

AI/機械学習プラットフォーム DataRobotの 使用方法

How to Use the AI/Machine Learning Platform DataRobot

伊地知晋平 DataRobot Inc. リージョナルディレクター、データサイエンス Shinpei Ijichi

」 緒言

従来、品質管理分野では、プロセスから得られたデータを 使って統計解析によりプロセスパラメータとアウトカムとの 間の関係性をモデル化し、品質向上や最適化を目的としたモ デルの活用が一般的に行われてきた。近年のIoTおよびコン ピュータ技術の発展の恩恵により、プロセスからますます多 種大量のデータを安価に獲得できるようになったため、製造 業の現場でもモデル化の手段として古典的な統計解析手法だ けでなく機械学習が用いられるようになってきており、日本 の産業界においてその流れは今後ますます加速するものと思 われる。

一方で、統計解析や機械学習を使いこなして成果を上げる ためには、Problem – Plan – Data – Analysis – Conclusionの 頭文字をとった PPDAC サイクル¹⁾ (図1)を何度も回す必要 がある。

この「何度も回す」が重要で、一般的にはデータを一回分



図1 PPDACサイクル

析すればそれでビジネス目標を達成できるだけの成果となる ケースは少なく、仮説立案~仮説検証を何度も繰り返しなが ら、プロセスパラメータとアウトカム間の関係性を帰納的に 詰めていく泥臭い作業が必要となることが多い。

しかしながら、統計学やデータサイエンスの非専門家に とって、RやPythonなどのプログラミング言語を用いてコー ディングを行いながら、何度もPPDACサイクルを回す作業 は非常に負荷の重いものであるし、本来彼らの専門領域での ドメイン知識を駆使して仮説を立案したり分析結果を解釈し たりしなければいけないのに、そのような本質的な業務に十 分な時間を割けない矛盾を感じることが少なくない(筆者は 過去製造業系企業でのデータ分析業務に従事していたが、一 つのモデルを実務で使えるレベルに持っていくために数週間 かかることも稀ではなかった)。

この問題に対する有力なソリューションが機械学習プロ セスの自動化技術である。筆者が所属するDataRobot Inc. は、2012年米国ボストンに設立以来機械学習自動化技術に特 化したソリューション開発に取り組み、2015年にこの分野 では世界で最も早くエンタープライズAIプラットフォーム [DataRobot Automated Machine Learning] を世に出した。 以降エンタープライズAIのグローバルリーダーとして「企 業のAI化を実現」を旗印に精力的な新技術開発に努め、現 在は全世界数百名のデータサイエンティストやソフトウェ ア技術者が日夜技術開発に携わっている。またAIプラット フォームをご提供するだけではなく、2019年からはお客様を AIで成功に導くためのカスタムサービスプログラム「AIサ クセス[®] | を通して、日本でも約2年間で数百を超えるAI開 発ユースケース (事例) に携わってきた経験・知見に基づく ノウハウを惜しみなくユーザー様に提供し、お客様の成功を ご支援させていただいている。



本稿では、「DataRobot Automated Machine Learning」(以 下DataRobot)を例に、データの準備からモデルの業務実装 までの一連のプロセスが高度に自動化されたエンタープライ ズAIプラットフォームの使用方法を紹介する。さらに、モデ ルを配備 (デプロイ)し運用を開始してからの監視・管理も 同プラットフォーム上で一元的に行えることを示す。なお、 本稿は2020年3月時点のDataRobot Ver.5.3 (クラウド版)を 元に書かれていることをお断りしておく。また、図2~図22 は全てDataRobotの製品画面よりキャプチャーした。



3.1 機械学習プロセス俯瞰

データ分析フレームワークとして世界的に認知されている PPDACサイクルは機械学習プロセスにも適用できる。

- a) Problem: 機械学習で解くべき問題の発見・定義とビジネスインパクトの試算
- b) Plan: 問題に合ったデータの収集計画を立て、これま でに分かっているドメイン知識を習得
- c) Data: データを準備し、データ形式を整える
- d) Analyze:機械学習を用いたモデリングを行い、同時 にインサイトを得る
- e) Conclusion: 機械学習モデルを解釈し実運用可否を判断する

機械学習プロセスにおいては、上記a)-e)を何回か回 した後、さらに成果物である機械学習モデルを実使用 環境に配備 (デプロイ)し、実運用へ進むまでのステッ プがある。

- f) Deployment: 実運用に進めると判断された機械学習モ デルを開発環境から運用環境に展開する
- g) Operationalize: 運用環境に実装された機械学習モデ
 ルを試験的に運用する
- h) Production & Monitoring: 機械学習モデルを本運用に 切り替え、そのパフォーマンスをトラックする

3.2 DataRobotの守備範囲

前節で俯瞰した機械学習プロセスは、「データを価値に変 換するプロセス」とも言い換えることができる。

2020年3月現在、DataRobotは、c)の一部からh)までの ステップをカバーしている。また2019年12月に発表された Paxata社の買収により、今後データ準備プロセスも高度な省 力化が進み、データが価値に変換されるEnd-to-Endプロセス の自動化・省力化はさらに広範囲に実現される。 DataRobotはデータ投入からベストモデルを作成・配備す るまでには合計数クリックを必要とするだけである。また、 モデルから様々なインサイトを得ることができるため、製造 業や医療業界では要因分析を主目的としてご使用いただいて いるケースもある。

4 DataRobotの使用法

では本章から具体的にDataRobotの使用方法と主要な機能 を説明する。

4.1 データの読み込み~探索的データ解析 (EDA)

図2にDataRobotのメイン画面を示す。DataRobotが推奨 するブラウザはGoogle Chromeである。

DataRobotにデータを投入する方法はいくつかあり、メイン 画面に表示される方法のいずれかを使用してデータをインポー トできるが、最も簡便かつよく使われているのはPC内のローカ ルファイルをドラッグ&ドロップする方法である。データの読み 込み時点で探索的データ解析 (EDA) も同時に行われる (図3)。

データの読み込みが終わると特徴量名がリスト表示され、 任意の特徴量のヒストグラム(数値データ)や棒グラフ(カ テゴリデータ)、基礎統計量などを表示・確認できる(図4)。

モデリングを始めるには、先に予測ターゲットをユーザー が指定する必要があり、ターゲット特徴量の名前(予測する データセットの列)を入力するか、特徴量セットの名前の横 にある「ターゲットとして使う」をクリックする。こうして 画面上部の予測ターゲットを指定する欄に特徴量名が入力さ れると右側の開始ボタンがアクティブになる(図5)。

機械学習を実行する際には読み込んだデータを学習デー タと検定データに分割し、交差検定のフォールド数やホー ルドアウトデータの割合などを設定しなければならないが、 DataRobotはそれらの条件を自動で設定する。そのままの条 件でモデリングを行っても問題ないが、もしユーザーが細か



図2 DataRobotのメイン画面



図3 探索的データ解析



図4 数値データ特徴量の分布と基礎統計量

く調整したい場合には「高度なオプション」メニューから様々 な条件をマニュアルで設定できる(図6)。

4.2 機械学習モデルの自動作成

「開始」ボタンを押すとDataRobotは機械学習モデリング を始め、実行中のモデルの進捗インジケーターが画面右側の ワーカーキューに表示される(図7)。

モデリング結果は (選択した最適化指標に基づいて) 精度 パフォーマンスが最も優れているモデルがリストの上位にラ ンク付けされた状態でモデルのリーダーボード (順位表) に 表示される。

DataRobotは読み込まれたデータの統計的な性質に基づい てブループリントと呼ばれる「前処理&アルゴリズム」のフ ローを概ね30-40種ほど作成し(「オートパイロット」モード の場合)、各ブループリントから高精度な機械学習モデルを自 動作成する。選択されるアルゴリズムにはDataRobot社で独 自に開発したものに加えてPython scikit-learnやRなど様々



図5 予測ターゲットを指定したところ

パーティション スマートダウンサンプリング 時系列	特徴量探索 特徴量制約 その他
パーティション手法を選択: ランダム パーティション特徴量 グループ 日付/周	前刻 層化抽出
行は、各パーティションにわたって同様のターゲット分布が	行われるように割り当てられます。
次の設定でモデルを生成:	
<u>交差検定</u> 教師-検定-ホールドアウト	
交差検定 (CV) の分割数: 2~50の間で定義されているCVフォールドの数。	 CVの分割 ホールドアウト
5	0% 50% 100%
ホールドアウトの割合(%): ホールドアウトセットに割り当てられたデータの割合(0% ~9%の間である必要があります)。 	
20	

図6 「高度なオプション」 画面



なオープンソースの機械学習ライブラリも網羅されている。 なぜ多様なモデルを作成するのかといえば大きく2つの理 由がある。一つは世の中に万能な機械学習アルゴリズムはな く、与えられたデータに対してどのアルゴリズムで作られたモ デルの精度パフォーマンスが最も良いかはやってみないとわ からないという事情による。もう一つは、精度の良い単体のモ デルを複数集めて組み合わせるともっと精度の良いモデル(ア ンサンブルモデルと呼ばれる)になる場合が多いからである。 DataRobotは様々な組み合わせ方法を用いたアンサンブルモデ ルまでも自動作成して、リーダーボード上に表示する。(図8で Blenderという単語が含まれるモデルがアンサンブルモデル)

4.3 機械学習モデルの解釈とインサイト獲得

DataRobotではリーダーボードに表示されたモデル毎に統 一のフォーマットで様々なインサイト(知見)を得られる。 代表的な機能を以下にまとめる。

- a) モデル精度の表示
 - (ア)分類モデル:ROC曲線、混同行列など
 - (イ)回帰モデル:予測値/実測値プロット、残差プロッ トなど
- b)アウトカム(目的変数)と関係性の高い特徴量の表示:
 特徴量のインパクト²⁾(図9)
- c)アウトカムと各特徴量との関係性表示:特徴量ごとの 作用³⁾(図10)
- d)予測結果に強く影響を与えた特徴量とその値を表示:
 予測の説明(図11)
- e) テキスト変数の単語とアウトカムとの関係の強さ: ワードクラウド (図12)

4.4 ベストモデルの選択

前述のように、リーダーボードには精度指標が良い順に多 数のモデルが表示されるので、それらのモデルからユーザー







図9 特徴量のインパクト

にとってベストなモデルを選択する。もしユーザーが精度パフォーマンスだけに関心があるのであれば、精度が最も良い モデルを選択すれば良い。多くの場合、複数の単独モデルを 組み合わせた Blender モデルの精度が最も良くなる。

モデルの実運用形態からリアルタイムに予測を行う必要 があるなどの理由で、予測スコアリング速度も重要な場合に は、予測精度と予測速度の両方が良いパフォーマンスを示す バランスの取れたモデルを選択する。一般的にはBlenderモ デルよりも単独モデルの予測スコアリング速度が速いので、 精度と速度のバランスを重視する場合には単独モデルが選 択されることが多い。DataRobotでは速度と精度のトレード オフをプロットしたチャート(図13)でバランスの取れたモ



図10 特徴量ごとの作用



図11 予測の説明



図12 ワードクラウド(日本語にも対応)

デルを確認できる。チャート中の点の一つ一つが精度(縦軸) 上位のモデルであり、それらの予測時間(横軸)の推定値に 従ってプロットされている(精度上位のモデル10個が右側に リストアップされている)。精度指標はリーダーボードから 任意に指定できる。

リーダーボードではホールドアウトデータに対する精度は デフォルトで非表示となっている。まず交差検定の精度指標 や「速度vs精度チャート」などを確認してユーザーにとって ベストなモデルを選定したら、その後に「ホールドアウトの 解除」を行う。ホールドアウトデータの精度も交差検定の精 度と比べて大きく乖離していなければ、学習時に使われてい ない全く未知のデータに対しての精度パフォーマンスも信用 できるとして、モデルの配備 (デプロイ) に進む。

このとき、ホールドアウトデータの精度を見て改めてモデ ルを選択しなおしてはいけない(モデルの再選択を行なって しまうとホールドアウトデータに過学習したモデルを選ぶ危 険性があるため)。

その後、例えば学習データの行数が小さい場合には、選択 したモデルだけに対して全データを使って再学習させること



図13 予測速度 vs 予測精度



図14 全データを使用してのモデル再学習

も簡単にできる(図14)。このように学習データ数を増やす とさらなる精度の向上が期待できる。

4.5 モデルの配備 (デプロイ) と予測の実施

モデルを運用する形態は予測を行う頻度や1回の予測量、 セキュリティ要件などによって変わるが、DataRobotではベ ストモデルを選択した後に、モデルを実運用環境へ様々な形 態で柔軟に配備できる。本稿ではDataRobot内の実運用環境 にモデルを配備する方法を紹介する(なお、モデルをJAVA スクリプトなどの形式で出力するオプション機能を利用すれ ば、ユーザーのシステムにそのモデルを組み込んで、システ ム内に閉じた形で予測スコアリングを行うことができる)。

方法は非常に簡単で、選ばれたベストモデルの「予測>デ プロイ」タブを選択し(図15)、表示されたオレンジ色のボタ ンを2回クリックするだけで短時間に完了する(図16)。

配備モデルでAPI (Application Programming Interface)を 利用して予測スコアリングを実施するには、DataRobotが自動 作成したPythonスクリプト(図17)を実行すれば良い。実際に スクリプトをユーザーのシステムに組み込んで、定期的に予測 スコアリングを自動で実施させている例もよく見られる。

評価 解釈 説明 予測 コンプライアンス	
予測を作成 デプロイ DataRobot Prime ダウンロード	
	を ナノロ1 観 域に表示しま9。
デプロイに名前を付ける	デプロイを説明
プリードアウト予測	
エンドポイント	
cfds-ccm-prod.orm.datarobot.com	
cfds-ccm-prod orm.datarobot.com ✓ しきい値 予測値に対して設定されたしさい値: 0.5〇 ① データドリフト追跡を有効化 ◎ ② DataRobotによる予測のセグメント分析の実行を許可 ◎	
cfds-ccm-prod orm.datarobot.com	

図15 モデルをデプロイ

ਊDataRobot データ モデル ຜ デプロイ インサイト	Jupyter リポジトリ アプリケーション AIカタ	ログ
< プリードアウト予測 DataRobot予測サーバー コーティング製品ブリードアウト故障_train	サービス ドリフト 精度 in.xlsx 😧 😧 😧	
無要 サービスの正常性 データドリフト 構定 インデグレーショ 概要 このデプロイメントの目的、コンテンツ、履歴を理解します。 サマリー	ョン 設定 コンテンツ	
名前 プリードアウト予測 🥒 エンドポイント https://clds.com.prod.om.datarobot.com 説明 なし 🥒	データセット コーティング 品コブリードアウト 故障_train.xlsx ターベット ブリードアウト モデル Nystroen Kernel SVM Classifier プロジェクト コーティング 製品ブリードアウト 故障_train.xlsx 予測に対するしるい値 0.5	

図16 実運用環境に配備 (デプロイ) されたモデル



4章ではDataRobotでモデルを自動作成し実運用環境へ配備(デプロイ)するまでの方法と機能を紹介した。本章では、 実運用を開始したモデルの監視と管理もDataRobot上で一元 的に行えることを紹介する。

5.1 実運用化されたモデルの管理

DataRobotで作成した機械学習モデルをいくつも業務プロ セスで実運用するようになると、それらのモデルのパフォー マンスを監視して、適切なタイミングでモデルの更新などの 介入を行う必要があるが、それらの監視・管理もデプロイイ ンベントリタブ (図18) から行える。

インベントリの上部には、信号機の色(緑、黄、赤)で色分 けされた正常性インジケーターが表示され、全てのアクティ ブな配備モデルの使用状況とステータスが要約されている。 そのサマリーの下には各配備モデルのステータスが下記項目 について表示され、a)~c)は同じく信号機の色で表示される。

- a)サービスの正常性:過去24時間におけるエラーの有無
- b) データドリフト
- c) 精度: 正解データが与えられると計算される
- d) アクティビティ:過去7日間の予測数パターン
- e)平均予測数/日
- f) 最新の予測:モデルに対して最後の予測を行ってから の時間



図17 自動作成された Python スクリプトの一部

5.2 監視項目の詳細

インベントリの任意の配備モデル名をクリックすると、そ の詳細情報(図16)が表示される。そこで「サービスの正常 性」、「データドリフト」、「精度」のタブを選択するとそのモ デルの過去のステータスがそれぞれ示される。

5.2.1 サービスの正常性

サービスの正常性タブ (図19) では、予測の合計数、予測 実行時間、レスポンス時間、システムエラーの割合、などが ダッシュボード表示される。

5.2.2 データドリフト

データドリフトタブ (図20、図21)は、配備したモデルの 特定の時間間隔における正常性を識別するのに役立つ、3つ のインタラクティブな視覚化機能を提供する。

図20左側の「特徴量ドリフトと特徴量の有用性チャート」 においては、チャート上部のスライダ(水色の横線)で設定 できる任意の期間で、有用性上位の特徴量を「実際に特徴量 値の分布がどれだけ変化したか(縦軸)」と「特徴量の有用性 (横軸)」の2軸に基づきプロットしている(各特徴量が一つ 一つの点)。

ドリフトの大きさによって特徴量を表すポイントの色は信 号機の色で表示される。任意のポイントをクリックすると図 20右側の「特徴量の詳細チャート」に、学習データで選択さ れた特徴量のレコードのスコアリングデータに対するパーセ ンテージ (分布)を表示させて確認することができる(薄い 青色の棒グラフが学習データの分布、濃い青色の棒グラフが 予測スコアリングに使われたデータの分布を示している)。



図18 デプロイインベントリタブ



図19 サービスの正常性



図20 データドリフト



図21 時間経過に伴う予測

図20は左側のチャートで有用性が高くデータドリフト が発生していると表示されたある特徴量のポイント(赤色 になっている)をクリックした際の表示例である。右側の チャートでは薄い色の棒グラフの分布と濃い色の棒グラフの 分布の形状が一部で大きく異なっている。

図21はデータドリフトチャートの下部に表示される「時間 経過に伴う予測」チャートである。縦軸が平均予測値、横軸 が時間を表し、回帰モデルの場合には予測値の中央から±40% の範囲を表す視覚的インジケーターもプロットされる。一番 左側の1ポイントは、モデル学習時のホールドアウトデータ で作成された予測値である。

各ポイントにカーソルを合わせるとそのポイントにおける 期間の開始点-終点、平均予測値、予測回数、10パーセンタイ ル点-90パーセンタイル点の値などを確認できる。もしある 期間の予測平均値が大きくシフトしていた場合には、その期 間に絞って原因を探るなどのアクションを取れる。

5.2.3 精度

精度タブ (図22) は、標準的な統計的手法と可視化によっ て、時間の経過に伴うモデルデプロイのパフォーマンスを表 示する。精度タブを使用するには、DataRobot外で収集され た予測値および実測値を含むデータセットをアップロードす る必要がある。

データドリフトチャートと同様に、精度チャートも上部の スライダ (水色の横線) で過去の任意の期間を設定できる。 上の1本 (緑色) の折れ線グラフは選択した精度指標値 (この 例ではLogLoss)の時間経過に伴う変化を示している。これ

精度 このデブロイの予測の構成を経時評価する モデル 現在 編詞 画次 〇 数 2019-11-20 2016-11-20	1428 Dytryk 2019-1228			2025-03-00
LogLoss Mil:NA 0.355	auc Re:nva 0.678	Kolmogorov-Smirnov Mili: N/A 0.265	Gini Norm RN: N/A 0.357	Rate@Top10% RR: N/A 0,279
時采列の稿度 040 055 025				
予測値&実測値 ○ R200 ● 7200 ■ R2000 ※ R2				
in the second		~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~		

図22 精度

らの指標は、配備 (デプロイ)前にモデルの評価に使用され る指標と同一である。グラフの上にあるメトリックス・タイ ルをクリックすると、表示が変更される。

下の2本 (水色&オレンジ色)の折れ線グラフには、データ セットの全体のタイムラインに沿って予測値と実測値が表示 される。グラフ上のポイントは、データドリフトタブに表示 されるポイントに一致している。したがって、いずれかのプ ロットのポイントの上にマウスを置くと同じ値が表示される (時間スライダが同じように設定されている場合)。

6 結語

以上、エンタープライズAIプラットフォームDataRobot の使用方法と主な機能を紹介した。DataRobotにデータを投 入してからベストモデルを実運用環境へ配備するには、予 測ターゲットの入力、開始ボタンのクリック、モデルデプロ イボタンのクリック (2回)とわずか数回のクリックで済む。 コーディングや機械学習アルゴリズムの知識が乏しい人で も、DataRobotを使えばデータサイエンティストが数週間か けなければできないようなタスクを数時間で行うことができ るため、本稿1章で述べたように様々な仮説を素早く検証し、 本当に実運用に堪えるAIアプリケーションを短期間に開発 することができる。

参考文献

- 総務省統計局,統計データ利活用センター: http://www.stat.go.jp/dstart/point/seminar1/01.html, (accessed 2021-7-28).
- 2) DataRobotブログ, (2019). https://blog.datarobot.com/jp/ permutation-importance
- 3) DataRobotブログ, (2019). https://blog.datarobot.com/ jp/2018/02/15/modelxray

(2020年4月1日受付)