

Spectrum™ Technology Platform

Version 12.0 SP1

Machine Learning-Handbuch



Inhalt

1 - Einführung

Machine Learning-Modul	5
Machine Learning-Workflow	6

2 - Binning

Einführung in das Binning	9
Konfigurieren von Standardoptionen	9
Binning-Ausgabe	10

3 - Binning Lookup

Einführung in Binning Lookup	12
Definieren von Binning-Eigenschaften	12
Binning-Ausgabe	12

4 - K-Means Clustering

Einführung in das „K-Means Clustering“	14
Definieren von Modelleigenschaften	14
Konfigurieren von Standardoptionen	15
Konfigurieren erweiterter Optionen	15
Modellausgabe	16

5 - Linear Regression

Einführung in Linear Regression	19
Definieren von Modelleigenschaften	19
Konfigurieren von Standardoptionen	20
Konfigurieren erweiterter Optionen	21
Modellausgabe	22

6 - Logistic Regression

Einführung in Logistic Regression	24
Definieren von Modelleigenschaften	24
Konfigurieren von Standardoptionen	25
Konfigurieren erweiterter Optionen	25
Modellausgabe	27

7 - Principal Component Analysis

Einführung in Principal Component Analysis	29
Definieren von Modelleigenschaften	29
Konfigurieren von Standardoptionen	30
Konfigurieren erweiterter Optionen	30
Modellausgabe	31

8 - Random Forest Classification

Einführung in Random Forest Classification	33
Definieren von Modelleigenschaften	33
Konfigurieren von Standardoptionen	34
Konfigurieren erweiterter Optionen	35
Modellausgabe	37

9 - Random Forest Regression

Einführung in Random Forest Regression	39
Definieren von Modelleigenschaften	39
Konfigurieren von Standardoptionen	40
Konfigurieren erweiterter Optionen	41
Modellausgabe	43

10 - Java Model Scoring

Einführung in Java Model Scoring	45
Definieren von Modelleigenschaften	45
Modellausgabe	46

11 - Machine

Learning-Modellverwaltung

Einführung in die Machine	
Learning-Modellverwaltung	48
Die Registerkarte „Modelldetail“	49

1 - Einführung

In this section

Machine Learning-Modul	5
Machine Learning-Workflow	6

Machine Learning-Modul

Das Spectrum™ Technology Platform Machine Learning-Modul bietet die Möglichkeit, ein Binning für numerische Daten durchzuführen, überwachte und unüberwachte Machine Learning-Modelle anzupassen und Daten in diesen Modellen zu bewerten.

Anmerkung: Das Machine Learning-Modul wird nur unter Windows- und Linux-Betriebssystemen unterstützt.

Binning

Beim Binning werden Datensätze für eine kontinuierliche Variable in Gruppen (Bins) aufgeteilt, ohne dass dabei Zielinformationen berücksichtigt werden. Sie können das unbeaufsichtigte Binning mit einer der beiden folgenden Methoden durchführen: mit Bins vom Typ „equal-width“ oder mit Bins vom Typ „equal-frequency“.

Binning Lookup

Binning Lookup wendet vorher definiertes Binning auf neue Daten an und verwendet dabei vorhandene Bins, die in Datenflüssen erstellt wurden, die den Binning-Schritt verwenden.

K-Means Clustering

Beim „K-Means Clustering“ werden Modelle auf der Grundlage des analytischen Clusterings erstellt. Dabei wird eine Reihe von Datensätzen basierend auf Datenwerten in Cluster mit ähnlichen Datensätzen segmentiert.

Linear Regression

Mit Linear Regression werden Modelle aus Datasets erstellt, die im Hinblick auf Eingabevariablen kontinuierliche Ziele verwenden.

Logistic Regression

Mit Logistic Regression werden Modelle aus Datasets erstellt, die im Hinblick auf Eingabevariablen binäre Ziele verwenden.

Principal Component Analysis

Principal Component Analysis (PCA) ist ein statistisches Verfahren, das einen Beobachtungssatz von möglicherweise korrelierten Variablen in einen Wertesatz von linear nicht korrelierten Variablen (prinzipielle Komponenten) umwandelt.

Random Forest Classification

Mit Random Forest Classification werden Modelle aus Datasets erstellt, die im Hinblick auf Eingabevariablen binäre oder multinomiale Ziele verwenden.

Random Forest Regression

Mit Random Forest Regression werden Modelle aus Datasets erstellt, die im Hinblick auf Eingabevariablen kontinuierliche Ziele verwenden.

Java Model Scoring

Dieses Feature bewertet mithilfe der Formel, die beim Anpassen eines Machine Learning-Modells erstellt wird, neue Daten.

Machine Learning-Modellverwaltung

Die „Machine Learning“-Modellverwaltung ermöglicht Ihnen die Verwaltung aller Machine Learning-Modelle auf Ihrem Spectrum™ Technology Platform-Server. Sie können Modelle verfügbar machen, die Verfügbarkeit von Modellen aufheben oder Modelle löschen. Zusätzlich können Sie zu jedem Modell detaillierte Informationen anzeigen und zwei beliebige Modelle des gleichen Typs miteinander vergleichen.

Anmerkung: Das Machine Learning-Modul verwendet eine zugrunde liegende H2O.ai-Bibliothek für Modellierungsalgorithmen und Java Model Scoring.

Machine Learning-Workflow

Ein typischer Machine Learning-Workflow umfasst folgende Schritte, die in mindestens einem Datenfluss ausgeführt werden:

1. Greifen Sie über andere Spectrum-Module, z. B. Data Integration, auf die Daten zu.
2. Bereiten Sie die Daten mit Schritten aus anderen Spectrum-Modulen vor, z. B. denen im Data Integration-, im Data Quality- und im Core-Modul.
3. Passen Sie ein Machine Learning-Modell an, führen Sie den Datenfluss aus und überprüfen Sie anschließend die Inhalte auf der Registerkarte „Modellausgabe“ im Modellschritt. Anschließend können Sie das Modell bei Bedarf anpassen und den Datenfluss erneut ausführen. Im Anschluss müssen Sie die vollständige Ausgabe der Modellbewertung im Tool für die „Machine Learning“-Modellverwaltung überprüfen. Sie können jeweils ein Modell überprüfen oder zwei Modelle miteinander vergleichen.
4. (Optional) Wenn das Modell für die Bewertung von Daten verwendet wird, müssen Sie das Modell im Tool für die „Machine Learning“-Modellverwaltung verfügbar machen. Mit diesem Tool wird das Modell für den „Java Model Scoring“-Schritt zur Verfügung gestellt.
 - a. Erstellen Sie mithilfe der oben beschriebenen Schritte 1 – 2 einen Spectrum™ Technology Platform-Datenfluss, und ersetzen Sie Schritt 3 dann durch den „Java Model Scoring“-Schritt.

- Richten Sie diesen Datenfluss so ein, dass er im Batchmodus ausgeführt wird, um eine Datei mit Modellbewertungen aufzufüllen, die auf aktualisierte Daten angewendet werden (die Felder, die als X-Felder oder Eingaben verwendet werden, werden in den Schritten 1–2 als natürlicher Bestandteil des Handeltreibens aktualisiert).
- b. Verwenden Sie alternativ zum Bewerten von bedarfsgesteuerten Daten einen Webservice in Spectrum™ Technology Platform. Beispiel: Greifen Sie auf die Website zu, rufen Sie die Kunden-ID und die Modelleingaben ab, bewerten Sie diese Eingaben, und geben Sie die Bewertung an einen Prozess zurück, der Webinhalte für Ihren Kunden anpasst.
5. (Optional) Sie können Modellbewertungen auch in einer „Data Hub“-Diagrammdatenbank als Entitätseigenschaft, auf Karten oder in CES-Anwendungen bereitstellen.

2 - Binning

In this section

Einführung in das Binning	9
Konfigurieren von Standardoptionen	9
Binning-Ausgabe	10

Einführung in das Binning

Der „Binning“-Schritt führt das so genannte unbeaufsichtigte Binning durch, bei dem eine kontinuierliche Variable in Gruppen (Bins) unterteilt wird. Dabei werden keine objektiven Informationen berücksichtigt. Die erfassten Daten beinhalten Bereiche, Mengen und Prozentsätze von Werten der einzelnen Bereiche.

Zu den Vorteilen bei der Durchführung des Binnings zählen folgende:

- Es lässt zu, dass Datensätze mit fehlenden Daten in das Modell eingeschlossen werden.
- Es steuert bzw. verringert die Auswirkungen von Ausreißern innerhalb des Modells.
- Es löst das Problem, dass die Merkmale verschiedene Skalierungen aufweisen. Dadurch können die Gewichtungen der Koeffizienten im Endmodell miteinander verglichen werden.

Sie können beim unbeaufsichtigten Binning der Spectrum™ Technology Platform Bins vom Typ „equal-width“ verwenden, bei denen die Daten in gleich große Bins unterteilt werden, oder in Bins vom Typ „equal-frequency“, bei denen die Daten in Gruppen aufgeteilt werden, die in etwa die gleiche Anzahl von Datensätzen enthalten. Im „Binning“-Schritt werden Bins vom Typ „equal-width“ als „Equal Range“-Bins und Bins vom Typ „equal-frequency“ als „Equal Population“-Bins bevorzugt.

Sie können über Anweisungen auf der Befehlszeile eine Liste von Binnings anzeigen und Binnings löschen. Weitere Informationen erhalten Sie unter „Machine Learning-Modul“ im Abschnitt [Administrationsumgebung](#) des Administratorhandbuchs.

Konfigurieren von Standardoptionen

1. Wählen Sie aus, ob Sie den **Binning-Stil** „equal-range“ oder „equal-population“ durchführen möchten.
2. Wählen Sie unter **Nullwert-Bin** aus, wie Sie mit leeren Bin-Feldern umgehen möchten, die unbekannte Werte aufgrund von fehlenden Daten darstellen. Wählen Sie **Höchste** aus, um dem höchsten Bin Nullwerte zuzuweisen, und wählen Sie **Niedrigste** aus, um dem niedrigsten Bin Nullwerte zuzuweisen. Der niedrigste Bin enthält immer 1.
3. Klicken Sie auf **Interne Ziel-Bins** und geben Sie die Anzahl der Bins ein, die Sie zwischen den End-Bins einfügen möchten. Wenn Sie das Binning vom Typ „equal-range“ durchführen, können Sie diesen Verarbeitungstyp oder **Bin-Breite** auswählen, jedoch nicht beides. Wenn Sie das Binning vom Typ „equal-population“ durchführen, können Sie nur interne Bins verarbeiten.
4. Wenn Sie das Binning vom Typ „equal-range“ durchführen und statt der Verarbeitung eines internen Bins diesen Verarbeitungstyp auswählen möchten, klicken Sie auf **Bin-Breite** und geben Sie die Anzahl der Einheiten an, die in den einzelnen Bins enthalten sein sollen.

5. Klicken Sie bei jedem Feld, dessen Daten beim Binning eingeschlossen werden sollen, auf **Einschließen**. Beachten Sie, dass nur numerische Felder in dieser Liste angezeigt werden.
6. Klicken Sie auf **OK**, um Ihre Einstellungen zu speichern.

Binning-Ausgabe

Der „Binning“-Schritt verfügt über zwei Ausgabeports. Am ersten Port werden für jedes ausgewählte Eingabefeld alle Eingabefelder und ein gebinntes Feld ausgegeben. Beispiel: Wenn die Eingabe die Felder „Name“, „Alter“ und „Einkommen“ enthält und Sie das Binning für „Alter“ und „Einkommen“ durchführen, enthält die Ausgabe des ersten Ports folgende Felder:

- Name
- Alter
- Binned_Age
- Einkommen
- Binned_Income

Am zweiten Port werden für jedes ausgewählte Eingabefeld vier Informationstypen ausgegeben. Beispiel: Wenn Sie das Binning für „Alter“ durchführen, enthält die Ausgabe des zweiten Ports folgende Felder:

- Age_Bins
- Age_BinValue
- Age_Count
- Age_Percentage

3 - Binning Lookup

In this section

Einführung in Binning Lookup	12
Definieren von Binning-Eigenschaften	12
Binning-Ausgabe	12

Einführung in Binning Lookup

Binning Lookup wendet vorher definiertes Binning auf neue Daten an und verwendet dabei vorhandene Bins, die in Datenflüssen erstellt wurden, die den **Binning**-Schritt verwenden.

Definieren von Binning-Eigenschaften

1. Klicken Sie unter **Primäre Schritte/Bereitgestellte Schritte/Machine Learning** auf den **Binning Lookup**-Schritt und ziehen Sie ihn auf die Arbeitsfläche. Platzieren Sie ihn an einer beliebigen Stelle im Datenfluss und verbinden Sie ihn mit anderen Schritten. Beachten Sie, dass die Eingabephase die Daten enthalten muss, die gebinnt werden sollen. Für die Binning-Ausgabe ist ein Ausgabeschritt erforderlich. Sie können optional einen zweiten Ausgabeschritt verbinden, um die Binning-Zusammenfassung zu erfassen.
2. Wählen Sie den entsprechenden **Binning-Namen** aus dem Dropdown-Menü aus. Das sind die Namen vorhandener Bins, die von den Binning-Schritt verwendenden Datenflüssen erstellt wurden.
3. Die Felder **Binning-Typ** und **Beschreibung** werden mit den Bins des Binning-Namens importiert, den Sie in Schritt 2 ausgewählt haben, und sind daher nicht bearbeitbar.
4. Das Raster **Eingaben** zeigt alle Felder, die für das Binning im Binning-Schritt eingeschlossen wurden, zusammen mit dem Datentyp.
5. Klicken Sie auf **OK**, um Ihre Einstellungen zu speichern.

Binning-Ausgabe

Auf dieser Registerkarte werden die Felder und Datentypen angezeigt, die von Ihrem „Binning Lookup“-Schritt gebinnt werden. Weitere Informationen zur Ausgabe, die durch die Verwendung eines Binning-Schritts generiert wird, finden Sie unter **Binning-Ausgabe** auf Seite 10. Wir haben nun die Option verfügbar gemacht, gebinnte Felder aus der Registerkarte der Binning-Ausgabe zu bearbeiten. Sie können „Spectrum Binning-Felder“ verwenden, um ein Binning-Feld umzubenennen. Sie haben nun auch die Option, Binning-Felder einzuschließen/auszuschließen.

4 - K-Means Clustering

In this section

Einführung in das „K-Means Clustering“	14
Definieren von Modelleigenschaften	14
Konfigurieren von Standardoptionen	15
Konfigurieren erweiterter Optionen	15
Modellausgabe	16

Einführung in das „K-Means Clustering“

Beim „K-Means Clustering“ werden Modelle auf der Grundlage des analytischen Clusterings erstellt. Dabei wird eine Reihe von Datensätzen basierend auf Datenwerten in Cluster mit ähnlichen Datensätzen segmentiert.

Sie müssen zunächst die Registerkarte „Modelleigenschaften“ ausfüllen, um Ihr Modell erstellen zu können. Auf den Registerkarten „Standardoptionen“ und „Erweiterte Optionen“ werden genügend Standardeinstellungen zum Abschließen eines Auftrags bereitgestellt. Sie können diese Einstellungen jedoch auch ändern und an Ihre Bedürfnisse anpassen. Anschließend führen Sie Ihren Auftrag aus, und auf der Registerkarte „Modellausgabe“ wird eine eingeschränkte Version der Details zur resultierenden Modellausgabe angezeigt. Das Modell wird auf dem Spectrum™ Technology Platform-Server gespeichert und die vollständige Ausgabe ist im Tool für die „Machine Learning“-Modellverwaltung verfügbar.

Definieren von Modelleigenschaften

1. Klicken Sie unter **Primäre Schritte/Bereitgestellte Schritte/Machine Learning** auf den **K-Means Clustering**-Schritt und ziehen Sie ihn auf die Arbeitsfläche. Platzieren Sie ihn an einer beliebigen Stelle im Datenfluss und verbinden Sie ihn mit anderen Schritten. Beachten Sie, dass der Eingabeschritt die Datenquelle sein muss, die Eingabevariablenfelder für Ihr Modell enthält. Ein Ausgabeschritt ist nicht erforderlich, es sei denn, Sie wählen auf der Registerkarte „Standardoptionen“ die Eingabedatenoption „Bewerten“ aus. Sie können auch eine Verbindung zu einem Ausgabeschritt herstellen, wenn Sie Ihre Ausgabe unabhängig von dem Tool für die Machine Learning-Modellverwaltung erfassen möchten.
2. Doppelklicken Sie auf den „K-Means“-Schritt, um das Dialogfeld „**K-Means Clustering**“-Optionen anzuzeigen.
3. Geben Sie einen **Modellnamen** ein, wenn Sie nicht den Standardnamen verwenden möchten.
4. Optional: Aktivieren Sie das Kästchen **Überschreiben**, um das vorhandene Modell mit neuen Daten zu überschreiben.
5. Geben Sie die **Anzahl der Cluster** ein, die in Ihrem Modell enthalten sein sollen, wenn diese von der Standardanzahl (5) abweicht.
6. Optional: Geben Sie eine **Beschreibung** des Modells ein.
7. Klicken Sie bei jedem Feld, dessen Daten zum Modell hinzugefügt werden sollen, auf **Einschließen**.
8. Geben Sie über das Dropdown-Menü **Modelltyp** an, ob das Eingabefeld als numerisches Feld, als kategorisches Feld oder als DateTime-Feld verwendet werden soll.

9. Klicken Sie auf **OK**, um das Modell und die Konfiguration zu speichern, oder fahren Sie mit der nächsten Registerkarte fort.

Konfigurieren von Standardoptionen

1. Lassen Sie **Eingabefelder standardisieren** aktiviert, um die numerischen Spalten zu standardisieren, damit diese keine Mittelwert- und Einheitenvarianz aufweisen.
Wenn Sie keine Standardisierung verwenden, können die Ergebnisse Komponenten enthalten, die von Variablen dominiert werden, die statt richtiger Beiträge relativ zu anderen Attributen in der Skalierung eine größere Varianz zu haben scheinen.
2. Aktivieren Sie **Eingabedaten bewerten**, um eine Spalte für die Modellvorhersage (Punktzahl) für Eingabedaten hinzuzufügen.
3. Aktivieren Sie **Anzahl der Cluster schätzen**, damit der Algorithmus „K-Means“ versucht, die Anzahl der in Ihrem Modell enthaltenen Cluster zu bestimmen. Auch wenn Sie die Anzahl der gewünschten Cluster auf der Registerkarte „Modelleigenschaften“ angeben, kann bei der Routine während der Verarbeitung festgestellt werden, dass für die Daten eine andere Anzahl von Clustern geeigneter wäre.
4. Geben Sie für den **Prozentsatz für Trainingsdaten** einen Wert zwischen 1 und 100 an, wenn die Eingabedaten zufällig in Stichproben für Trainings- und Testdaten aufgeteilt werden.
5. Geben Sie den Wert 100 abzüglich der in Schritt 5 eingegebenen Menge als **Prozentsatz für Testdaten** ein.
6. Geben Sie eine Ziffer als **Ausgangswert für Stichprobe** ein, um sicherzustellen, dass die Darstellung der Daten bei jeder Datenflussausführung gleich ist, wenn diese in Test- und Trainingsdaten aufgeteilt werden. Deaktivieren Sie dieses Feld, damit die Aufteilung bei jeder Datenflussausführung beliebig erfolgt.
7. Klicken Sie auf **OK**, um das Modell und die Konfiguration zu speichern, oder fahren Sie mit der nächsten Registerkarte fort.

Konfigurieren erweiterter Optionen

1. Lassen Sie **Konstante Felder ignorieren** aktiviert, damit Felder übersprungen werden, die für die einzelnen Datensätze die gleichen Werte enthalten.
2. Lassen Sie **Seed für Algorithmus** aktiviert, und geben Sie einen numerischen Ausgangswert ein, um sicherzustellen, dass die Darstellung der Daten bei jeder Datenflussausführung gleich ist, wenn diese in Test- und Trainingsdaten aufgeteilt werden. Deaktivieren Sie dieses Feld, damit die Aufteilung bei jeder Datenflussausführung beliebig erfolgt.

3. Wählen Sie im Dropdown-Menü **Init** den richtigen Initialisierungsmodus aus.

Furthest	Initialisiert den ersten Zentroid zufällig, den zweiten Zentroid initialisiert der Modus jedoch anschließend so, dass es der davon am weitesten entfernte Datenpunkt ist. Initialisiert die Zentroide so, dass sie gut verteilt sind.
Plus-Plus	Initialisiert das Clusterzentrum, bevor mit den standardmäßigen „k-means“-Optimierungsiterationen fortgefahren wird. Bei der „k-means++“-Initialisierung wird garantiert, dass der Algorithmus die Lösung „ $O(\log k)$ competitive“ für die optimale „k-means“-Lösung findet.
Random	Standardeinstellung. Wählt Cluster K zufällig aus der Gruppe der Beobachtungen N aus, damit die einzelnen Beobachtungen gleichermaßen die Möglichkeit haben, ausgewählt zu werden.

4. Lassen Sie **Seed für N-fach** aktiviert, und geben Sie einen numerischen Ausgangswert ein, um sicherzustellen, dass die Darstellung der Daten bei jeder Datenflussausführung gleich ist, wenn diese in Test- und Trainingsdaten aufgeteilt werden. Deaktivieren Sie dieses Feld, damit die Aufteilung bei jeder Datenflussausführung beliebig erfolgt.
5. Aktivieren Sie **N-fach** und geben Sie die Anzahl der Folds ein, wenn Sie eine Kreuzvalidierung durchführen.
6. Aktivieren Sie **Foldzuweisung**, und wählen Sie aus der Dropdown-Liste aus, ob Sie eine Kreuzvalidierung durchführen. Dieses Feld ist nur anwendbar, wenn Sie unter **N-fach** einen Wert eingegeben haben.

Auto	Standardeinstellung. Lässt zu, dass der Algorithmus automatisch eine Option auswählt; derzeit wird „Random“ verwendet.
Modulo	Teilt das Dataset gleichmäßig auf die Folds auf und ist nicht vom Ausgangswert abhängig.
Random	Teilt die Daten zufällig in „N-fach“-Bestandteile ein; diese Einstellung ist für umfangreiche Datasets am besten geeignet.

7. Aktivieren Sie **Maximale Iterationen** und geben Sie die Anzahl der Trainingsiterationen ein, die erfolgen sollen.
8. Klicken Sie auf **OK**, um das Modell und die Konfiguration zu speichern, oder fahren Sie mit der nächsten Registerkarte fort.

Modellausgabe

Auf dieser Registerkarte werden die Metriken angezeigt, mit denen Sie das angepasste Modell bewerten. Diese Felder können nicht bearbeitet werden. Die Spalte „Training“ enthält immer Daten. Wenn Sie auf der Registerkarte „Standardoptionen“ eine Aufteilung in Training/Test ausgewählt

haben, wird die Spalte „Test“ ebenfalls aufgefüllt. Eine Ausnahme besteht, wenn Sie auf der Registerkarte „Erweiterte Optionen“ eine Validierung vom Typ „N-fach“ ausgewählt haben. In diesem Fall wird die Spalte „N-fach“ aufgefüllt. Klicken Sie auf die Schaltfläche **Ausgabe**, um die Ausgabe erneut zu generieren, und klicken Sie auf **Modelldetails**, um in dem Tool für die „Machine Learning“-Modellverwaltung die gesamte Ausgabe anzuzeigen.

5 - Linear Regression

In this section

Einführung in Linear Regression	19
Definieren von Modelleigenschaften	19
Konfigurieren von Standardoptionen	20
Konfigurieren erweiterter Optionen	21
Modellausgabe	22

Einführung in Linear Regression

Mit Linear Regression können Sie Machine Learning durchführen, indem Sie Modelle aus Datasets erstellen, die im Hinblick auf Eingabevariablen kontinuierliche Ziele verwenden.

Sie müssen zunächst die Registerkarte „Modelleigenschaften“ ausfüllen, um Ihr Modell erstellen zu können. Auf den Registerkarten „Standardoptionen“ und „Erweiterte Optionen“ werden genügend Standardeinstellungen zum Abschließen eines Auftrags bereitgestellt. Sie können diese Einstellungen jedoch auch ändern und an Ihre Bedürfnisse anpassen. Anschließend führen Sie Ihren Auftrag aus, und auf der Registerkarte „Modellausgabe“ wird eine eingeschränkte Version des resultierenden Modells angezeigt. Die vollständige Ausgabe ist im Tool für die „Machine Learning“-Modellverwaltung verfügbar.

Definieren von Modelleigenschaften

1. Klicken Sie unter **Primäre Schritte/Bereitgestellte Schritte/Machine Learning** auf den **Linear Regression**-Schritt und ziehen Sie ihn auf die Arbeitsfläche. Platzieren Sie ihn an einer beliebigen Stelle im Datenfluss und verbinden Sie ihn mit anderen Schritten. Beachten Sie, dass der Eingabeschritt die Datenquelle sein muss, die Ziel- und Eingabevariablenfelder für Ihr Modell enthält. Ein Ausgabeschritt ist nicht erforderlich, es sei denn, Sie wählen auf der Registerkarte „Standardoptionen“ die Eingabedatenoption „Bewerten“ aus. Sie können auch eine Verbindung zu einem Ausgabeschritt herstellen, wenn Sie Ihre Ausgabe unabhängig von dem Tool für die Machine Learning-Modellverwaltung erfassen möchten.
2. Doppelklicken Sie auf den „Linear Regression“-Schritt, um das Dialogfeld „**Linear Regression**“-**Optionen** anzuzeigen.
3. Geben Sie einen **Modellnamen** ein, wenn Sie nicht den Standardnamen verwenden möchten.
4. Aktivieren Sie das Kontrollkästchen **Überschreiben**, um das vorhandene Modell mit neuen Daten zu überschreiben.
5. Klicken Sie auf das Dropdown-Menü **Zielfeld**, und wählen Sie ein numerisches Feld aus.
6. Geben Sie eine **Beschreibung** des Modells ein.
7. Klicken Sie bei jedem Feld, dessen Daten zum Modell hinzugefügt werden sollen, auf **Einschließen**. Achten Sie darauf, das Feld einzuschließen, das Sie als Zielfeld ausgewählt haben.
8. Geben Sie über das Dropdown-Menü **Modelltyp** an, ob die einzelnen Eingabefelder als numerisches Feld, als kategorisches Feld oder als DateTime-Feld verwendet werden sollen.
9. Klicken Sie auf **OK**, um das Modell und die Konfiguration zu speichern, oder fahren Sie mit der nächsten Registerkarte fort.

Konfigurieren von Standardoptionen

1. Lassen Sie **Eingabefelder standardisieren** aktiviert, um die numerischen Spalten zu standardisieren, damit diese keine Mittelwert- und Einheitenvarianz aufweisen.
Wenn Sie keine Standardisierung verwenden, können die Ergebnisse Komponenten enthalten, die von Variablen dominiert werden, die statt richtiger Beiträge relativ zu anderen Attributen in der Skalierung eine größere Varianz zu haben scheinen.
2. Aktivieren Sie **Eingabedaten bewerten**, um eine Spalte für die Modellvorhersage (Punktzahl) für Eingabedaten hinzuzufügen.
3. Wählen Sie eine **Verknüpfungsfunktion** aus der Dropdown-Liste aus. Dies spezifiziert die Verknüpfung zwischen zufälligen und systematischen Komponenten. Es besagt, in welcher Beziehung der erwartete Wert der Antwort mit dem linearen Prädiktor erklärender Variablen steht.

Identität Sagt sinnlose „Wahrscheinlichkeiten“ von unter null oder größer als eins voraus. Wird manchmal für binominale Daten verwendet, um ein lineares Wahrscheinlichkeitsmodell zu erzielen.

$$g(p) = p$$

Invers Berechnet die Umkehrung von Verknüpfungsfunktionen für reelle Schätzungen.

$$g(\mu_i) = 1/\mu_i$$

Protokoll Zählt Vorkommen in einer festen Zeitspanne und einem festen Bereich.

$$g(\mu_i) = \log(\mu_i)$$

4. Geben Sie an, wie mit fehlenden Daten umgegangen werden soll, indem Sie **Überspringen** aktivieren oder **Mittelwerte zuschreiben**, wodurch der Mittelwert für fehlende Daten hinzugefügt wird.
5. Geben Sie für den **Prozentsatz für Trainingsdaten** einen Wert zwischen 1 und 100 an, wenn die Eingabedaten zufällig in Stichproben für Trainings- und Testdaten aufgeteilt werden.
6. Geben Sie den Wert 100 abzüglich der in Schritt 5 eingegebenen Menge als **Prozentsatz für Testdaten** ein.
7. Geben Sie eine Ziffer als **Ausgangswert für Stichprobe** ein, um sicherzustellen, dass die Darstellung der Daten bei jeder Datenflussausführung gleich ist, wenn diese in Test- und Trainingsdaten aufgeteilt werden. Deaktivieren Sie dieses Feld, damit die Aufteilung bei jeder Datenflussausführung beliebig erfolgt.
8. Klicken Sie auf **OK**, um das Modell und die Konfiguration zu speichern, oder fahren Sie mit der nächsten Registerkarte fort.

Konfigurieren erweiterter Optionen

1. Lassen Sie **Konstante Felder ignorieren** aktiviert, damit Felder übersprungen werden, die für die einzelnen Datensätze die gleichen Werte enthalten.
2. Aktivieren Sie **p-Werte berechnen**, um p-Werte für die Parameterschätzungen zu berechnen.
3. Aktivieren Sie **Kollineare Spalte entfernen**, damit kollineare Spalten während der Modellerstellung automatisch entfernt werden. Dies führt zu einem Koeffizienten von 0 im zurückgegebenen Modell.

Diese Option muss aktiviert werden, wenn **p-Werte berechnen** ebenfalls aktiviert ist.

4. Lassen Sie **Konstanten Begriff einschließen (abfangen)** aktiviert, um einen konstanten Begriff im Modell einzuschließen (abzufangen).

Dieses Feld muss aktiviert werden, wenn **Kollineare Spalte entfernen** ebenfalls aktiviert ist.

5. Wählen Sie einen **Solver** aus der Dropdown-Liste aus. Beachten Sie, dass „CoordinateDescent“ und „CoordinateDescentNaive“ derzeit experimentell sind.

Auto	Solver wird basierend auf Eingabedaten und Parametern bestimmt.
CoordinateDescent	IRLSM mit der Version der Kovarianzaktualisierungen der zyklischen Koordinate, die aus der innersten Schleife stammt.
CoordinateDescentNaive	IRLSM mit der Version der naiven Aktualisierungen der zyklischen Koordinate, die aus der innersten Schleife stammt.
IRLSM	Ideal für Probleme mit einer geringen Anzahl von Prädiktoren oder Lambda-Suchvorgänge mit L1-Penalty.
LBFGS	Ideal für Datasets mit vielen Spalten.

6. Lassen Sie **Seed für N-fach** aktiviert, und geben Sie einen numerischen Ausgangswert ein, um sicherzustellen, dass die Darstellung der Daten bei jeder Datenflussausführung gleich ist, wenn diese in Test- und Trainingsdaten aufgeteilt werden. Deaktivieren Sie dieses Feld, damit die Aufteilung bei jeder Datenflussausführung beliebig erfolgt.
7. Aktivieren Sie **N-fach** und geben Sie die Anzahl der Folds ein, wenn Sie eine Kreuzvalidierung durchführen.
8. Klicken Sie auf **Foldzuweisung**, und wählen Sie aus der Dropdown-Liste aus, ob Sie eine Kreuzvalidierung durchführen. Dieses Feld ist nur anwendbar, wenn Sie unter **N-fach** einen Wert eingegeben haben und **Foldfeld** nicht angegeben ist.

Auto	Lässt zu, dass der Algorithmus automatisch eine Option auswählt; derzeit wird „Random“ verwendet.
Modulo	Teilt das Dataset gleichmäßig auf die Folds auf und ist nicht vom Ausgangswert abhängig.

Random Teilt die Daten zufällig in „N-fach“-Bestandteile ein; diese Einstellung ist für umfangreiche Datasets am besten geeignet.

9. Wenn Sie eine Kreuzvalidierung durchführen, aktivieren Sie **Foldfeld** und wählen Sie aus der Dropdown-Liste das Feld aus, das die Foldindexzuweisung für die Kreuzvalidierung enthält. Dieses Feld ist nur anwendbar, wenn Sie unter **N-fach** und **Foldzuweisung** keinen Wert eingegeben haben.
10. Aktivieren Sie **Maximale Iterationen** und geben Sie die Anzahl der Trainingsiterationen ein, die erfolgen sollen.
11. Aktivieren Sie **Ziel-Epsilon** und geben Sie den Schwellenwert für die Konvergenz an. Dieser Wert muss zwischen 0 und 1 liegen. Wenn der Zielwert geringer ist als dieser Schwellenwert, wird das Modell konvergiert.
12. Aktivieren Sie **Beta-Epsilon** und geben Sie den Schwellenwert für die Konvergenz an. Dieser Wert muss zwischen 0 und 1 liegen. Wenn der Zielwert geringer ist als dieser Schwellenwert, wird das Modell konvergiert. Wenn die L1-Normalisierung der aktuellen Beta-Änderung unter diesem Schwellenwert liegt, sollten Sie die Verwendung der Konvergenz in Erwägung ziehen.
13. Klicken Sie auf **OK**, um das Modell und die Konfiguration zu speichern, oder fahren Sie mit der nächsten Registerkarte fort.

Modellausgabe

Auf dieser Registerkarte werden die Metriken angezeigt, mit denen Sie das angepasste Modell bewerten. Diese Felder können nicht bearbeitet werden. Die Spalte „Training“ enthält immer Daten. Wenn Sie auf der Registerkarte „Standardoptionen“ eine Aufteilung in Training/Test ausgewählt haben, wird die Spalte „Test“ ebenfalls aufgefüllt. Eine Ausnahme besteht, wenn Sie auf der Registerkarte „Erweiterte Optionen“ eine Validierung vom Typ „N-fach“ ausgewählt haben. In diesem Fall wird die Spalte „N-fach“ aufgefüllt.

Nachdem Sie Ihren Auftrag ausgeführt haben, wird das resultierende Modell auf dem Spectrum™ Technology Platform-Server gespeichert. Klicken Sie auf die Schaltfläche **Ausgabe**, um die Ausgabe erneut zu generieren, und klicken Sie auf **Modelldetails**, um in dem Tool für die „Machine Learning“-Modellverwaltung die gesamte Ausgabe anzuzeigen.

6 - Logistic Regression

In this section

Einführung in Logistic Regression	24
Definieren von Modelleigenschaften	24
Konfigurieren von Standardoptionen	25
Konfigurieren erweiterter Optionen	25
Modellausgabe	27

Einführung in Logistic Regression

Mit Logistic Regression können Sie Machine Learning durchführen, indem Sie Modelle aus Datasets erstellen, die im Hinblick auf Eingabevariablen binäre Ziele verwenden.

Sie müssen zunächst die Registerkarte „Modelleigenschaften“ ausfüllen, um Ihr Modell erstellen zu können. Auf den Registerkarten „Standardoptionen“ und „Erweiterte Optionen“ werden genügend Standardeinstellungen zum Abschließen eines Auftrags bereitgestellt. Sie können diese Einstellungen jedoch auch ändern und an Ihre Bedürfnisse anpassen. Anschließend führen Sie Ihren Auftrag aus, und auf der Registerkarte „Modellausgabe“ wird eine eingeschränkte Version des resultierenden Modells angezeigt. Die vollständige Ausgabe ist im Tool für die „Machine Learning“-Modellverwaltung verfügbar.

Definieren von Modelleigenschaften

1. Klicken Sie unter **Primäre Schritte/Bereitgestellte Schritte/Machine Learning** auf den **Logistic Regression**-Schritt und ziehen Sie ihn auf die Arbeitsfläche. Platzieren Sie ihn an einer beliebigen Stelle im Datenfluss und verbinden Sie ihn mit anderen Schritten. Beachten Sie, dass der Eingabeschritt die Datenquelle sein muss, die Ziel- und Eingabevariablenfelder für Ihr Modell enthält. Ein Ausgabeschritt ist nicht erforderlich, es sei denn, Sie wählen auf der Registerkarte „Standardoptionen“ die Eingabedatenoption „Bewerten“ aus. Sie können auch eine Verbindung zu einem Ausgabeschritt herstellen, wenn Sie Ihre Ausgabe unabhängig von dem Tool für die Machine Learning-Modellverwaltung erfassen möchten.
2. Doppelklicken Sie auf den „Logistic Regression“-Schritt, um das Dialogfeld „**Logistic Regression**“-**Optionen** anzuzeigen.
3. Geben Sie einen **Modellnamen** ein, wenn Sie nicht den Standardnamen verwenden möchten.
4. Aktivieren Sie das Kontrollkästchen **Überschreiben**, um das vorhandene Modell mit neuen Daten zu überschreiben.
5. Klicken Sie auf das Dropdown-Menü **Zielfeld** und wählen Sie „Kategorisch“ aus.
6. Geben Sie eine **Beschreibung** des Modells ein.
7. Klicken Sie bei jedem Feld, dessen Daten zum Modell hinzugefügt werden sollen, auf **Einschließen**. Achten Sie darauf, das Feld einzuschließen, das Sie als Zielfeld ausgewählt haben.
8. Geben Sie über das Dropdown-Menü **Modelltyp** an, ob die einzelnen Eingabefelder als numerisches Feld, als kategorisches Feld oder als DateTime-Feld verwendet werden sollen.
9. Klicken Sie auf **OK**, um das Modell und die Konfiguration zu speichern, oder fahren Sie mit der nächsten Registerkarte fort.

Konfigurieren von Standardoptionen

1. Lassen Sie **Eingabefelder standardisieren** aktiviert, um die numerischen Spalten zu standardisieren, damit diese keine Mittelwert- und Einheitenvarianz aufweisen.
Wenn Sie keine Standardisierung verwenden, können die Ergebnisse Komponenten enthalten, die von Variablen dominiert werden, die statt richtiger Beiträge relativ zu anderen Attributen in der Skalierung eine größere Varianz zu haben scheinen.
2. Aktivieren Sie **Eingabedaten bewerten**, um eine Spalte für die Modellvorhersage (Punktzahl) für Eingabedaten hinzuzufügen.
3. Aktivieren Sie **Vorherig**, wenn die Daten erfasst wurden und die Bedeutung der Antwort nicht die Realität widerspiegelt. Geben Sie anschließend die vorherige Wahrscheinlichkeit für $p(y=1)$ in das Textfeld ein.
4. Geben Sie an, wie mit fehlenden Daten umgegangen werden soll, indem Sie **Überspringen** aktivieren oder **Mittelwerte zuschreiben**, wodurch der Mittelwert für fehlende Daten hinzugefügt wird.
5. Geben Sie für den **Prozentsatz für Trainingsdaten** einen Wert zwischen 1 und 100 an, wenn die Eingabedaten zufällig in Stichproben für Trainings- und Testdaten aufgeteilt werden.
6. Geben Sie den Wert 100 abzüglich der in Schritt 5 eingegebenen Menge als **Prozentsatz für Testdaten** ein.
7. Geben Sie eine Ziffer als **Ausgangswert für Stichprobe** ein, um sicherzustellen, dass die Darstellung der Daten bei jeder Datenflussausführung gleich ist, wenn diese in Test- und Trainingsdaten aufgeteilt werden. Deaktivieren Sie dieses Feld, damit die Aufteilung bei jeder Datenflussausführung beliebig erfolgt.
8. Klicken Sie auf **OK**, um das Modell und die Konfiguration zu speichern, oder fahren Sie mit der nächsten Registerkarte fort.

Konfigurieren erweiterter Optionen

1. Lassen Sie **Konstante Felder ignorieren** aktiviert, damit Felder übersprungen werden, die für die einzelnen Datensätze die gleichen Werte enthalten.
2. Lassen Sie **p-Werte berechnen** aktiviert, um p-Werte für die Parameterschätzungen zu berechnen.
3. Lassen Sie **Kollineare Spalte entfernen** aktiviert, damit kollineare Spalten während der Modellerstellung automatisch entfernt werden. Dies führt zu einem Koeffizienten von 0 im zurückgegebenen Modell.
Diese Option muss aktiviert werden, wenn **p-Werte berechnen** ebenfalls aktiviert ist.

4. Lassen Sie **Konstanten Begriff einschließen (abfangen)** aktiviert, um einen konstanten Begriff im Modell einzuschließen (abzufangen).

Dieses Feld muss aktiviert werden, wenn **Kollineare Spalte entfernen** ebenfalls aktiviert ist.

5. Wählen Sie einen **Solver** aus der Dropdown-Liste aus. Beachten Sie, dass „CoordinateDescentNaive“ und „CoordinateDescentNaive“ derzeit experimentell sind.

Auto Solver wird basierend auf Eingabedaten und Parametern bestimmt.

CoordinateDescentNaive IRLSM mit der Version der Kovarianzaktualisierungen der zyklischen Koordinate, die aus der innersten Schleife stammt.

CoordinateDescentNaive IRLSM mit der Version der naiven Aktualisierungen der zyklischen Koordinate, die aus der innersten Schleife stammt.

IRLSM Ideal für Probleme mit einer geringen Anzahl von Prädiktoren oder Lambda-Suchvorgänge mit L1-Penalty.

L_BFGS Ideal für Datasets mit vielen Spalten.

6. Lassen Sie **Seed für N-fach** aktiviert, und geben Sie einen numerischen Ausgangswert ein, um sicherzustellen, dass die Darstellung der Daten bei jeder Datenflussausführung gleich ist, wenn diese in Test- und Trainingsdaten aufgeteilt werden. Deaktivieren Sie dieses Feld, damit die Aufteilung bei jeder Datenflussausführung beliebig erfolgt.

7. Aktivieren Sie **N-fach** und geben Sie die Anzahl der Folds ein, wenn Sie eine Kreuzvalidierung durchführen.

8. Aktivieren Sie **Foldzuweisung**, und wählen Sie aus der Dropdown-Liste aus, ob Sie eine Kreuzvalidierung durchführen. Dieses Feld ist nur anwendbar, wenn Sie unter **N-fach** einen Wert eingegeben haben und **Foldfeld** nicht angegeben ist.

Auto Lässt zu, dass der Algorithmus automatisch eine Option auswählt; derzeit wird „Random“ verwendet.

Modulo Teilt das Dataset gleichmäßig auf die Folds auf und ist nicht vom Ausgangswert abhängig.

Random Teilt die Daten zufällig in „N-fach“-Bestandteile ein; diese Einstellung ist für umfangreiche Datasets am besten geeignet.

Stratified Schichtet die Folds basierend auf der Antwortvariable für Klassifizierungsprobleme. Verteilt Beobachtungen aus den verschiedenen Klassen gleichmäßig auf alle Datasets, wenn ein Dataset in Trainings- und Testdaten aufgeteilt wird. Dies kann nützlich sein, wenn viele Klassen vorhanden sind und das Dataset relativ klein ist.

9. Wenn Sie eine Kreuzvalidierung durchführen, aktivieren Sie **Foldfeld** und wählen Sie aus der Dropdown-Liste das Feld aus, das die Foldindexzuweisung für die Kreuzvalidierung enthält.

Dieses Feld ist nur anwendbar, wenn Sie unter **N-fach** und **Foldzuweisung** keinen Wert eingegeben haben.

10. Aktivieren Sie **Maximale Iterationen** und geben Sie die Anzahl der Trainingsiterationen ein, die erfolgen sollen.
11. Aktivieren Sie **Ziel-Epsilon** und geben Sie den Schwellenwert für die Konvergenz an. Dieser Wert muss zwischen 0 und 1 liegen. Wenn der Zielwert geringer ist als dieser Schwellenwert, wird das Modell konvergiert.
12. Aktivieren Sie **Beta-Epsilon** und geben Sie den Schwellenwert für die Konvergenz an. Dieser Wert muss zwischen 0 und 1 liegen. Wenn der Zielwert geringer ist als dieser Schwellenwert, wird das Modell konvergiert. Wenn die L1-Normalisierung der aktuellen Beta-Änderung unter diesem Schwellenwert liegt, sollten Sie die Verwendung der Konvergenz in Erwägung ziehen.
13. Klicken Sie auf **OK**, um das Modell und die Konfiguration zu speichern, oder fahren Sie mit der nächsten Registerkarte fort.

Modellausgabe

Auf dieser Registerkarte werden die Metriken angezeigt, mit denen Sie das angepasste Modell bewerten. Diese Felder können nicht bearbeitet werden. Die Spalte „Training“ enthält immer Daten. Wenn Sie auf der Registerkarte „Standardoptionen“ eine Aufteilung in Training/Test ausgewählt haben, wird die Spalte „Test“ ebenfalls aufgefüllt. Eine Ausnahme besteht, wenn Sie auf der Registerkarte „Erweiterte Optionen“ eine Validierung vom Typ „N-fach“ ausgewählt haben. In diesem Fall wird die Spalte „N-fach“ aufgefüllt.

Nachdem Sie Ihren Auftrag ausgeführt haben, wird das resultierende Modell auf dem Spectrum™ Technology Platform-Server gespeichert. Klicken Sie auf die Schaltfläche **Ausgabe**, um die Ausgabe erneut zu generieren, und klicken Sie auf **Modelldetails**, um in dem Tool für die „Machine Learning“-Modellverwaltung die gesamte Ausgabe anzuzeigen.

7 - Principal Component Analysis

In this section

Einführung in Principal Component Analysis	29
Definieren von Modelleigenschaften	29
Konfigurieren von Standardoptionen	30
Konfigurieren erweiterter Optionen	30
Modellausgabe	31

Einführung in Principal Component Analysis

Principal Component Analysis (PCA) ist ein statistisches Verfahren, das einen Beobachtungssatz von möglicherweise korrelierten Variablen in einen Wertesatz von linear nicht korrelierten Variablen (prinzipielle Komponenten) umwandelt.

Sie müssen zunächst die Registerkarte „Modelleigenschaften“ ausfüllen, um Ihr Modell erstellen zu können. Auf den Registerkarten „Standardoptionen“ und „Erweiterte Optionen“ werden genügend Standardeinstellungen zum Abschließen eines Auftrags bereitgestellt. Sie können diese Einstellungen jedoch auch ändern und an Ihre Bedürfnisse anpassen. Anschließend führen Sie Ihren Auftrag aus, und auf der Registerkarte „Modellausgabe“ wird eine eingeschränkte Version des resultierenden Modells angezeigt. Die vollständige Ausgabe ist im Tool für die „Machine Learning“-Modellverwaltung verfügbar. Wenn Sie mit der Ausgabe Ihres Modells zufrieden sind, können Sie es verfügbar machen und in einem Bewertungsdatenfluss verwenden.

Definieren von Modelleigenschaften

1. Klicken Sie unter **Primäre Schritte/Bereitgestellte Schritte/Machine Learning** auf den **PCA-Optionen**-Schritt und ziehen Sie ihn auf die Arbeitsfläche. Platzieren Sie ihn an einer beliebigen Stelle im Datenfluss und verbinden Sie ihn mit anderen Schritten. Beachten Sie, dass der Eingabeschritt die Datenquelle sein muss, die die prinzipiellen Komponenten Ihres Modells enthält. Ein Ausgabeschritt ist nicht erforderlich, jedoch können Sie eine Verbindung zu einem Ausgabeschritt herstellen, wenn Sie Ihre Ausgabe unabhängig von dem Tool für die Machine Learning-Modellverwaltung erfassen möchten.
2. Doppelklicken Sie auf den „PCA Options“-Schritt, um das Dialogfeld **PCA-Optionen** anzuzeigen.
3. Geben Sie einen **Modellnamen** ein, wenn Sie nicht den Standardnamen verwenden möchten.
4. Optional: Aktivieren Sie das Kästchen **Überschreiben**, um das vorhandene Modell mit neuen Daten zu überschreiben.
5. Geben Sie die Anzahl von **prinzipiellen Komponenten** ein, die Ihr Modell enthalten soll.
6. Optional: Geben Sie eine **Beschreibung** des Modells ein.
7. Klicken Sie in der Tabelle **Eingaben** bei jedem Feld, dessen Daten zum Modell hinzugefügt werden sollen, auf „Einschließen“.
8. Geben Sie über das Dropdown-Menü **Modelldatentyp** an, ob das Eingabefeld als kategorisches, DateTime-, numerisches, Zeichenfolgen- oder uniqueid-Feld verwendet werden soll.
9. Klicken Sie auf **OK**, um das Modell und die Konfiguration zu speichern, oder fahren Sie mit der nächsten Registerkarte fort.

Konfigurieren von Standardoptionen

1. Lassen Sie **Alle Faktorstufen verwenden** deaktiviert, um die erste prinzipielle Komponente zu überspringen, die über die größte Streuung in den Daten verfügt. Aktivieren Sie dieses Kontrollkästchen, um die erste prinzipielle Komponente zu behalten.
2. Aktivieren Sie **Eingabedaten bewerten**, um eine Spalte für die Modellvorhersage (Punktzahl) für Eingabedaten hinzuzufügen.
3. Wählen Sie die geeignete **Transformation** für die Trainingsdaten aus.

Demean	Subtrahiert den Mittelwert jeder Spalte.
Descal	Dividiert durch die Standardabweichung jeder Spalte.
Keine	
Normalisieren	Subtrahiert den Mittelwert jeder Spalte und dividiert jede Spalte durch ihren Bereich (Maximum minus Minimum).
Standardisieren	Verwendet keine Mittelwert- und Einheitenvarianz. Standardeinstellung.
4. Geben Sie an, wie mit **Fehlenden Daten** umgegangen werden soll, indem Sie **Überspringen** aktivieren oder **Mittelwerte zuschreiben**, wodurch der Mittelwert für fehlende Daten hinzugefügt wird.
5. Klicken Sie auf **OK**, um das Modell und die Konfiguration zu speichern, oder fahren Sie mit der nächsten Registerkarte fort.

Konfigurieren erweiterter Optionen

1. Lassen Sie **Konstante Felder ignorieren** aktiviert, damit Felder übersprungen werden, die für die einzelnen Datensätze die gleichen Werte enthalten.
2. Wählen Sie eine **PCA-Methode** aus der Dropdown-Liste aus. Beachten Sie, dass „GLRM“ und „Potenz“ derzeit experimentell sind.

GLRM	Passt ein generalisiertes niederrangiges Modell mit L2-Verlustfunktion und ohne Regularisierung an. Löst für die SVD mithilfe lokaler Matrixalgebra. Diese Option ist nur aktiviert, wenn Sie auf der Registerkarte „Standardoptionen“ Alle Faktorstufen verwenden aktiviert haben.
GramSVD	Verwendet eine verteilte Berechnung der Gram-Matrix, gefolgt von einer lokalen SVD mithilfe des JAMA-Pakets.

- Potenz** Berechnet die SVD mithilfe der Potenzmethode.
- Randomisiert** Verwendet die randomisierte Unterraumiterationsmethode.

3. Lassen Sie **Maximale Iterationen** deaktiviert, um eine unbegrenzte Anzahl von Trainingsiterationen einzustellen (Standard). Aktivieren Sie das Kontrollkästchen, um die Anzahl der Trainingsiterationen zu begrenzen.
4. Klicken Sie auf **OK**, um das Modell und die Konfiguration zu speichern, oder fahren Sie mit der nächsten Registerkarte fort.

Modellausgabe

Auf dieser Registerkarte werden die Metriken angezeigt, mit denen Sie das angepasste Modell bewerten. Diese Felder können nicht bearbeitet werden.

Nachdem Sie Ihren Auftrag ausgeführt haben, wird das resultierende Modell auf dem Spectrum™ Technology Platform-Server gespeichert. Klicken Sie auf die Schaltfläche **Ausgabe**, um die Ausgabe erneut zu generieren, und klicken Sie auf **Modelldetails**, um in dem Tool für die „Machine Learning“-Modellverwaltung die gesamte Ausgabe anzuzeigen.

8 - Random Forest Classification

In this section

Einführung in Random Forest Classification	33
Definieren von Modelleigenschaften	33
Konfigurieren von Standardoptionen	34
Konfigurieren erweiterter Optionen	35
Modellausgabe	37

Einführung in Random Forest Classification

Mit Random Forest Classification können Sie Machine Learning durchführen, indem Sie Modelle aus Datasets erstellen, die im Hinblick auf Eingabevariablen binäre oder multinomiale Ziele verwenden.

Sie müssen zunächst die Registerkarte „Modelleigenschaften“ ausfüllen, um Ihr Modell erstellen zu können. Auf den Registerkarten „Standardoptionen“ und „Erweiterte Optionen“ werden genügend Standardeinstellungen zum Abschließen eines Auftrags bereitgestellt. Sie können diese Einstellungen jedoch auch ändern und an Ihre Bedürfnisse anpassen. Anschließend führen Sie Ihren Auftrag aus, und auf der Registerkarte „Modellausgabe“ wird eine eingeschränkte Version des resultierenden Modells angezeigt. Die vollständige Ausgabe ist im Tool für die „Machine Learning“-Modellverwaltung verfügbar.

Anmerkung: Klicken Sie [hier](#), um weitere Informationen zu Random Forest Classification und zugehörige Optionen zu erhalten.

Definieren von Modelleigenschaften

1. Klicken Sie unter **Primäre Schritte/Bereitgestellte Schritte/Machine Learning** auf den **Random Forest Classification**-Schritt und ziehen Sie ihn auf die Arbeitsfläche. Platzieren Sie ihn an einer beliebigen Stelle im Datenfluss und verbinden Sie ihn mit anderen Schritten. Beachten Sie, dass der Eingabeschritt die Datenquelle sein muss, die Ziel- und Eingabevariablenfelder für Ihr Modell enthält. Ein Ausgabeschritt ist nicht erforderlich, es sei denn, Sie wählen auf der Registerkarte „Standardoptionen“ die Eingabedatenoption „Bewerten“ aus. Sie können auch eine Verbindung zu einem Ausgabeschritt herstellen, wenn Sie Ihre Ausgabe unabhängig von dem Tool für die Machine Learning-Modellverwaltung erfassen möchten.
2. Doppelklicken Sie auf den „Random Forest Classification“-Schritt, um das Dialogfeld „**Random Forest Classification**“-Optionen anzuzeigen.
3. Geben Sie einen **Modellnamen** ein, wenn Sie nicht den Standardnamen verwenden möchten.
4. Aktivieren Sie das Kontrollkästchen **Überschreiben**, um das vorhandene Modell mit neuen Daten zu überschreiben.
5. Klicken Sie auf das Dropdown-Menü **Zielfeld**, und wählen Sie ein numerisches Feld aus.
6. Klicken Sie auf **Multinomiale Ebenen**, um drei oder mehr vorhandene Kategorien in das Zielfeld einzugeben. Beachten Sie, dass eine Aktivierung dieses Felds das Feld **Eingabedaten bewerten** deaktiviert.
7. Geben Sie eine **Beschreibung** des Modells ein.

8. Klicken Sie bei jedem Feld, dessen Daten zum Modell hinzugefügt werden sollen, auf **Einschließen**. Achten Sie darauf, das Feld einzuschließen, das Sie als Zielfeld ausgewählt haben.
9. Geben Sie über das Dropdown-Menü **Modelldatentyp** an, ob die einzelnen Eingabefelder als numerisches Feld, als kategorisches Feld oder als DateTime-Feld verwendet werden sollen.
10. Klicken Sie auf **OK**, um das Modell und die Konfiguration zu speichern, oder fahren Sie mit der nächsten Registerkarte fort.

Konfigurieren von Standardoptionen

1. Geben Sie die maximale **Anzahl von Strukturen** in Ihrem Modell ein.
2. Geben Sie die **Maximale Tiefe** ein: die maximale Anzahl von Ebenen, die Ihr Modell enthalten soll.
3. Geben Sie die **Minimale Anzahl von Zeilen** ein: die minimale Anzahl von Zeilen (oder Datensätzen), die Ihr Modell enthalten soll.
4. Geben Sie die **Anzahl numerischer Bins** ein: die Anzahl von Bins, die das Histogramm erstellen und dann am besten Punkt aufteilen soll.
5. Geben Sie die **Anzahl von Bins auf höchster Ebene** ein: die minimale Anzahl von Bins, die Sie auf der Stammebene haben möchten.
6. Geben Sie die **Anzahl kategorischer Bins** ein: die maximale Anzahl von Bins, die das Histogramm erstellen und dann am besten Punkt aufteilen soll.
7. Aktivieren Sie **Abtaste**, und geben Sie den Prozentsatz der Zeilen ein, die als Stichprobe in jeder Struktur verwendet werden soll. Dies kann ein Wert zwischen 0,0 und 999 sein.
8. Aktivieren Sie **Spaltenabtastrate pro Struktur**, und geben Sie die Spaltenabtastrate für die einzelnen Strukturen ein. Dies kann ein Wert zwischen 0,0 und 1,0 sein.
9. Aktivieren Sie **Spalten in jeder Ebene**, und geben Sie die relative Änderung der Spaltenabtastrate für jede Ebene ein. Gültige Werte liegen zwischen 1,0 und der Zahl des ausgewählten Eingabeprediktors. Der Standardwert ist 1,0.
10. Aktivieren Sie **Eingabedaten bewerten**, um eine Spalte für die Modellvorhersage (Punktzahl) für Eingabedaten hinzuzufügen.
11. Geben Sie für den **Prozentsatz für Trainingsdaten** einen Wert zwischen 1 und 100 an, wenn die Eingabedaten zufällig in Stichproben für Trainings- und Testdaten aufgeteilt werden.
12. Geben Sie den Wert 100 abzüglich der in Schritt 5 eingegebenen Menge als **Prozentsatz für Testdaten** ein.
13. **Ausgangswert für Stichprobe**, um sicherzustellen, dass die Darstellung der Daten bei jeder Datenflussausführung gleich ist, wenn diese in Test- und Trainingsdaten aufgeteilt werden. Deaktivieren Sie dieses Feld, damit die Aufteilung bei jeder Datenflussausführung beliebig erfolgt.

14. Klicken Sie auf **OK**, um das Modell und die Konfiguration zu speichern, oder fahren Sie mit der nächsten Registerkarte fort.

Konfigurieren erweiterter Optionen

1. Lassen Sie **Konstante Felder ignorieren** aktiviert, damit Felder übersprungen werden, die für die einzelnen Datensätze die gleichen Werte enthalten.
2. Aktivieren Sie **Klassen ausgleichen**, um die Klassenverteilung auszugleichen und entweder für die Mehrheitsklassen ein Undersampling oder für die Minderheitsklassen ein Oversampling durchzuführen.
3. Wählen Sie einen **Histogrammtyp** aus.

Auto Für Buckets wird ein Binning vom Minimum bis zum Maximum in Schritten von $(\max - \min) / N$ durchgeführt. Verwenden Sie diese Option, um den Histogrammtyp für das Auffinden optimaler Teilungspunkte anzugeben.

QuantilesGlobal Buckets haben die gleiche Population. Dies berechnet `nbins` Quantile für jede numerische (nicht binäre) Spalte. Dann wird jeder Bucket (zwischen zwei Quantilen) einheitlich angepasst (zufällig für Reste), sodass sich insgesamt `nbins_top_level` Bins ergeben.

Random Der Algorithmus nimmt Stichproben von $N-1$ Punkten von Minimum bis Maximum und verwendet die sortierte Liste, um die beste Teilung zu finden.

RoundRobin Der Algorithmus wechselt durch alle Histogrammtypen (einer pro Struktur).

UniformAdaptive Jedes Feature wird per Binning einem Bucket zugeordnet, sodass sich Buckets mit gleicher Schrittgröße (nicht Population) ergeben. Dies ist die schnellste Methode, kann aber zu ungenaueren Aufteilungen führen, wenn die Verteilung sehr verzerrt ist.

4. Wählen Sie eine **Kategorische Codierung** aus.

Auto Führt automatisch eine `Enum`-Codierung durch.

Binary Konvertiert Kategorien in Ganzzahlen, dann in Binärwerte, und weist jeder Ziffer eine separate Spalte zu. Codiert die Daten in weniger Dimensionen, jedoch werden Entfernungen etwas verzerrt.

Anmerkung: Pro kategorischem Feature können nicht mehr als 32 Spalten vorhanden sein.

Eigen k Spalten pro kategorischem Feature, behält nur Projektionen einer 1-aus- n -codierten Matrix auf k -dimensionalen Eigen-Raum bei.

- Enum** Wechselt durch alle Histogrammtypen (einer pro Struktur).
- OneHotExplicit** Pro Kategorie ist eine Spalte vorhanden, wobei „1“ oder „0“ in jeder Zelle anzeigen, ob die Zeile die Kategorie dieser Spalte enthält.

5. Lassen Sie **Seed für Algorithmus und n-fach** aktiviert, und geben Sie einen numerischen Ausgangswert ein, um sicherzustellen, dass die Darstellung der Daten bei jeder Datenflussausführung gleich ist, wenn diese in Test- und Trainingsdaten aufgeteilt werden. Deaktivieren Sie dieses Feld, damit die Aufteilung bei jeder Datenflussausführung beliebig erfolgt.
6. Aktivieren Sie **N-fach** und geben Sie die Anzahl der Folds ein, wenn Sie eine Kreuzvalidierung durchführen.
7. Aktivieren Sie **Foldzuweisung**, und wählen Sie aus der Dropdown-Liste aus, ob Sie eine Kreuzvalidierung durchführen. Dieses Feld ist nur anwendbar, wenn Sie unter **N-fach** einen Wert eingegeben haben und **Foldfeld** nicht angegeben ist.

Auto Lässt zu, dass der Algorithmus automatisch eine Option auswählt; derzeit wird „Random“ verwendet.

Modulo Teilt das Dataset gleichmäßig auf die Folds auf und ist nicht vom Ausgangswert abhängig.

Random Teilt die Daten zufällig in „N-fach“-Bestandteile ein; diese Einstellung ist für umfangreiche Datasets am besten geeignet.

Stratified Schichtet die Folds basierend auf der Antwortvariable für Klassifizierungsprobleme. Verteilt Beobachtungen aus den verschiedenen Klassen gleichmäßig auf alle Datasets, wenn ein Dataset in Trainings- und Testdaten aufgeteilt wird. Dies kann nützlich sein, wenn viele Klassen vorhanden sind und das Dataset relativ klein ist.

8. Wenn Sie eine Kreuzvalidierung durchführen, aktivieren Sie **Foldfeld** und wählen Sie aus der Dropdown-Liste das Feld aus, das die Foldindexzuweisung für die Kreuzvalidierung enthält. Dieses Feld ist nur anwendbar, wenn Sie unter **N-fach** und **Foldzuweisung** keinen Wert eingegeben haben.
9. Aktivieren Sie **Runden stoppen**, um das Training zu beenden, wenn die Option „Stopping_metric“ sich nicht in der angegebenen Anzahl von Trainingsrunden verbessert, und geben Sie die Anzahl nicht erfolgreicher Trainingsrunden ein, die absolviert werden, bevor gestoppt werden soll. Um diese Funktion zu deaktivieren, geben Sie „0“ an. Die Metrik wird anhand der Überprüfungsdaten berechnet (falls vorhanden), ansonsten werden Trainingsdaten verwendet.
10. Wählen Sie eine **Abbruchmetrik**, um festzulegen, wann die Erstellung neuer Strukturen eingestellt werden soll.

AUC Fläche unter ROC-Kurve.

Anmerkung: Gilt nur für binomiale Modelle.

Auto	Standardwert ist <i>Abweichung</i> .
Lifftopgroup	Beste 1 %.
Logloss	Logarithmischer Abfall.
Meanperclasserror	Die Fehlklassifizierungsrate.
Misclassification	Der Wert von $(1 - (\text{korrekte Vorhersagen}/\text{gesamte Vorhersagen})) * 100$.
MSE	Mittlerer quadratischer Fehler, berücksichtigt sowohl Streuung als auch Tendenz des Prädiktors.
RMSE	Wurzel aus dem mittleren quadratischen Fehler; misst die Differenz zwischen Werten (Stichproben- und Populationswerte), die von einem Modell oder einem Schätzwert vorhergesagt wurden, und tatsächlich beobachteten Werten. Auch Quadratwurzel von MSE.

11. Aktivieren Sie **Abbruchtoleranz**, und geben Sie einen Wert ein, um die relative Toleranz für den metrikbasierten Abbruch des Trainings zu spezifizieren, wenn die Verbesserung geringer ist als dieser Wert. Dieses Feld ist nur aktiviert, wenn Sie **Runden stoppen** aktiviert haben.
12. Aktivieren Sie **Minimale Aufteilungsverbesserung**, und geben Sie einen Wert ein, um die minimale relative Verbesserung in der Verringerung des quadratischen Fehlers anzugeben, bei der eine Aufteilung durchgeführt werden soll. Wenn diese Option richtig ausgeführt wird, kann die Überanpassung verringert werden. Optimale Werte bewegen sich im Bereich von $1e-10$ bis $1e-3$. Dieses Feld ist nur aktiviert, wenn Sie **Runden stoppen** aktiviert haben.
13. Klicken Sie auf **OK**, um das Modell und die Konfiguration zu speichern, oder fahren Sie mit der nächsten Registerkarte fort.

Modellausgabe

Auf dieser Registerkarte werden die Metriken angezeigt, mit denen Sie das angepasste Modell bewerten. Diese Felder können nicht bearbeitet werden. Die Spalte „Training“ enthält immer Daten. Wenn Sie auf der Registerkarte „Standardoptionen“ eine Aufteilung in Training/Test ausgewählt haben, wird die Spalte „Test“ ebenfalls aufgefüllt. Eine Ausnahme besteht, wenn Sie auf der Registerkarte „Erweiterte Optionen“ eine Validierung vom Typ „N-fach“ ausgewählt haben. In diesem Fall wird die Spalte „N-fach“ aufgefüllt.

Nachdem Sie Ihren Auftrag ausgeführt haben, wird das resultierende Modell auf dem Spectrum™ Technology Platform-Server gespeichert. Klicken Sie auf die Schaltfläche **Ausgabe**, um die Ausgabe erneut zu generieren, und klicken Sie auf **Modelldetails**, um in dem Tool für die „Machine Learning“-Modellverwaltung die gesamte Ausgabe anzuzeigen.

9 - Random Forest Regression

In this section

Einführung in Random Forest Regression	39
Definieren von Modelleigenschaften	39
Konfigurieren von Standardoptionen	40
Konfigurieren erweiterter Optionen	41
Modellausgabe	43

Einführung in Random Forest Regression

Mit Random Forest Regression können Sie Machine Learning durchführen, indem Sie Modelle aus Datasets erstellen, die im Hinblick auf Eingabevariablen kontinuierliche Ziele verwenden.

Sie müssen zunächst die Registerkarte „Modelleigenschaften“ ausfüllen, um Ihr Modell erstellen zu können. Auf den Registerkarten „Standardoptionen“ und „Erweiterte Optionen“ werden genügend Standardeinstellungen zum Abschließen eines Auftrags bereitgestellt. Sie können diese Einstellungen jedoch auch ändern und an Ihre Bedürfnisse anpassen. Anschließend führen Sie Ihren Auftrag aus, und auf der Registerkarte „Modellausgabe“ wird eine eingeschränkte Version des resultierenden Modells angezeigt. Die vollständige Ausgabe ist im Tool für die „Machine Learning“-Modellverwaltung verfügbar.

Anmerkung: Klicken Sie [hier](#), um weitere Informationen zu Random Forest Regression und zugehörige Optionen zu erhalten.

Definieren von Modelleigenschaften

1. Klicken Sie unter **Primäre Schritte/Bereitgestellte Schritte/Machine Learning** auf den **Random Forest Regression**-Schritt und ziehen Sie ihn auf die Arbeitsfläche. Platzieren Sie ihn an einer beliebigen Stelle im Datenfluss und verbinden Sie ihn mit anderen Schritten. Beachten Sie, dass der Eingabeschritt die Datenquelle sein muss, die Ziel- und Eingabevariablenfelder für Ihr Modell enthält. Ein Ausgabeschritt ist nicht erforderlich, es sei denn, Sie wählen auf der Registerkarte „Standardoptionen“ die Eingabedatenoption „Bewerten“ aus. Sie können auch eine Verbindung zu einem Ausgabeschritt herstellen, wenn Sie Ihre Ausgabe unabhängig von dem Tool für die Machine Learning-Modellverwaltung erfassen möchten.
2. Doppelklicken Sie auf den „Random Forest Regression“-Schritt, um das Dialogfeld „**Random Forest Regression**“-Optionen anzuzeigen.
3. Geben Sie einen **Modellnamen** ein, wenn Sie nicht den Standardnamen verwenden möchten.
4. Optional: Aktivieren Sie das Kästchen **Überschreiben**, um das vorhandene Modell mit neuen Daten zu überschreiben.
5. Klicken Sie auf das Dropdown-Menü **Zielfeld**, und wählen Sie ein numerisches Feld aus.
6. Optional: Geben Sie einen **Beschreibung** des Modells ein.
7. Klicken Sie bei jedem Feld, dessen Daten zum Modell hinzugefügt werden sollen, auf **Einschließen**. Achten Sie darauf, das Feld einzuschließen, das Sie als Zielfeld ausgewählt haben.
8. Geben Sie über das Dropdown-Menü **Modelltyp** an, ob die einzelnen Eingabefelder als numerisches Feld, als kategorisches Feld oder als DateTime-Feld verwendet werden sollen.

9. Klicken Sie auf **OK**, um das Modell und die Konfiguration zu speichern, oder fahren Sie mit der nächsten Registerkarte fort.

Konfigurieren von Standardoptionen

1. Geben Sie die maximale **Anzahl von Strukturen** in Ihrem Modell ein. Der Standardwert ist 50.
2. Geben Sie die **Maximale Tiefe** ein: die maximale Anzahl von Ebenen, die Ihr Modell enthalten soll. Der Standardwert ist 5.
3. Geben Sie die **Minimale Anzahl von Zeilen** ein: die minimale Anzahl von Zeilen (oder Datensätzen), die Ihr Modell enthalten soll. Der Standardwert ist 10.
4. Geben Sie die **Anzahl numerischer Bins** ein: die Anzahl von Bins, die das Histogramm erstellen und dann am besten Punkt aufteilen soll. Der Standardwert ist 20.
5. Geben Sie die **Anzahl von Bins auf höchster Ebene** ein: die minimale Anzahl von Bins, die Sie auf der Stammebene haben möchten. Der Standardwert ist 1024.
6. Geben Sie die **Anzahl kategorischer Bins** ein: die maximale Anzahl von Bins, die das Histogramm erstellen und dann am besten Punkt aufteilen soll. Der Standardwert ist 1024.
7. Aktivieren Sie **Abtaste**, und geben Sie den Prozentsatz der Zeilen ein, die als Stichprobe in jeder Struktur verwendet werden soll. Dies kann ein Wert zwischen 0,0 und 1,0 sein.
8. Aktivieren Sie **Spaltenabtaste pro Struktur**, und geben Sie die Spaltenabtaste für die einzelnen Strukturen ein. Dies kann ein Wert zwischen 0,0 und 1,0 sein.
9. **Spalten in jeder Ebene** gibt die Spalten an, die zufällig in jeder Ebene ausgewählt werden. Wenn diese Option nicht aktiviert ist, wird der Standardwert von -1 verwendet. Die Anzahl der Variablen ist dann die Wurzel der Anzahl der Spalten für die Klassifizierung und $p/3$ die Regression (wobei „p“ die Anzahl der Prädiktoren ist). Wenn Sie die Option aktivieren, können Sie einen Wert gleich oder größer als 1 angeben. Der angegebene Wert darf nicht größer sein als die Anzahl der Prädiktoren.
10. Aktivieren Sie **Eingabedaten bewerten**, um eine Spalte für die Modellvorhersage (Punktzahl) für Eingabedaten hinzuzufügen.
11. Geben Sie für den **Prozentsatz für Trainingsdaten** einen Wert zwischen 1 und 100 an, wenn die Eingabedaten zufällig in Stichproben für Trainings- und Testdaten aufgeteilt werden.
12. Geben Sie den Wert 100 abzüglich der in Schritt 5 eingegebenen Menge als **Prozentsatz für Testdaten** ein.
13. **Ausgangswert für Stichprobe**, um sicherzustellen, dass die Darstellung der Daten bei jeder Datenflussausführung gleich ist, wenn diese in Test- und Trainingsdaten aufgeteilt werden. Deaktivieren Sie dieses Feld, damit die Aufteilung bei jeder Datenflussausführung beliebig erfolgt.
14. Klicken Sie auf **OK**, um das Modell und die Konfiguration zu speichern, oder fahren Sie mit der nächsten Registerkarte fort.

Konfigurieren erweiterter Optionen

1. Lassen Sie **Konstante Felder ignorieren** aktiviert, damit Felder übersprungen werden, die für die einzelnen Datensätze die gleichen Werte enthalten.

2. Wählen Sie einen **Histogrammtyp** aus.

Auto Für Buckets wird ein Binning vom Minimum bis zum Maximum in Schritten von $(\text{max-min})/N$ durchgeführt. Verwenden Sie diese Option, um den Histogrammtyp für das Auffinden optimaler Teilungspunkte anzugeben.

QuantilesGlobal Buckets haben die gleiche Population. Dies berechnet `nbins` Quantile für jede numerische (nicht binäre) Spalte. Dann wird jeder Bucket (zwischen zwei Quantilen) einheitlich angepasst (zufällig für Reste), sodass sich insgesamt `nbins_top_level` Bins ergeben.

Random Der Algorithmus nimmt Stichproben von $N-1$ Punkten von Minimum bis Maximum und verwendet die sortierte Liste, um die beste Teilung zu finden.

RoundRobin Der Algorithmus wechselt durch alle Histogrammtypen (einer pro Struktur).

UniformAdaptive Jedes Feature wird per Binning einem Bucket zugeordnet, sodass sich Buckets mit gleicher Schrittgröße (nicht Population) ergeben. Dies ist die schnellste Methode, kann aber zu ungenaueren Aufteilungen führen, wenn die Verteilung sehr verzerrt ist.

3. Wählen Sie eine **Kategorische Codierung** aus.

Auto Führt automatisch eine `Enum`-Codierung durch.

Binary Konvertiert Kategorien in Ganzzahlen, dann in Binärwerte, und weist jeder Ziffer eine separate Spalte zu. Codiert die Daten in weniger Dimensionen, jedoch werden Entfernungen etwas verzerrt.

Anmerkung: Pro kategorischem Feature können nicht mehr als 32 Spalten vorhanden sein.

Eigen k Spalten pro kategorischem Feature, behält nur Projektionen einer 1-aus-n-codierten Matrix auf k -dimensionalen Eigen-Raum bei.

Enum Wechselt durch alle Histogrammtypen (einer pro Struktur).

OneHotExplicit Pro Kategorie ist eine Spalte vorhanden, wobei „1“ oder „0“ in jeder Zelle anzeigen, ob die Zeile die Kategorie dieser Spalte enthält.

4. Lassen Sie **Seed für Algorithmus und n-fach** aktiviert, und geben Sie einen numerischen Ausgangswert ein, um sicherzustellen, dass die Darstellung der Daten bei jeder

Datenflussausführung gleich ist, wenn diese in Test- und Trainingsdaten aufgeteilt werden. Deaktivieren Sie dieses Feld, damit die Aufteilung bei jeder Datenflussausführung beliebig erfolgt.

5. Aktivieren Sie **N-fach** und geben Sie die Anzahl der Folds ein, wenn Sie eine Kreuzvalidierung durchführen.
6. Aktivieren Sie **Foldzuweisung**, und wählen Sie aus der Dropdown-Liste aus, ob Sie eine Kreuzvalidierung durchführen. Dieses Feld ist nur anwendbar, wenn Sie unter **N-fach** einen Wert eingegeben haben und **Foldfeld** nicht angegeben ist.

Auto Lässt zu, dass der Algorithmus automatisch eine Option auswählt; derzeit wird „Random“ verwendet.

Modulo Teilt das Dataset gleichmäßig auf die Folds auf und ist nicht vom Ausgangswert abhängig.

Random Teilt die Daten zufällig in „N-fach“-Bestandteile ein; diese Einstellung ist für umfangreiche Datasets am besten geeignet.

7. Wenn Sie eine Kreuzvalidierung durchführen, aktivieren Sie **Foldfeld** und wählen Sie aus der Dropdown-Liste das Feld aus, das die Foldindexzuweisung für die Kreuzvalidierung enthält. Dieses Feld ist nur anwendbar, wenn Sie unter **N-fach** und **Foldzuweisung** keinen Wert eingegeben haben.
8. Aktivieren Sie **Runden stoppen**, um das Training zu beenden, wenn die Option „Stopping_metric“ sich nicht in der angegebenen Anzahl von Trainingsrunden verbessert, und geben Sie die Anzahl nicht erfolgreicher Trainingsrunden ein, die absolviert werden, bevor gestoppt werden soll. Um diese Funktion zu deaktivieren, geben Sie „0“ an. Die Metrik wird anhand der Überprüfungsdaten berechnet (falls vorhanden), ansonsten werden Trainingsdaten verwendet.
9. Wählen Sie eine **Abbruchmetrik**, um festzulegen, wann die Erstellung neuer Strukturen eingestellt werden soll.

Auto Standardwert ist *Abweichung*.

Abweichung Mittlere Restabweichung; identisch mit MSE.

MAE Mittlerer absoluter Fehler; die Differenz zwischen zwei kontinuierlichen Variablen.

MSE Mittlerer quadratischer Fehler, berücksichtigt sowohl Streuung als auch Tendenz des Prädiktors.

RMSE Wurzel aus dem mittleren quadratischen Fehler; misst die Differenz zwischen Werten (Stichproben- und Populationswerte), die von einem Modell oder einem Schätzwert vorhergesagt wurden, und tatsächlich beobachteten Werten. Auch Quadratwurzel von MSE.

RMSLE Wurzel des mittleren quadratischen logarithmischen Fehlers; misst das Verhältnis zwischen vorhergesagtem und tatsächlichem Wert.

10. Aktivieren Sie **Abbruchtoleranz**, und geben Sie einen Wert ein, um die relative Toleranz für den metrikbasierten Abbruch des Trainings zu spezifizieren, wenn die Verbesserung geringer ist als dieser Wert.
11. Aktivieren Sie **Minimale Aufteilungsverbesserung**, und geben Sie einen Wert ein, um die minimale relative Verbesserung in der Verringerung des quadratischen Fehlers anzugeben, bei der eine Aufteilung durchgeführt werden soll. Wenn diese Option richtig ausgeführt wird, kann die Überanpassung verringert werden. Optimale Werte bewegen sich im Bereich von $1e-10$ bis $1e-3$. Dieses Feld ist nur aktiviert, wenn Sie **Runden stoppen** aktiviert haben.
12. Klicken Sie auf **OK**, um das Modell und die Konfiguration zu speichern, oder fahren Sie mit der nächsten Registerkarte fort.

Modellausgabe

Auf dieser Registerkarte werden die Metriken angezeigt, mit denen Sie das angepasste Modell bewerten. Diese Felder können nicht bearbeitet werden. Die Spalte „Training“ enthält immer Daten. Wenn Sie auf der Registerkarte „Standardoptionen“ eine Aufteilung in Training/Test ausgewählt haben, wird die Spalte „Test“ ebenfalls aufgefüllt. Eine Ausnahme besteht, wenn Sie auf der Registerkarte „Erweiterte Optionen“ eine Validierung vom Typ „N-fach“ ausgewählt haben. In diesem Fall wird die Spalte „N-fach“ aufgefüllt.

Nachdem Sie Ihren Auftrag ausgeführt haben, wird das resultierende Modell auf dem Spectrum™ Technology Platform-Server gespeichert. Klicken Sie auf die Schaltfläche **Ausgabe**, um die Ausgabe erneut zu generieren, und klicken Sie auf **Modelldetails**, um in dem Tool für die „Machine Learning“-Modellverwaltung die gesamte Ausgabe anzuzeigen.

10 - Java Model Scoring

In this section

Einführung in Java Model Scoring	45
Definieren von Modelleigenschaften	45
Modellausgabe	46

Einführung in Java Model Scoring

Java Model Scoring ermöglicht Ihnen die Bewertung neuer Daten mithilfe der Formel, die beim Anpassen eines Machine Learning-Modells erstellt wird.

Anmerkung: Modelle müssen zunächst über die Machine Learning-Modellverwaltung verfügbar gemacht werden, bevor sie im „Java Model Scoring“-Schritt verfügbar werden. Weitere Informationen finden Sie unter [Einführung in die Machine Learning-Modellverwaltung](#) auf Seite 48.

Sie müssen zwei Registerkarten im Dialogfeld „**Java Model Scoring**“-**Optionen** ausfüllen, um Ihre Daten bewerten zu können. Identifizieren Sie zunächst das Modell und dessen Typ, und stellen Sie anschließend sicher, dass die Felder des Modells den Spectrum™ Technology Platform-Feldern ordnungsgemäß zugeordnet wurden. Im Anschluss können Sie die Ausgabe konfigurieren, indem Sie auswählen, welche Felder eingeschlossen werden sollen, und Ihren Auftrag ausführen. Die Registerkarte **Modellausgabe** enthält das Mapping für Datentypen der Spectrum™ Technology Platform und Ihr Modell.

Wenn Ihr Auftrag einen Schritt enthält, der die Ausgabe in einer Datei oder Tabelle erfasst, können Sie diese Ausgabe in einem nachfolgenden Datenfluss oder Webservice verwenden.

Definieren von Modelleigenschaften

1. Klicken Sie unter **Primäre Schritte/Bereitgestellte Schritte/Advanced Analytics** auf den **Java Model Scoring**-Schritt und ziehen Sie ihn auf die Arbeitsfläche. Platzieren Sie ihn an einer beliebigen Stelle im Datenfluss und verbinden Sie ihn mit Eingabe- und Ausgabeschritten. Beachten Sie, dass der Eingabeschritt die Datenquelle sein muss, die Ziel- und Eingabevariablenfelder für Ihr Modell enthält. Wenn Sie Ihren Auftrag im Batchmodus ausführen, benötigen Sie zudem einen Ausgabeschritt für die Erfassung von Modellbewertungen; andernfalls können Sie Daten mithilfe eines Spectrum™ Technology Platform-Webservice in Echtzeit bewerten.
2. Doppelklicken Sie auf den „Java Model Scoring“-Schritt, um das Dialogfenster „**Model Scoring**“-**Optionen** anzuzeigen.
3. Optional: Wählen Sie im Dropdown-Menü **Typfilter** den Typ eines Modells aus, das Sie bewerten.
4. Wählen Sie den **Typfilter** aus, der für die Bewertung des Modells verwendet wird.
5. Wählen Sie den **Modellnamen** aus dem Dropdown-Menü aus.
6. Geben Sie im Feld **Modelltyp** den Typ des Modells aus, das Sie bewerten.
7. Optional: Geben Sie eine **Beschreibung** des Modells ein.

8. Aktivieren Sie **Unbekannte kategorische Ebenen ignorieren**, um Daten in der Spalte „Predicted_Value“ für Zeilen mit einer kategorischen Ebene zurückzugeben, die nicht in den Daten zur Anpassung des Modells enthalten war. Wenn Sie dieses Kästchen deaktiviert lassen, wird in der Spalte „Predicted_Value“ für diese Zeilen „Null/NA“ zurückgegeben.
9. Die Tabelle **Eingaben** enthält Informationen zu den Eingabefeldern des Modells. Diese Felder und die zugehörigen Datentypen ordnet Spectrum automatisch Felder und Datentypen zu.
10. Klicken Sie auf **OK**, um diese Optionen zu speichern, oder fahren Sie mit der nächsten Registerkarte fort.

Modellausgabe

Die Tabelle **Ausgaben** enthält Informationen zu den Ausgabefeldern des Modells. Diese Felder und die zugehörigen Datentypen ordnet Spectrum automatisch Felder und Datentypen zu.

1. Klicken Sie bei jedem Feld, dessen Daten im Modell eingeschlossen werden sollen, auf **Einschließen**.
2. Klicken Sie auf **OK**, um das Modell zu speichern.

11 - Machine Learning-Modellverwaltung

In this section

Einführung in die Machine Learning-Modellverwaltung	48
Die Registerkarte „Modelldetail“	49

Einführung in die Machine Learning-Modellverwaltung

Auf der Registerkarte „Modellanalyse“ bei der Machine Learning-Modellverwaltung wird eine Liste mit allen Machine Learning-Modellen auf Ihrem Spectrum™ Technology Platform-Server angezeigt. Sie können diese Liste filtern, indem Sie eine Zeichenfolge in das Textfeld eingeben. Jedes Feld in der Tabelle wird dann nach dieser Zeichenfolge durchsucht.

In diesen Modellen können mehrere Vorgänge durchgeführt werden. Sie können Modelle verfügbar machen, die Verfügbarkeit von Modellen aufheben oder Modelle löschen. Mit den verfügbar gemachten Modellen werden im „Java Model Scoring“-Schritt mithilfe der beim Anpassen der Machine Learning-Modelle erstellten Formeln neue Daten bewertet. Zusätzlich können Sie zu jedem Modell detaillierte Informationen anzeigen. Welche Details zurückgegeben werden, hängt von dem Modelltyp ab, dessen Daten Sie anzeigen. Schließlich können Sie zwei beliebige Modelle des gleichen Typs miteinander vergleichen. Bei diesem Vergleich werden für alle Modelle, die Sie miteinander vergleichen, die gleichen Informationen nebeneinander angezeigt, die auf der Registerkarte „Modelldetail“ enthalten sind.

Zugreifen auf die Modellanalyse bei der Machine Learning-Modellverwaltung

Es gibt drei Möglichkeiten für den Zugriff auf die Machine Learning-Modellverwaltung:

- Verwenden Sie die Begrüßungsseite der Spectrum™ Technology Platform:
 - Öffnen Sie einen Webbrowser und navigieren Sie zur Spectrum™ Technology Platform-Begrüßungsseite unter:

http://<servername>:<port>

 Wenn Sie beispielsweise Spectrum™ Technology Platform auf einem Computer mit dem Namen „myspectrumplatform“ installiert haben und dieser den HTTP-Standardport 8080 verwendet, navigieren Sie zu:

http://myspectrumplatform:8080
 - Klicken Sie auf **Spectrum Machine Learning**.
 - Klicken Sie auf „**Machine Learning**“-**Datenbank öffnen**.
- Klicken Sie bei einem der Schritte zur Modellerstellung auf **Für Modelldetails hier klicken**.
- Verwenden Sie einen Webbrowser:
 - Öffnen Sie einen Webbrowser und navigieren Sie zur Seite „Machine Learning-Modellverwaltung“ der Spectrum™ Technology Platform unter:

http://<Servername>:<Port>/machinelearning






Wenn Sie beispielsweise Spectrum™ Technology Platform auf einem Computer mit dem Namen „myspectrumplatform“ installiert haben und dieser den HTTP-Standardport 8080 verwendet, navigieren Sie zu:

http://myspectrumplatform:8080/machinelearning

- Geben Sie einen gültigen Benutzernamen und das zugehörige Kennwort für die Spectrum™ Technology Platform ein.
- Klicken Sie auf die Registerkarte **Modellanalyse**, wenn das Tool geöffnet wird.

Vorgänge zur Modellverwaltungsanalyse

Führen Sie die folgenden Vorgänge durch, indem Sie ein Modell auswählen und auf die entsprechende Schaltfläche klicken:

	Machen Sie das Modell verfügbar, damit es für den „Java Model Scoring“-Schritt verfügbar ist. Wenn ein Modell nicht verfügbar gemacht wird, kann es nicht für Bewertungen verwendet werden.
	Heben Sie die Verfügbarkeit des Modells auf.
	Löschen Sie das Modell. Anmerkung: Ein verfügbar gemachtes Modell kann nicht gelöscht werden. Zu diesem Zeitpunkt ist jedoch keine inhärente Sicherheit vorhanden, durch die verhindert wird, dass ein Benutzer die Modelle eines anderen Benutzers löscht.
	Zeigen Sie Details zur Modellausgabe an. Sie können auch über die „K-Means Clustering“- und „Logistic Regression“-Schritte auf diese Informationen zugreifen, indem Sie auf der Registerkarte „Modellausgabe“ auf „Für Modelldetails hier klicken“ klicken.
	Vergleichen Sie die Modelle miteinander.

Die Registerkarte „Modelldetail“

Auf dem Bildschirm „Modelldetail“ werden folgende Informationen für ein Modell angezeigt:

- **Modellname:** Der Name des Modells
- **Modelltyp:** Der Typ des Machine Learning-Modells

- **Benutzer:** Der Benutzername der Person, die das Modell erstellt hat
- **Beschreibung:** Die Beschreibung des Modells, wenn bei der Erstellung des Modells eine Beschreibung angegeben wurde
- **Status:** Gibt an, ob das Modell verfügbar gemacht wurde oder ob die Verfügbarkeit aufgehoben wurde
- **Datenflussname:** Der Name des Datenflusses, der das Modell erzeugt hat
- **Erstellungszeit:** Das Datum und die Uhrzeit der Modellerstellung

Zusätzliche Details werden auf Grundlage des Modelltyps bereitgestellt.

„K-Means Clustering“-Details

Auf dem Bildschirm „Modelldetail“ werden folgende Informationen für „K-Means Clustering“-Modelle angezeigt:

Modellübersicht

- Anzahl der Zeilen
- Anzahl der Cluster
- Anzahl der kategorischen Spalten
- Anzahl der Iterationen
- Innerhalb der Cluster-Summe von Quadraten
- Gesamtsumme von Quadraten
- Zwischen der Cluster-Summe von Quadraten

Metriken

Stellt für folgende Informationen Trainings-, Test- und „N-fach“-Daten bereit:

- Innerhalb der Cluster-Gesamtsumme von Quadraten
- Gesamtsumme von Quadraten
- Zwischen der Cluster-Summe von Quadraten

Zentroidstatistik

Stellt für die einzelnen Zentroide folgende Trainings-, Test- und „N-fach“-Daten bereit:

- Size
- Innerhalb der Cluster-Summe von Quadraten

Cluster-Mittelwerte

Stellt für jeden Zentroid detaillierte Informationen bereit. Der Inhalt variiert je nach Eingabedaten. Ein Cluster stellt eine Gruppe von Beobachtungen aus einem Dataset dar, die gemäß eines bestimmten Clustering-Algorithmus als ähnlich identifiziert wurden

Standardisierte Cluster-Mittelwerte

Stellt für jeden Zentroid standardisierte Informationen bereit. Der Inhalt variiert je nach Eingabedaten.

Details zu „Logistic Regression“

Auf dem Bildschirm „Modelldetail“ werden folgende Informationen für „Logistic Regression“-Modelle angezeigt:

Metriken

Stellt für folgende Informationen Trainings-, Test- und „N-fach“-Daten bereit:

- Mean Squared Error (MSE)
- Root Mean Squared Error (RMSE)
- Anzahl der Beobachtungen
- R-Quadrat (R²)
- Logarithmischer Abfall (Logloss)
- Area under the curve (AUC)
- Gini-Koeffizient
- Mean Per Class Error
- Akaike-Informationskriterium (AIC)
- Restabweichung
- Abweichung von Null
- Freiheitsgrad von Null
- Restfreiheitsgrad

Maximum Metrics Threshold

Gibt den Training Maximum Metrics Threshold für Trainings-, Text- und „N-fach“-Daten mithilfe der folgenden Metriken an:

- max f1
- max f2
- max f0point5
- max accuracy
- max precision
- max recall
- max specificity
- max absolute_mcc
- max min_per_class_accuracy
- max mean_per_class_accuracy

Konfusionsmatrix

Stellt die Leistung eines Modells in einer Reihe von Trainings-, Test- und „N-fach“-Daten bereit, bei denen die tatsächlichen Werte bekannt sind.

Standardisiertes Koeffizientendiagramm

Zeigt die wichtigsten Prädiktoren an, indem der relative Wert der Koeffizienten angegeben wird. Dieser gibt an, wie stark sich das Ziel durch eine Änderung der Eingabe verändert.

GLM-Koeffizienten

Zeigt Koeffizienten für ein Generalized Linear-Modell, das Regressionsmodelle für Ergebnisse nach Exponentialverteilungen schätzt.

AUC-Kurven

Area under the curve (Fläche unter der Kurve); bestimmt, welches der Modelle die Klassen mithilfe der Trainings-, Test- und „N-fach“-Daten am besten vorhersagt.

Anhebungs-/Verstärkungskurven

Wertet die Fähigkeit des binären Klassifizierungsmodells aus, mithilfe der Trainings-, Test- und „N-Fold“-Daten Vorhersagen zu treffen.

Details zu Linear Regression

Auf dem Bildschirm „Modelldetail“ werden folgende Informationen für Linear Regression-Modelle angezeigt:

Metriken

Stellt für folgende Informationen Trainings-, Test- und „N-fach“-Daten bereit:

- Mean Squared Error (MSE)
- Root Mean Squared Error (RMSE)
- Anzahl der Beobachtungen
- R-Quadrat (R²)
- Mittlere Restabweichung
- Mittlerer absoluter Fehler (MAE)
- Wurzel aus dem mittleren quadratischen Fehler (RMSE)
- Akaike-Informationskriterium (AIC)
- Restabweichung
- Abweichung von Null
- Freiheitsgrad von Null
- Restfreiheitsgrad

Standardisiertes Koeffizientendiagramm

Zeigt die wichtigsten Prädiktoren an, indem der relative Wert der Koeffizienten angegeben wird. Dieser gibt an, wie stark eine Änderung des jeweiligen Prädiktorkoeffizientenwerts den Zielwert positiv oder negativ verändert. Zudem werden die besten 25 Koeffizienten im Modell dargestellt.

GLM-Koeffizienten

Zeigt Koeffizienten für ein Generalized Linear-Modell, das Regressionsmodelle für Ergebnisse nach Exponentialverteilungen schätzt.

Details zu Random Forest Regression

Auf dem Bildschirm „Modelldetail“ werden folgende Informationen für Random Forest Regression-Modelle angezeigt:

Metriken

Stellt für folgende Informationen Trainings-, Test- und „N-fach“-Daten bereit:

- Mean Squared Error (MSE)
- Root Mean Squared Error (RMSE)
- Anzahl der Beobachtungen
- R-Quadrat (R²)
- Mittlere Restabweichung
- Mittlerer absoluter Fehler (MAE)
- Wurzel aus dem mittleren quadratischen Fehler (RMSE)

Variablenwichtigkeiten

Stellt Wichtigkeitswerte für jede Variable unter Verwendung der folgenden Metriken bereit:

- Relative Wichtigkeit
- Skalierte Wichtigkeit
- Prozentsatz

Zudem werden die besten 25 Variablen im Modell dargestellt.

Details zu Random Forest Classification – Binomial

Auf dem Bildschirm „Modelldetail“ werden folgende Informationen für **binomiale** Random Forest Classification-Modelle angezeigt:

Metriken

Stellt für folgende Informationen Trainings-, Test- und „N-fach“-Daten bereit:

- Mean Squared Error (MSE)
- Root Mean Squared Error (RMSE)
- Anzahl der Beobachtungen
- R-Quadrat (R²)
- Logloss
- Area under the curve (AUC)

- Gini
- Mean Per Class Error

Maximum Metrics Threshold

Gibt den Training Maximum Metrics Threshold für Trainings-, Test- und „N-fach“-Daten mithilfe der folgenden Metriken an:

- max f1
- max f2
- max f0point5
- max accuracy
- max precision
- max recall
- max specificity
- max absolute_mcc
- max min_per_class_accuracy
- max mean_per_class_accuracy

Konfusionsmatrix

Stellt die Leistung eines Modells in einer Reihe von Trainings-, Test- und „N-fach“-Daten bereit, bei denen die tatsächlichen Werte bekannt sind.

Variablenwichtigkeiten

Stellt Wichtigkeitswerte für jede Variable unter Verwendung der folgenden Metriken bereit:

- Relative Wichtigkeit
- Skalierte Wichtigkeit
- Prozentsatz

Zudem werden die besten 25 Variablen im Modell dargestellt.

AUC-Kurven

Area under the curve (Fläche unter der Kurve); bestimmt, welches der Modelle die Klassen mithilfe der Trainings-, Test- und „N-fach“-Daten am besten vorhersagt.

Anhebungs-/Verstärkungskurven

Wertet die Fähigkeit des binären Klassifizierungsmodells aus, mithilfe der Trainings-, Test- und „N-Fold“-Daten Vorhersagen zu treffen.

Details zu Random Forest Classification – Multinomial

Auf dem Bildschirm „Modelldetail“ werden folgende Informationen für **multinomiale** Random Forest Classification-Modelle angezeigt:

Metriken

Stellt für folgende Informationen Trainings-, Test- und „N-fach“-Daten bereit:

- Mean Squared Error (MSE)
- Root Mean Squared Error (RMSE)
- Anzahl der Beobachtungen
- R-Quadrat (R²)
- Logloss
- Mean Per Class Error

Konfusionsmatrix

Stellt die Leistung eines Modells in einer Reihe von Trainings-, Test- und „N-fach“-Daten bereit, bei denen die tatsächlichen Werte bekannt sind.

Variablenwichtigkeiten

Stellt Wichtigkeitswerte für jede Variable unter Verwendung der folgenden Metriken bereit:

- Relative Wichtigkeit
- Skalierte Wichtigkeit
- Prozentsatz

Zudem werden die besten 25 Variablen im Modell dargestellt.

Details zu Principal Component Analysis

Auf dem Bildschirm „Modelldetail“ werden folgende Informationen für PCA-Modelle angezeigt:

Bedeutung der Komponenten

Zeigt die prinzipiellen Komponenten in der Reihenfolge der Wichtigkeit an, die auf folgenden Metriken basieren:

- Standardabweichung
- Streuungsverhältnis
- Kumulatives Verhältnis

Drehung

Stellt die Matrix der Variablenbelastungen dar, die Gewichtung mit der jede standardisierte Originalvariable multipliziert werden sollte, um die Komponentenpunktzahl zu erhalten.

Notices

© 2017 Pitney Bowes Software Inc. Alle Rechte vorbehalten. MapInfo und Group 1 Software sind Marken von Pitney Bowes Software Inc. Alle anderen Marken und Markenzeichen sind Eigentum ihrer jeweiligen Besitzer.

USPS® Urheberrechtshinweise

Pitney Bowes Inc. wurde eine nicht-ausschließliche Lizenz erteilt, die die Veröffentlichung und den Verkauf von ZIP + 4® Postleitzahl-Datenbanken auf optischen und magnetischen Medien genehmigt. Folgende Marken sind Markenzeichen des United States Postal Service: CASS, CASS Certified, DPV, eLOT, FASTforward, First-Class Mail, Intelligent Mail, LACS^{Link}, NCOA^{Link}, PAVE, PLANET Code, Postal Service, POSTNET, Post Office, RDI, Suite^{Link}, United States Postal Service, Standard Mail, United States Post Office, USPS, ZIP Code, und ZIP + 4. Hierbei handelt es sich jedoch nicht um eine vollständige Liste der Marken, die zum United States Postal Service gehören.

Pitney Bowes Inc. ist nicht-exklusiver Lizenznehmer von USPS® für die Verarbeitungsprozesse von NCOA^{Link}®.

Die Preisgestaltung jeglicher Pitney Bowes Softwareprodukte, -optionen und -dienstleistungen erfolgt nicht durch USPS® oder die Regierung der Vereinigten Staaten. Es wird auch keine Regulierung oder Genehmigung der Preise durch USPS® oder die US-Regierung durchgeführt. Bei der Verwendung von RDI™-Daten zur Berechnung von Paketversandkosten wird die Entscheidung, welcher Paketlieferdienst genutzt wird, nicht von USPS® oder der Regierung der Vereinigten Staaten getroffen.

Datenbereitstellung und Hinweise

Hier verwendete Datenprodukte und Datenprodukte, die in Software-Anwendungen von Pitney Bowes verwendet werden, sind durch verschiedene Markenzeichen und mindestens eines der folgenden Urheberrechte geschützt:

- © Copyright United States Postal Service. Alle Rechte vorbehalten.
- © 2014 TomTom. Alle Rechte vorbehalten. TomTom und das TomTom Logo sind eingetragene Marken von TomTom N.V.
- © 2016 HERE

Fuente: INEGI (Instituto Nacional de Estadística y Geografía)

Basierend auf elektronischen Daten © National Land Survey Sweden.

- © Copyright United States Census Bureau
- © Copyright Nova Marketing Group, Inc.

Teile dieses Programms sind urheberrechtlich geschützt durch © Copyright 1993-2007 Nova Marketing Group Inc. Alle Rechte vorbehalten.

- © Copyright Second Decimal, LLC
- © Copyright Canada Post Corporation

Diese CD-ROM enthält Daten einer urheberrechtlich geschützten Datenerfassung der Canada Post Corporation.

© 2007 Claritas, Inc.

Das Geocode Address World Dataset enthält lizenzierte Daten des GeoNames-Projekts (www.geonames.org), die unter den Bedingungen der Creative Commons Attribution License („Attribution License“) bereitgestellt werden. Die Attribution License können Sie unter <http://creativecommons.org/licenses/by/3.0/legalcode> einsehen. Ihre Nutzung der GeoNames-Daten (wie im Spectrum™ Technology Platform Nutzerhandbuch beschrieben) unterliegt den Bedingungen der Attribution License. Bei Konflikten zwischen Ihrer Vereinbarung mit Pitney Bowes Software, Inc. und der Attribution License hat die Attribution License lediglich bezüglich der Nutzung von GeoNames-Daten Vorrang.



3001 Summer Street
Stamford CT 06926-0700
USA

www.pitneybowes.com