

Oracle Machine Learning for R

エンタープライズ・データに対してRを使用可能にする、
Oracle Machine Learningのコンポーネント

2023年4月、バージョン1.0

Copyright © 2023, Oracle and/or its affiliates

公開

本書の目的

本書では、Oracle DatabaseおよびOracle Autonomous Database用のOracle Machine Learning for Rの概要を説明しています。本書は、Oracle Machine Learning for Rのビジネス上の利点の評価と、データ・サイエンス・プロジェクトおよびITプロジェクトの計画立案を支援することのみを目的としています。

対象読者

本書は、Oracle DatabaseまたはOracle Autonomous Databaseに保存されたデータや、そこからアクセスできるデータを使用するデータの準備、機械学習、ソリューション・デプロイメントに対し、Rを大規模で使用するに関心のあるすべての方を対象としています。

目次

本書の目的	2
対象読者	2
はじめに	4
企業向けのR	4
透過レイヤー	5
インデータベースMLアルゴリズム	6
埋込みR実行	6
データベース処理	9
結論	10
参考資料	11
OML4Rのドキュメント	11
Oracle Machine Learning	11
ROracle	11

はじめに

Oracle Machine Learning for R (OML4R) は、Oracle DatabaseとOracle Autonomous Databaseで使用できるOracle Machine Learningのコンポーネントで、オープンソースのスクリプト言語および環境であるRをエンタープライズ・データに対して使用できるようにするものです。Oracle Machine Learning for Rは、Rを使用したエンド・ツー・エンドの分析プロセスにデータベース中心の環境を提供します。ユーザーは、データ探索とデータ準備のために、データベース内に存在するデータ上にあるオーバーロードされたR関数を呼び出します。ここでは、データベースを高パフォーマンスのコンピューティング環境として利用します。また、Rユーザーが強力なインデータベース機械学習アルゴリズムをRから呼び出して、アプリケーションをサポートするユーザー定義のR関数をデプロイすることもできます。

企業向けのR

Rを通じてデータベースを高パフォーマンスのコンピューティング環境として利用するには、背後で透過的にSQLに変換されるオーバーロードされた機能とともに、R *data.frame* プロキシ・オブジェクトを通じてデータベース表やデータベース・ビューを操作します。これにより、データをデータベース内に保持したままで、使い慣れたR構文を使用してデータベース・データを操作できるため、より迅速かつ大規模にデータの探索、準備、分析を行うことが可能になります。



また、OML4Rでは自然なRインタフェースが提供されるため、インデータベース機械学習モデルを構築でき、それらのモデルをデータのスコアリングで使用できます。機械学習技術には、分類、リグレーション、クラスタリング、属性評価、異常検出、相関ルール、時系列、特徴抽出などが含まれます。インデータベース・アルゴリズムは、アルゴリズム別の自動データ準備、パーティション化されたモデル一式、および統合テキスト・マイニングをサポートします。Oracle Machine Learning for Rを使用すると、より多くのデータに対してさらに多くのモデルを構築し、大量データのスコアリングを素早く実行できます。データ・サイエンティストが生産性向上の恩恵を受けると同時に、エキスパートでないユーザーも強力なインデータベース・アルゴリズムを使用しやすくなります。

Oracle Autonomous Databaseでは、SQLやPythonと並んで、Oracle Machine Learning Notebooks経由でRを使用し、指定されたタスクに最適な言語を選択できます。OML4Rクライアント・パッケージを使用すると、RエンジンからオンプレミスのOracle DatabaseやOracle Database Cloud Serviceデータベースに接続できます。サードパーティ・クライアントからOracle Autonomous DatabaseへのOML4Rアクセスは、今後のリリースで導入されます。

埋込みR実行を使用することで、Rユーザーはユーザー定義のR関数を開発でき（サードパーティ・パッケージを含む）、それらを本番デプロイメント用のデータベース・スクリプト・リポジトリで管理できます。その後、これらのユーザー定義関数は、Oracle Autonomous Database上のSQLインタフェースおよびRESTエンドポイント経由、およびOracle Database上のSQLインタフェース経由で呼び出すことができます。たとえば、ユーザー定義のR関数をデータパラレルおよびタスクパラレルの方法で実行して、データベースによって起動および管理される複数のRエンジンを利用することで、ネイティブRモデルの大規模なスコアリングを実行することもできます。これらのユーザー定義関数の結果には、構造化された結果とイメージ結果の両方が含まれる可能性があります。

Rオブジェクトは、フラット・ファイル内に管理する代わりに、データ・ストア機能を使用してデータベース内に保存することもできます。データ・ストアは埋込みR実行と連携して、ユーザー定義関数に1つまたは複数のRオブジェクト（ネイティブR機械学習モデルなどの複雑なオブジェクトを含む）を渡します。

スクリプト・リポジトリ、データ・ストア、および埋込み実行機能により、データ・サイエンス・チーム内の共同作業が容易になるため、データ・サイエンティストからアプリケーション開発者へデータ・サイエンス作業成果を手軽に引き継ぐことができ、即時のデプロイメントが可能になります。

透過レイヤー

OML4Rの透過レイヤーにより、インデータベース・パフォーマンスを利用するための基盤が提供されます。プロキシ・オブジェクトとオーバーロードされたR関数（透過的にSQLに変換される）を使用して、データベース表の列索引、問合せ最適化、並列処理、インメモリ・キャッシング、表レベルのパーティション化を利用できるため、データベース表を操作する際のパフォーマンスを大幅に向上させることが可能です。

ここで、`ore.create`関数を使用してR `data.frame`（サンプルでは`df`）から表を作成し、使用したい表の名前（この場合は'`BOSTON`'）を指定する方法を見てみましょう。

```
ore.create(df, table = 'BOSTON')
ore.sync(table = 'BOSTON')
```

```
dim(BOSTON)
head(BOSTON)
summary(BOSTON)
min(BOSTON$age)
split(BOSTON, BOSTON$chas)
```

すでに表がデータベース内に存在する場合は、`ore.sync`関数を使用して、プロキシ・オブジェクトを取得したい対象の表またはビューの名前（もしくはSQL問合せ）を指定できます。このプロキシ・オブジェクト（サンプルでは`BOSTON`と命名）はその後、`dplyr`パッケージからの関数を含む、`dim`、

head、*summary*、*split*などの使い慣れたR関数を呼び出すために使えるようになります。これらのオーバーロードされた関数により、インデータベース処理用に機能がSQLに変換されます。

インデータベースMLアルゴリズム

OML4Rは、Rの計算式仕様を使用して、自然なR APIをインデータベース・アルゴリズムに提供します。OML4Rは、以下のように、分類、リグレーション、クラスタリング、属性評価、異常検出、相関ルール、時系列、および特徴抽出のためのインデータベース・アルゴリズムをサポートします。インデータベース・アルゴリズムは、アルゴリズム別の自動データ準備、パーティション化されたモデル一式、および統合テキスト・マイニングをサポートします。

R APIを介して生成されたインデータベース・モデルは、SQL APIを使用したアクセス、個別管理のためにフラット・ファイルとしてエクスポート、他のOracle DatabaseインスタンスやAutonomous Databaseインスタンスへのインポート、およびOracle Machine Learning Servicesへのデプロイが可能です。

分類 <ul style="list-style-type: none">• デシジョン・ツリー• ロジスティック回帰 (GLM)• Naive Bayes• ニューラル・ネットワーク• サポート・ベクター・マシン (SVM)• ランダム・フォレスト• XGBoost (21c)	クラスタリング <ul style="list-style-type: none">• 階層型 k-means• 直行パーティショニング• 期待値最大化• EMを使用した混合ガウス・モデル	マーケット・バスケット分析 <ul style="list-style-type: none">• Apriori - 相関ルール
回帰 <ul style="list-style-type: none">• 一般化線形モデル (GLM)• ニューラル・ネットワーク• サポート・ベクター・マシン (SVM)• XGBoost (21c)	属性評価 <ul style="list-style-type: none">• 最小記述長• ランダム・フォレスト	特徴抽出 <ul style="list-style-type: none">• 非負値行列因子分解• 主成分分析• 特異値分解• 明示的セマンティック分析
	異常検出 <ul style="list-style-type: none">• 1クラスのサポート・ベクター・マシン	時系列 <ul style="list-style-type: none">• 単純指数平滑法• 二重指数平滑法• 三重指数平滑法

埋込みR実行

埋込みRの実行により、データベース環境によって起動および管理されるRエンジンで、ユーザー定義のR関数を実行できます。この機能により、ユーザー定義のR関数内にデータベース・データを自動的にロードできます。関数の結果には、イメージと構造化コンテンツの両方が含まれる可能性があり、構造化コンテンツは`data.frame`プロキシ・オブジェクト経由でアクセス可能な表として返されます。ユーザー定義のR関数は、データベース・スクリプト・リポジトリ内に保管して管理できます。

Oracle Autonomous Databaseでは、サード・パーティ・パッケージをユーザー定義関数内で使用して、モデルの構築、データのスコアリング、視覚化の生成などを実行できます。管理者は、conda環境を作成してOracle Object Storageにアップロードします。そこから、非管理者ユーザーは、Oracle Machine Learning Notebooksで直接使用したり、埋込み実行のSQLまたはREST呼出しで使用したりするために、特定のconda環境をダウンロードしてアクティブ化できます。Oracle Databaseでは、適切なデータベース・マシン・アクセス権限を持つ

ユーザーは、サードパーティ・パッケージをカスタム・インストールでき、ローカルのR環境から、または埋込み実行のSQL呼出しでそれらを使用できます。

また、埋込みRの実行インフラストラクチャでは、データ・パラレルおよびタスク・パラレルで関数を実行できるため、ネイティブのR機械学習モデルを使用したデータ・スコアリングなどの“Embarrassingly Parallel（自明な並列性）”のユースケースにも対応できます。

埋込みR実行により、ユーザー定義のR関数の本番デプロイが容易になるため、SQLインタフェース（Oracle Autonomous Databaseの場合はRESTインタフェース）から関数を呼び出してアプリケーション統合を実現できます。

- **doEval**はもっとも単純で、関数を単一のRエンジンで実行して、`data.frame`や、Rグラフィック・エンジンから生成されるイメージを含む関数の結果を返します。
- **tableApply**も単一のRエンジンで関数を実行しますが、プロキシ・オブジェクトによってR `data.frame`として参照されるデータを関数の最初の引数として渡します。
- **rowRply**は、データを行チャンクにパーティション化することで並列処理を可能にします。データベース環境によって起動および制御される1つまたは複数のRエンジンを使用し、各チャンクに対してユーザー定義のR関数を実行します。
- **groupApply**もデータの並列処理を可能にします。データベース・データは、`index`パラメータにより指定された列でパーティション化されます。`parallel`引数の指定に従って1つまたは複数のRエンジンを使用し、各パーティションに対してユーザー定義のR関数を実行します。
- **indexApply**は、タスクの並列処理を可能にするために、ユーザー定義のR関数を指定された回数実行し、`index`の値を最初の引数として関数に渡し、`parallel`引数の指定に従って1つまたは複数のRエンジンを使用します。

以下のサンプルでは、R APIを使用してOracle Databaseでユーザー定義のR関数を呼び出しています。この関数は、`iris`データセットを使用してSpeciesを予測するランダム・フォレスト分類モデルを構築し、“`RF-model-iris`”というデータ・ストアにモデルを保管します。Oracle Databaseでは、`randomForest`パッケージがおのおののデータベース・サーバー・マシン・ノードのRエンジンにインストールされます。

```
buildRFmodel <- function(dat,dsname) {
  library(randomForest)
  dat$Species <- as.factor(dat$Species)
  mod <- randomForest(Species ~ ., data=dat)
  ore.save(mod,name=dsname,overwrite = TRUE)
  TRUE
}
ore.connect(...database credentials...)
ore.scriptCreate('buildRFmodel', buildRFmodel)

ore.sync(table='IRIS') # ore.frameプロキシ・オブジェクトを取得
```

```
ore.tableApply(IRIS, FUN.NAME='buildRFmodel',
               dsname= 'RF-model-iris',
               ore.connect=TRUE)
```

関数を定義したら、*ore.connect*を使用してデータベースに接続し、ユーザー定義のR関数をデータベース・スクリプト・リポジトリに保存します。

*Autonomous Database*では、データベースは自動で接続されるため、明示的な接続手順がないことに注意してください。

次に、データベースにある*IRIS*表へのプロキシ・オブジェクトを取得します。これには、R変数*IRIS*を通じてアクセスできます。続いて、*tableApply*関数を使用します。この関数は、プロキシ・オブジェクト、スクリプト・リポジトリに保管された関数の名前、モデルを保管するデータ・ストア・インスタンスに使用するための名前、そしてデータ・ストア機能を使用しているために関数内でのデータベースへの接続を自動的に確立する状態を取ります。

以下のサンプルでは、REST APIとOracle Autonomous Database上の対応するSQL APIを使用して、ユーザー定義関数*buildRFmodel*を呼び出しています。

```
$ curl -i -X POST --header "Authorization:Bearer ${token}"
  --header 'Content-Type: application/json'
  --header 'Accept: application/json' \
  -d '{"input": "IRIS", "parameters": {"dsname": "RF-model-iris"}}' \
  "${omlserver}/oml/api/r-scripts/v1/table-apply/buildRFmodel"
```

```
SELECT * FROM table(rqTableEval(
  inp_nam => 'IRIS',
  par_lst => '{"dsname": "RF-model-iris"}',
  out_fmt => 'XML',
  scr_name => 'buildRFmodel');
```


SQL インタフェースを使用して、スクリプト・リポジトリにユーザー定義のR関数を保管することもできます。これにより、データサイエンティストとR開発者から本番環境の管理者への引継ぎが手軽に行えます。たとえば、以下のサンプルでは、`sys.rqScriptCreate`関数を使用して対応するスコアリング関数を保管しています。

```
BEGIN
  sys.rqScriptCreate('scoreDataRF',
    'function(dat, dsname) {
      library(randomForest)
      ore.load(name=dsname)
      dat$Prediction <- predict(mod, newdata = dat)
      dat[,c("Species","Prediction")]
    },FALSE, TRUE); -- not global and enable overwrite
END;
```

スコアリングは本質的に自明な並列性のユースケースであるため、`ore.rowApply`を使用してスコアリング関数を呼び出して、一度に処理する行数（サンプルでは10）と使用するRエンジンの数（3）を指定できます。また、`FUN.VALUE`パラメータに、返す結果の構造体を指定します。このサンプルでは、結果として単一のプロキシ・オブジェクト`data.frame`を返します。

```
PRED = ore.rowApply(IRIS,
  FUN.NAME = 'scoreDataRF',
  rows = 10,
  parallel = 3,
  dsname = 'RF-model-iris',
  ore.connect = TRUE,
  FUN.VALUE = data.frame(Species=character(),
                        Prediction=character(),
                        stringsAsFactors=FALSE))

class(PRED)      # ore.frameプロキシ・オブジェクトを返す
ore.create(PRED,table = 'BATCH_SCORES') # persist as table
with(BATCH_SCORES, table(Species, Prediction))
```

これは、SQLおよびREST APIを使用して呼び出すこともできます。

データベース処理

Oracle Machine Learning for Rは、データベース内に存在するデータ上でのRスクリプトの実行において、データサイエンティストやより一般的なRユーザーをサポートします。これにより、使い慣れたR関数の使用をサポートしながら、クライアントへのデータ移動を回避します。

OML4Rは、選択されたR関数をオーバーロードし、これらの関数を透過的にSQLへ変換します。これにより、データ移動がなくなると同時に、セキュリティが維持され、データベース表の列索引、問合せ最適化、並列処理、それにストレージレベルのパーティション化も利用できます。Oracle Machine Learning for Rは、並列処理、スケーラビリティ、およびメモリ最適化のために再設計

された広範囲にわたるインデータベース機械学習の技術とアルゴリズムをサポートし、マルチノード・クラスタや、Oracle Exadataなどの最適化されたハードウェアを利用します。

データサイエンティストとRユーザーが、インタラクティブにデータを分析してRスクリプトを開発する間、データと結果をデータベース内に保持できるので、データ移動による待機時間が短縮または解消され、さらにスケーラブルなデータ処理が実現します。ファイル抽出からデータをロードすることに慣れている場合は、OML4Rによって、データベース中心のアーキテクチャへの移行が容易になります。このようなアーキテクチャにより、コストと複雑さが軽減される一方で、生産性と信頼性が向上します。さらに、データをデスクトップからデータベースへ簡単に移動できます。

OML4Rは、選択されたオペレーティング・システムのOracle Autonomous Databaseバージョン19cおよび21c、およびOracle Databaseバージョン18c以降でサポートされます。OML4Rは、Oracle Autonomous Databaseサブスクリプション、Oracle Databaseライセンス、およびOracle Database Cloud Serviceと同梱されています。

結論

Oracle Autonomous DatabaseまたはOracle Databaseによって管理されるデータ、もしくはそれらにアクセスできるデータがある場合は、データベースのパフォーマンスとスケーラビリティを利用して、データサイエンティストや他のRユーザーによるデータの探索、準備、分析が可能になります。さらには、インデータベース・アルゴリズムでネイティブR APIを使用して、大規模に機械学習モデルを構築することもできます。ユーザー定義のR関数をデプロイしてシステムが提供する機能を利用し、REnジンを開始および制御することで、データをロードし、データの並列処理とタスクの並列処理を可能にします。これらはすべて、サードパーティ・パッケージを使用して実行できます。

Oracle Cloud Free TierでAlways Freeクラウド・サービスを使用して、Oracle Autonomous Database上でOracle Machine Learning for Rを無料でお試ください。今すぐ[こちら](#)でアカウントを作成してください。

また、Oracle Databaseで使用するためのOML4Rを[こちらから](#)ダウンロードできます。

参考資料

OML4Rのドキュメント

<https://docs.oracle.com/en/database/oracle/machine-learning/oml4r>

Oracle Machine Learning

データ・サイエンス・プロジェクトをサポートする製品ファミリで、機械学習、統計分析、データ準備、視覚化を提供する、SQL/Python/R用のAPIと次のコンポーネントで構成されています。

Oracle Data Minerは、分析ワークフローを作成するためのドラッグ・アンド・ドロップ式ユーザー・インタフェースです。Oracle Machine Learning Servicesは、RESTベースのモデル管理とデプロイ、データとモデルの監視を提供します。Oracle Machine Learning AutoMLユーザー・インタフェースでは、ノーコードで機械学習モデルを作成できます。

<https://oracle.com/data-science/machine-learning>

ROracle

ROracleは、DBIに基づくOracle Database接続用のRパッケージです。Oracle Machine Learning for Rでは、データベース接続時と埋込みR実行関数内でROracleを使用します。ROracleは、データベース表での挿入、更新、および削除処理にも使用できます。

<https://cran.r-project.org/web/packages/ROracle/index.html>

<https://www.oracle.com/database/technologies/roracle-downloads.html>

Connect with us

+1.800.ORACLE1までご連絡いただくか、**oracle.com**をご覧ください。北米以外の地域では、**oracle.com/contact**で最寄りの営業所をご確認いただけます。

 blogs.oracle.com

 facebook.com/oracle

 twitter.com/oracle

Copyright © 2023, Oracle and/or its affiliates. All rights reserved. 本文書は情報提供のみを目的として提供されており、ここに記載されている内容は予告なく変更されることがあります。本文書は、その内容に誤りがないことを保証するものではなく、また、口頭による明示的保証や法律による黙示的保証を含め、商品性ないし特定目的適合性に関する黙示的保証および条件などのいかなる保証および条件も提供するものではありません。オラクルは本文書に関するいかなる法的責任も明確に否認し、本文書によって直接的または間接的に確立される契約義務はないものとします。本文書はオラクルの書面による許可を前もって得ることなく、いかなる目的のためにも、電子または印刷を含むいかなる形式や手段によっても再作成または送信することはできません。

本デバイスは、連邦通信委員会のルールに基づいた認可を未取得です。認可を受けるまでは、このデバイスの販売またはリースを提案することも、このデバイスを販売またはリースすることもありません。

OracleおよびJavaはOracleおよびその子会社、関連会社の登録商標です。その他の名称はそれぞれの会社の商標です。

IntelおよびIntel XeonはIntel Corporationの商標または登録商標です。すべてのSPARC商標はライセンスに基づいて使用されるSPARC International, Inc.の商標または登録商標です。AMD、Opteron、AMDロゴおよびAMD Opteronロゴは、Advanced Micro Devicesの商標または登録商標です。UNIXは、The Open Groupの登録商標です。0120

免責事項：データ・シートにこの免責事項の記載が必要かどうか分からない場合は、収益認識方針を参照してください。本書の内容と免責事項の要件についてさらに質問がある場合は、REVREC_US@oracle.com宛てに電子メールでご連絡ください。