オプション]タブでテストとトレーニングの分割を選択した場合は、[テスト]列にもデータが設定されます。ただし、[高度なオプション]タブでNフォールド検証を選択した場合を除きます。その場合は、[Nフォールド]列にデータが設定されます。

ジョブを実行すると、結果として得られたモデルが Spectrum™ Technology Platform サーバーに格納されます。出力を再生成するには [出力] ボタンをクリックし、出力全体を Machine Learning モデル管理ツールに表示するには [モデルの詳細] をクリックします。

## 出力ポート

Random Forest Regression ステージには、2つのオプション出力ポート、モデルスコアポートとモデルメトリクスポートが含まれます。これらのポートの機能は、ステージの基本オプションや高度なオプションの設定完了時の選択内容と入力情報によって決まります。例えば、[高度なオプション]タブで [Nフォールド] フィールドをオンにしてNフォールド検証の実施を選択した場合、モデルメトリクスポートによって生成された出力メトリクスの [Nフォールド]列にデータが設定されます。また、Nフォールド検証を実施しないことを選択した場合、[Nフォールド]列は空欄になります。同様に、[基本オプション]タブで [入力データを記録] フィールドをオンにすると、[モデルスコアポート]がアクティブになります。

## モデルスコアポート

[基本オプション]タブの [入力データを記録] チェックボックスをオンにした場合、モデルの作成時に Random Forest Regression に予測値の計算が指示され、出力データでそのスコアに対して [Predicted\_Value]列が追加されます。このポートには、どんな種類のシンクでも接続できます。例えば、Write to File ステージや Write to Null ステージなどです。

## モデル メトリクス ポート

モデル メトリクス ポートを使用するには、以下の手順に従います。

モデル メトリクス ポートを使用すると、モデル評価メトリクスをデータ ファイルに出力できます。これは、Spectrum™ Technology Platform の内部または外側で生成された多数のモデルを比較したり、メトリクスに関するその他のデータ処理タスクを実行したりするのに役立ちます。

1. Random Forest Regression ステージを使用しているデータフローを開きます。
2. Write to File ステージまたは別のデータ出力ステージを 2 番目の出力ポートに接続します。
3. ジョブを実行します。
4. ステップ 3 (53ページ) の代替: ステップ 2 (53ページ) で追加したシンク ステージに Random Forest Regression ステージを接続しているチャンネルにインスペクション ポイントを追加します。そのためには、チャンネルを右クリックし、[インスペクション ポイントの追加] を選択します。その後、Enterprise Designer ツールバーの [現在のフローのインスペクション] ボタン (🔍) をクリックします。インスペクションが実行され、以下に示すような結果が表示されます。

Creation Time	Flow Name	Metrics	Model Name	Model Type	N Fold	Test	Training
10/11/2018 2:37:00 AM	LinearRegressionTest1...	MSE	LinearRegressionTest1...	Linear Regression	0.0050019841090508	0.00321814876650744	1.31531988837612E-10
10/11/2018 2:37:00 AM	LinearRegressionTest1...	RMSE	LinearRegressionTest1...	Linear Regression	0.0707247064967455	0.0567287296394643	1.14687396359675E-05
10/11/2018 2:37:00 AM	LinearRegressionTest1...	Number of observations	LinearRegressionTest1...	Linear Regression	33	17	33
10/11/2018 2:37:00 AM	LinearRegressionTest1...	R2	LinearRegressionTest1...	Linear Regression	-0.366538763835926	0.0268719304981138	0.999999964065547
10/11/2018 2:37:00 AM	LinearRegressionTest1...	Mean residual deviance	LinearRegressionTest1...	Linear Regression	0.0050019841090508	0.00321814876650744	1.31531988837612E-10

# 9 - Machine Learning

## モデル管理

### このセクションの構成

---

Machine Learning モデル管理へのアクセス	55
モデル評価	56
ビニング管理	64

## Machine Learning モデル管理へのアクセス

Machine Learning モデル管理には、次の 3 つの方法でアクセスできます。

- Spectrum™ Technology Platform ようこそページを使用します。
  - Web ブラウザを起動し、次の Spectrum™ Technology Platform の Welcome ページを開きます。

サーバー名:ポート

例えば、Spectrum™ Technology Platform が "myspectrumplatform" という名前のコンピュータにインストールされており、デフォルト ポート 8080 を使用している場合は、次のアドレスに移動します。

myspectrumplatform:8080

- **[Spectrum Machine Learning]** をクリックします。
- **[Machine Learning モデル管理を開く]** をクリックします。
- いずれかのモデル構築ステージから **[モデルの詳細についてはここをクリック]** をクリックします。
- Web ブラウザを使用して次の手順を実行します。

- Web ブラウザを起動し、以下の Spectrum™ Technology Platform の Machine Learning モデル管理ページを開きます。

サーバー名:ポート/machinelearning

例えば、Spectrum™ Technology Platform が "myspectrumplatform" という名前のコンピュータにインストールされており、デフォルト ポート 8080 を使用している場合は、次のアドレスに移動します。

myspectrumplatform:8080/machinelearning

- 有効な Spectrum™ Technology Platform ユーザ名とパスワードを入力します。

## モデル評価






### モデル評価の概要

Machine Learning モデル管理の [モデル評価] タブには、Spectrum™ Technology Platform サーバー上にある機械学習モデルの全一覧が表示されます。テキスト ボックスに文字列を入力することによって、この一覧にフィルタを適用することができます。その文字列によって、テーブルのすべてのフィールドが検索されます。

これらのモデルに対して、複数の操作が実行できます。モデルのインポート、エクスポート、エクスポート、アンエクスポート、削除が可能です。エクスポートされたモデルは **Java Model Scoring** ステージで、機械学習モデルの適合を行った時に作成された式を使用して、新しいデータをスコアリングするために使用されます。また、各モデルの詳細情報を表示できます。詳細情報は、データを表示するモデルのタイプによって異なります。最後に、同じタイプの任意の 2 つのモデルを比較できます。比較を実行すると、比較する各モデルに対して [モデルの詳細] タブに表示されるのと同じ情報が、左右に並んで表示されます。

### モデル評価の操作

モデルを選択して該当するボタンをクリックすることによって、以下の操作を実行します。

	モデル出力の詳細を表示します。K-Means Clustering ステージと Logistic Regression ステージからのこの情報は、[モデル出力] タブの [モデルの詳細についてはここをクリック] をクリックすることによっても参照できます。
	モデルを比較します。
	モデルを特定のパスからインポートします。必要に応じて、同じ名前の既存のモデルを上書きするかどうかを選択します。
	モデルを特定のパスにエクスポートします。必要に応じて、同じ名前の既存のモデルを上書きするかどうかを選択します。
	モデルをエクスポートして、Java Model Scoring ステージで使用できるようにします。エクスポートされていないモデルを、スコアリングに使用することはできません。





モデルをアンエクスポートします。



モデルを削除します。

注：エクスポートされているモデルを削除することはできません。ただし現時点では、他のユーザのモデルを削除できないようにするためのセキュリティ機能は装備されていません。

## [モデルの詳細] タブ

[モデルの詳細] 画面には、すべてのモデルに関する以下の情報が表示されます。

モデル名	モデル名
モデル タイプ	機械学習モデルのタイプ
ユーザ	モデルを作成したユーザのユーザ名
説明	モデルの説明 (作成時に記述された場合)
ステータス	モデルがエクスポートされているかどうか
データフロー名	モデルを生成したデータフローの名前
作成時間	モデルが作成された日時

モデル タイプに応じて、その他の詳細情報が表示されます。

### K-Means Clustering の詳細情報

[モデルの詳細] 画面には、K-Means Clustering モデルに関する以下の情報が表示されます。

#### モデル サマリ

以下の項目に対するトレーニング データを提供します。

- 行数
- クラスタ数
- カテゴリ列数
- 反復回数
- クラスタ内平方和
- 総平方和
- クラスタ間平方和

### メトリクス

以下の項目に対するトレーニング、テスト、N フォールド データを提供します。

- 総クラス内平方和
- 総平方和
- クラス間平方和

### セントロイド統計

各中心点 (セントロイド) に対する以下のトレーニング、テスト、N フォールド データを提供します。

- サイズ
- クラス内平方和

### クラス平均

各中心点の詳細情報を提供します。内容は入力データによって異なります。クラスとは、特定のクラスリング アルゴリズムに基づいて類似と識別された、データ セットからのオブザベーションのグループです。

### 標準化されたクラス平均

各中心点の正規化情報を提供します。内容は入力データによって異なります。

## Logistic Regression の詳細情報

[モデルの詳細] 画面には、Logistic Regression モデルに関する以下の情報が表示されます。

### メトリクス

以下の項目に対するトレーニング、テスト、N フォールド データを提供します。

- 平均二乗誤差 (MSE)
- 平方根平均二乗誤差 (RMSE)
- オブザベーション数
- R の二乗 (R<sup>2</sup>)
- 対数損失 (Logloss)
- 曲線下面積 (AUC)
- 適合率-再現率曲線下面積 (PR AUC)
- ジニ係数
- クラスあたり平均誤差
- 赤池情報量基準 (AIC)
- ラムダ

- 残差逸脱度
- Null 逸脱度
- Null 自由度
- 残差自由度

### 最大メトリクスしきい値

以下のメトリクスを使用するトレーニング、テスト、N フォールド データに対する、トレーニング最大メトリクスしきい値を提供します。

- f1 最大値
- f2 最大値
- f0point5 最大値
- 最大正確度
- 最大適合率
- 最大再現率
- 最大特異度
- absolute\_mcc 最大値
- min\_per\_class\_accuracy 最大値
- mean\_per\_class\_accuracy 最大値

### 混同行列

真 (true) の値が既知の一連のトレーニング、テスト、N フォールド データに対するモデルのパフォーマンスを表します。

### 標準化係数チャート

入力がどれだけ変化すると目標が変化するかを表す、係数の相対値を提供することによって、最も重要な予測因子を示します。

### GLM 係数

指数分布に従う結果の回帰モデルを推定する、一般化線形モデル (GLM: Generalized Linear Model) の係数を示します。

### AUC 曲線

曲線下面積 (AUC)。使用モデルの中で、トレーニング、テスト、N フォールド データを使用して最も正確にクラスを予測するものを判定します。

### リフト/ゲイン曲線

トレーニング、テスト、N フォールド データを使用してバイナリ分類モデルの予測能力を評価します。

### Logistic Regression の詳細情報

[モデルの詳細] 画面には、Linear Regression モデルに関する以下の情報が表示されます。

#### メトリクス

以下の項目に対するトレーニング、テスト、N フォールド データを提供します。

- 平均二乗誤差 (MSE)
- 平方根平均二乗誤差 (RMSE)
- オブザベーション数
- R の二乗 (R2)
- 平均残差逸脱度
- 平均絶対誤差 (MAE)
- 二乗対数平均平方根誤差 (RMSLE)
- 赤池情報量基準 (AIC)
- ラムダ
- 残差逸脱度
- Null 逸脱度
- Null 自由度
- 残差自由度

#### 標準化係数チャート

特定の予測係数値の変化により目標値がどれだけ変化 (正または負の変化) するかを表す係数の相対値を提供することによって、最も重要な予測因子を示します。さらに、モデルの上位 25 の係数をグラフで示します。

#### GLM 係数

指数分布に従う結果の回帰モデルを推定する、一般化線形モデル (GLM: Generalized Linear Model) の係数を示します。

### Random Forest Regression の詳細情報

[モデルの詳細] 画面には、Random Forest Regression モデルに関する以下の情報が表示されます。

### メトリクス

以下の項目に対するトレーニング、テスト、N フォールド データを提供します。

- 平均二乗誤差 (MSE)
- 平方根平均二乗誤差 (RMSE)
- オブザベーション数
- R の二乗 (R<sup>2</sup>)
- 平均残差逸脱度
- 平均絶対誤差 (MAE)
- 二乗対数平均平方根誤差 (RMSLE)

### 変数重要度

以下のメトリクスを使用して、各変数の重要度の値を提供します。

- 相対的重要度
- 小数点以下桁数を含む重要度
- 比率

さらに、モデルの上位 25 の変数をグラフで示します。

### Random Forest Classification の詳細情報 — 二項

[モデルの詳細] 画面には、Random Forest Classification の二項モデルに関する以下の情報が表示されます。

### メトリクス

以下の項目に対するトレーニング、テスト、N フォールド データを提供します。

- 平均二乗誤差 (MSE)
- 平方根平均二乗誤差 (RMSE)
- オブザベーション数
- R の二乗 (R<sup>2</sup>)
- Logloss
- 曲線下面積 (AUC)
- 適合率-再現率曲線下面積 (PR AUC)
- ジニ
- クラスあたり平均誤差

### 最大メトリクスしきい値

以下のメトリクスを使用するトレーニング、テスト、N フォールド データに対する、トレーニング最大メトリクスしきい値を提供します。

- f1 最大値
- f2 最大値
- f0point5 最大値
- 最大正確度
- 最大適合率
- 最大再現率
- 最大特異度
- absolute\_mcc 最大値
- min\_per\_class\_accuracy 最大値
- mean\_per\_class\_accuracy 最大値

### 混同行列

真 (true) の値が既知の一連のトレーニング、テスト、N フォールド データに対するモデルのパフォーマンスを表します。

### 変数重要度

以下のメトリクスを使用して、各変数の重要度の値を提供します。

- 相対的重要度
- 小数点以下桁数を含む重要度
- 比率

さらに、モデルの上位 25 の変数をグラフで示します。

### AUC 曲線

曲線下面積 (AUC)。使用モデルの中で、トレーニング、テスト、N フォールド データを使用して最も正確にクラスを予測するものを判定します。

### リフト/ゲイン曲線

トレーニング、テスト、N フォールド データを使用してバイナリ分類モデルの予測能力を評価します。

## Random Forest Classification の詳細情報 — 多項

[モデルの詳細] 画面には、Random Forest Classification の多項モデルに関する以下の情報が表示されます。

## メトリクス

以下の項目に対するトレーニング、テスト、N フォールド データを提供します。

- 平均二乗誤差 (MSE)
- 平方根平均二乗誤差 (RMSE)
- オブザベーション数
- R の二乗 (R<sup>2</sup>)
- Logloss
- クラスあたり平均誤差

## 混同行列

真 (true) の値が既知の一連のトレーニング、テスト、N フォールド データに対するモデルのパフォーマンスを表します。

## 変数重要度

以下のメトリクスを使用して、各変数の重要度の値を提供します。

- 相対的重要度
- 小数点以下桁数を含む重要度
- 比率

さらに、モデルの上位 25 の変数をグラフで示します。

## 主成分分析の詳細情報

[モデルの詳細] 画面には、主成分分析 (PCA) モデルに関する以下の情報が表示されます。

## コンポーネントの重要度

主要コンポーネントを、以下のメトリクスに基づき重要度順に表示します。

- 標準偏差
- 寄与率
- 累積寄与率

## 回転

変数負荷量の行列をグラフで示します。コンポーネントのスコアを算出するには、正規化された元の各変数にこの重みを掛ける必要があります。

## ビンニング管理






### ビンニング管理の概要

Machine Learning モデル管理の [ビンニング管理] タブには、Spectrum™ Technology Platform サーバー上にあるビンニングの全一覧が表示されます。テキスト ボックスに文字列を入力することによって、この一覧にフィルタを適用することができます。その文字列によって、テーブルのすべてのフィールドが検索されます。

ビンニングに対して、複数の操作を実行できます。ビンニングのインポート、エクスポート、エクスポート、アンエクスポート、削除が可能です。エクスポートしたビンニングは、Binning Lookup ステージによって、過去に定義したビンニングを新しいデータに適用するために使用されます。

### ビンニング管理操作

ビンニングを選択して該当するボタンをクリックすることによって、以下の操作を実行します。

	ビンニングをインポートします。同じ名前の既存のビンニングを上書きするかどうかを適宜選択します。
	ビンニングをエクスポートします。同じ名前の既存のビンニングを上書きするかどうかを適宜選択します。
	Binning Lookup ステージで使用できるようにビンニングを公開 (エクスポート) します。エクスポートされていないビンニングを、検索に使用することはできません。
	ビンニングをアンエクスポートします。
	ビンニングを削除します。エクスポートされているビンニングは、削除できません。 注: ユーザは、いずれのユーザが作成したビンニングも削除可能です。現時点では、ユーザ固有の権限はありません。



# 10 - Data Science Demonstration Flows

## このセクションの構成

---

はじめに	66
教師あり学習: 貸付返済不能予測	66
教師なし学習: セグメンテーション	67

## はじめに

Spectrum Data Science には、Machine Learning モジュールや Analytics Scoring モジュール、またモデリング用のデータを準備するモジュールが含まれています。これらのデモは、データの準備、モデリング、モデルスコアリングの例を示しています。段階的な手順に従って独自のデータフローを作成したり、用意されているデータフローをリファレンスとして使用したりできます。

## 教師あり学習: 貸付返済不能予測

### ■ 教師あり学習のデモをダウンロードする

Data Science の教師あり学習デモは、Lending Club データを使用して貸付返済不能予測を実施します。このデモでは、Spectrum™ Technology Platform の Data Science ソリューションの機能を Enterprise Designer で示すための複数のファイルを利用します。

Spectrum\_DataScience\_Supervised\_Learning.zip には、以下のファイルが含まれています。

- Spectrum\_DataScience\_Supervised\_Learning.pdf — 単一カテゴリザのデータフロー、スコアリング データフロー、すべてのサポートを作成して使用する方法を紹介しているドキュメントです。
- Data.zip — 必須の入力ファイル、テストファイル、トレーニングファイルが付属のデータフローごとに用意されています。
  - loan.csv
  - LoanStats\_2016Q1.csv
  - LoanStats\_2016Q2.csv
  - LoanStats\_2016Q3.csv
  - testData.txt
  - testDataCollege.txt
  - testDataStable.txt
  - testDataThankful.txt
  - trainData.txt
  - trainDataCollege.txt
  - trainDataStable.txt
  - trainDataThankful.txt

- training.xml
- trainingCollege.xml
- trainingStable.xml
- trainingThanks.xml
- Lending\_Club\_Demo\_DF\_(V12.1).zip — Spectrum™ Technology Platform 12.1 用のデータフロー
  - LendingClub\_2007\_2016Q12\_v121\_MultipleCategorizers.df
  - LendingClub\_2007\_2016Q1Q2\_v121\_SingleCategorizer.df
  - LendingClub\_2016Q3\_v121\_SingleCategorizer\_Scoring.df
- Lending\_Club\_Demo\_DF\_(V12.2).zip — Spectrum™ Technology Platform 12.2 用のデータフロー
  - LendingClub\_2007\_2016Q12\_v122\_MultipleCategorizers.df
  - LendingClub\_2007\_2016Q1Q2\_v122\_SingleCategorizer.df
  - LendingClub\_2016Q3\_v122\_SingleCategorizer\_Scoring.df
- ReadMe.txt — これまでに述べたファイルに関する大まかな説明と手順です。

ドキュメントの手順ごとの詳細な説明に従って独自のデータフローを作成できます。付属のデータフローを参考にして、各ステージおよびデータフローを全体としてどのように完成させればよいか確認することもできます。

## 教師なし学習: セグメンテーション

### 教師なし学習のデモをダウンロードする

Data Science の教師なし学習デモは、Consumer Expenditure データを使用してセグメンテーションを実施します。このデモでは、Spectrum™ Technology Platform の Data Science ソリューションの機能を Enterprise Designer で示すための複数のファイルを利用します。

Spectrum\_DataScience\_Unsupervised\_Learning.zip には、以下のファイルが含まれています。

- Spectrum\_DataScience\_Unsupervised\_Learning.pdf — プライマリ データフロー、サブフロー、スコアリング データフロー、およびすべてのサポート ファイルを作成および使用する方法を紹介しているドキュメント
- Data.zip — 付属の各データフロー用の必須の入力ファイルと出力ファイル
  - Input フォルダ — 付属の各データフロー用の必須の入力ファイル

- **Output** フォルダ — 付属の各データフロー用の必須の出力ファイル
- **PythonBased** フォルダ — プライマリ データフローの **Group Statistics** および **Transformer** ステージの代替としてオプションの **Python** 処理を使用するための必須の入力ファイルと出力ファイル
- **Consumer\_Expenditure\_Demo\_DF\_(v12.1).zip** — **Spectrum™ Technology Platform 12.1** 用のデータフロー
  - `ConsumerExpenditure_v121_sampleandcluster.df`
  - `ConsumerExpenditure_v121_sampleandcluster_subflow.df`
  - `ConsumerExpenditure_v121_score.df`
  - `ConsumerExpenditure_v121_subflow.df`
  - **PythonBased** フォルダ — プライマリ データフローの **Group Statistics** および **Transformer** ステージの代替としてオプションの **Python** 処理を使用するための、必須のデータフロー、プロセスフロー、バッチ スクリプト、**Python** スクリプト、およびドキュメント
- **Consumer\_Expenditure\_Demo\_DF\_(v12.2).zip** — **Spectrum™ Technology Platform 12.2** 用のデータフロー
  - `ConsumerExpenditure_v122_sampleandcluster.df`
  - `ConsumerExpenditure_v122_sampleandcluster_subflow.df`
  - `ConsumerExpenditure_v122_score.df`
  - `ConsumerExpenditure_v122_subflow.df`
  - **PythonBased** フォルダ — プライマリ データフローの **Group Statistics** および **Transformer** ステージの代替としてオプションの **Python** 処理を使用するための、必須のデータフロー、プロセスフロー、バッチ スクリプト、**Python** スクリプト、およびドキュメント
- `ReadMe.txt` — これまでに述べたファイルに関する大まかな説明と手順です。

ドキュメントの手順ごとの詳細な説明に従って独自のデータフローを作成できます。付属のデータフローを参考にして、各ステージおよびデータフローを全体としてどのように完成させればよいか確認することもできます。

# 著作権に関する通知

© 2019 Pitney Bowes. All rights reserved. MapInfo および Group 1 Software は Pitney Bowes Software Inc. の商標です。その他のマークおよび商標はすべて、それぞれの所有者の資産です。

### USPS® 情報

Pitney Bowes Inc. は、ZIP + 4® データベースを光学および磁気媒体に発行および販売する非独占的ライセンスを所有しています。CASS、CASS 認定、DPV、eLOT、FASTforward、First-Class Mail、Intelligent Mail、LACS<sup>Link</sup>、NCOA<sup>Link</sup>、PAVE、PLANET Code、Postal Service、POSTNET、Post Office、RDI、Suite<sup>Link</sup>、United States Postal Service、Standard Mail、United States Post Office、USPS、ZIP Code、および ZIP + 4 の各商標は United States Postal Service が所有します。United States Postal Service に帰属する商標はこれに限りません。

Pitney Bowes Inc. は、NCOA<sup>Link</sup>® 処理に対する USPS® の非独占的ライセンスを所有しています。

Pitney Bowes Software の製品、オプション、およびサービスの価格は、USPS® または米国政府によって規定、制御、または承認されるものではありません。RDI™ データを利用して郵便送料を判定する場合に、使用する郵便配送業者の選定に関するビジネス上の意思決定が USPS® または米国政府によって行われることはありません。

### データ プロバイダおよび関連情報

このメディアに含まれて、Pitney Bowes Software アプリケーション内で使用されるデータ製品は、各種商標によって、および次の 1 つ以上の著作権によって保護されています。

© Copyright United States Postal Service. All rights reserved.

© 2014 TomTom. All rights reserved. TomTom および TomTom ロゴは TomTom N.V の登録商標です。

© 2016 HERE

Fuente: INEGI (Instituto Nacional de Estadística y Geografía)

電子データに基づいています。© National Land Survey Sweden.

© Copyright United States Census Bureau

© Copyright Nova Marketing Group, Inc.

このプログラムの一部は著作権で保護されています。© Copyright 1993-2007 by Nova Marketing Group Inc. All Rights Reserved

© Copyright Second Decimal, LLC

© Copyright Canada Post Corporation

この CD-ROM には、Canada Post Corporation が著作権を所有している編集物からのデータが収録されています。

© 2007 Claritas, Inc.

Geocode Address World データ セットには、  
<http://creativecommons.org/licenses/by/3.0/legalcode> に存在するクリエイティブ コモンズ アトリビューション ライセンス (「アトリビューション ライセンス」) の下に提供されている GeoNames Project ([www.geonames.org](http://www.geonames.org)) からライセンス供与されたデータが含まれています。お客様による GeoNames データ (Spectrum™ Technology Platform ユーザ マニュアルに記載) の使用は、アトリビューションライセンスの条件に従う必要があります。お客様と Pitney Bowes Software, Inc. との契約と、アトリビューション ライセンスの間に矛盾が生じる場合は、アトリビューションライセンスのみに基づいてそれを解決する必要があります。お客様による GeoNames データの使用に関しては、アトリビューション ライセンスが適用されるためです。



3001 Summer Street  
Stamford CT 06926-0700  
USA

[www.pitneybowes.com](http://www.pitneybowes.com)