

データの本質を読み解くための機械学習 ～MATLAB でデータ解析の課題に立ち向かう～

MathWorks Japan
アプリケーションエンジニア部
アプリケーションエンジニア
井原 瑞希

Buzzwords...

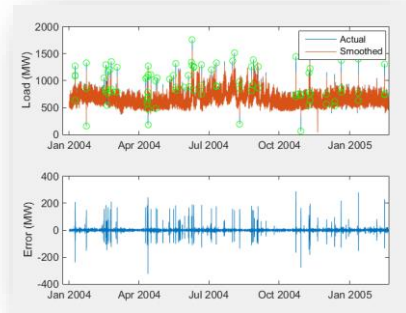
IoT
人工知能 / AI
ビッグデータ

データ解析

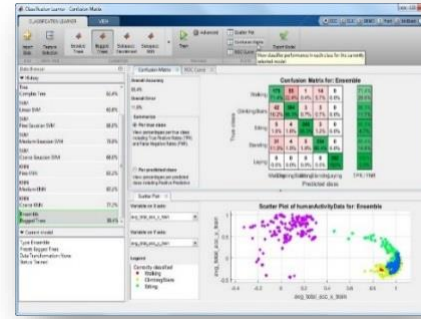
データ解析ワークフロー

	1	2	3	4
	Date	CAPTL	CENTRL	DUNWOD
1	01-Jan-2004 00:00:00	1015	1651	618
2	01-Jan-2004 01:00:00	927	1562	568
3	01-Jan-2004 02:00:00	891	1507	541
4	01-Jan-2004 03:00:00	NaN	1440	517
5	01-Jan-2004 04:00:00	NaN	1434	499
6	01-Jan-2004 05:00:00	NaN	1449	496
7	01-Jan-2004 06:00:00	NaN	1490	524
8	01-Jan-2004 07:00:00	NaN	1525	526
9	01-Jan-2004 08:00:00	960	1529	518
10	01-Jan-2004 09:00:00	1046	1628	541
11	01-Jan-2004 10:00:00	1111	1706	570

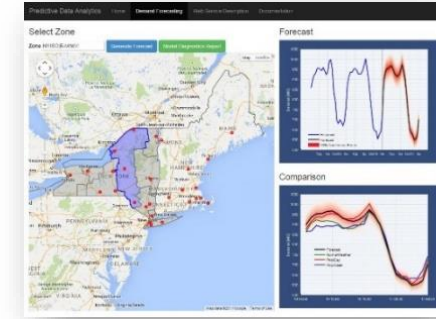
データへの
アクセスと探索



データの前処理

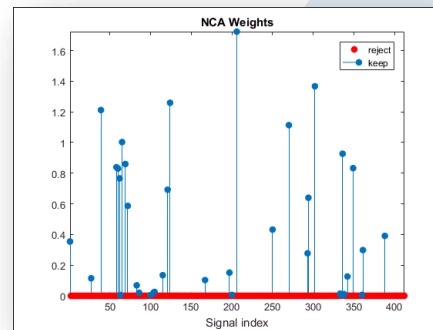


予測モデルの構築

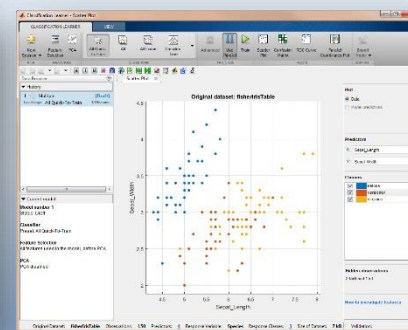


システムへの統合

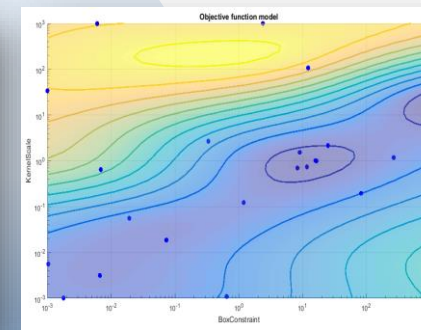
機械学習



特徴選択

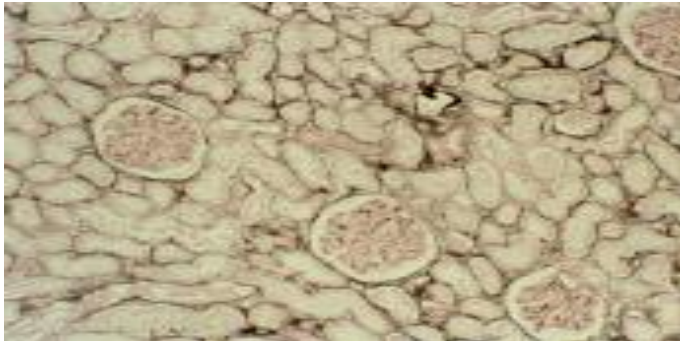
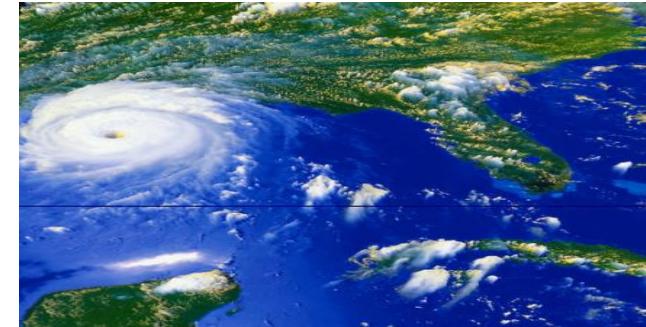
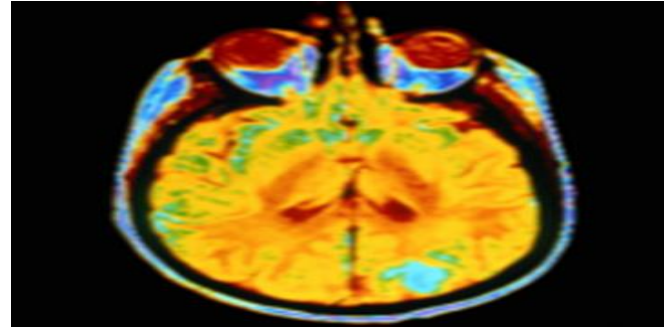


モデルの選択



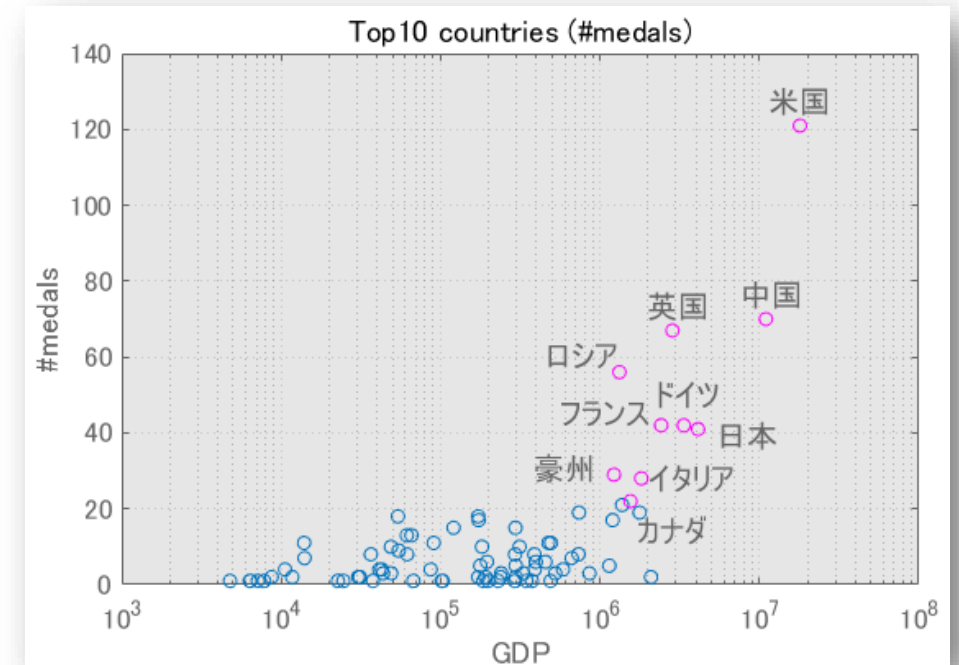
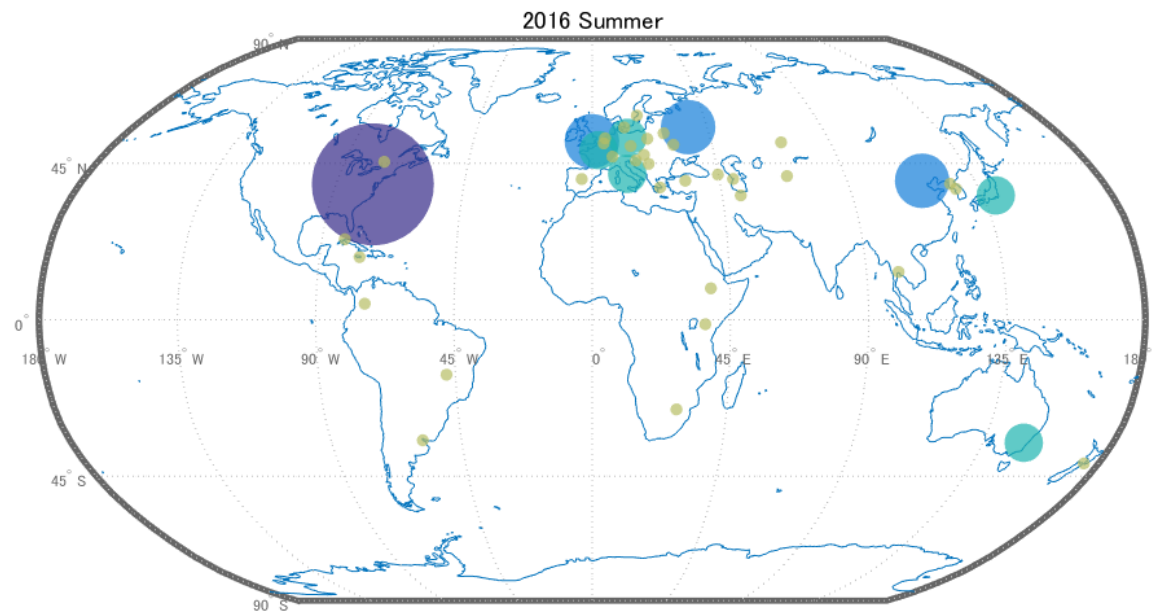
ファイン
チューニング

Machine Learning is Everywhere!



機械学習とは

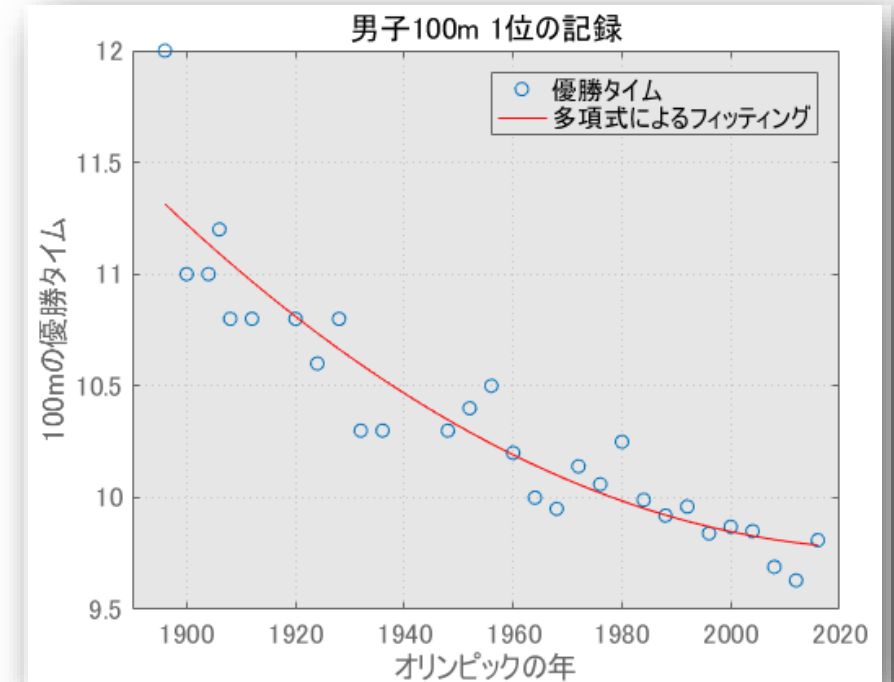
- 機械学習 “ではない” データ解析
 - データそのものから直接わかること



>> SimpleDALive

機械学習とは

- 機械学習の定義
 - データから直接観測できないパターンやルールを、モデルを元にして**機械**的 (自動的) に**学習**すること
- なぜ機械学習を使うのか
 - 未知のサンプルに対する予測が可能
 - 予測に必要な情報を残し、冗長な情報を省くことが可能



1940年、2020年のそれぞれの男子100mの1位の予想タイム

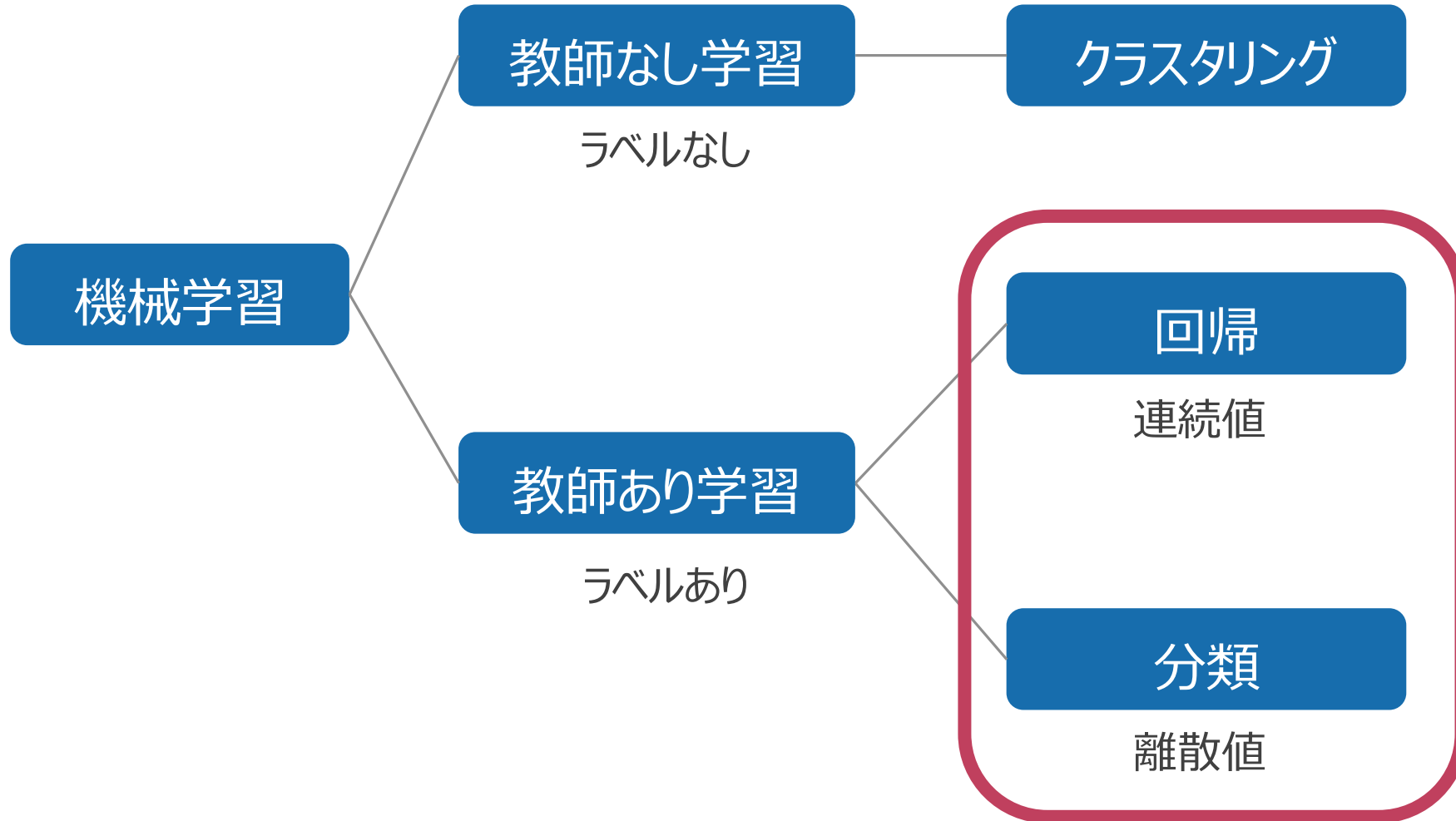
```
fitresult(1940)
```

```
ans = 10.4666
```

```
fitresult(2020)
```

```
ans = 9.7795
```

機械学習とは



本日のトピック

■ 回帰分析

- 回帰分析ワークフロー
- 回帰モデルの見方
- サンプル数が少ないときの回帰分析

■ 分類

- 分類のワークフロー
- コーティングを簡単にする方法
- 機械学習の課題 – ハイパーパラメータの探索

■ 機械学習によるビッグデータ解析

- メモリに収まりきれないビッグデータを扱う場合

回帰

分類

X

ビッグデータ解析

本日のトピック

■ 回帰分析

- 回帰分析ワークフロー
- 回帰モデルの見方
- サンプル数が少ないときの回帰分析

■ 分類

- 分類のワークフロー
- コーティングを簡単にする方法
- 機械学習の課題 – ハイパーパラメータの探索

■ 機械学習によるビッグデータ解析

- メモリに収まりきれないビッグデータを扱う場合

回帰

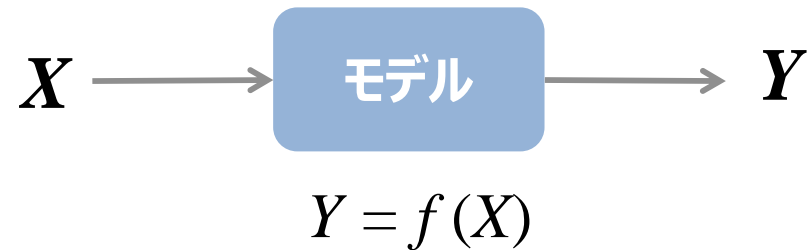
分類

X

ビッグデータ解析

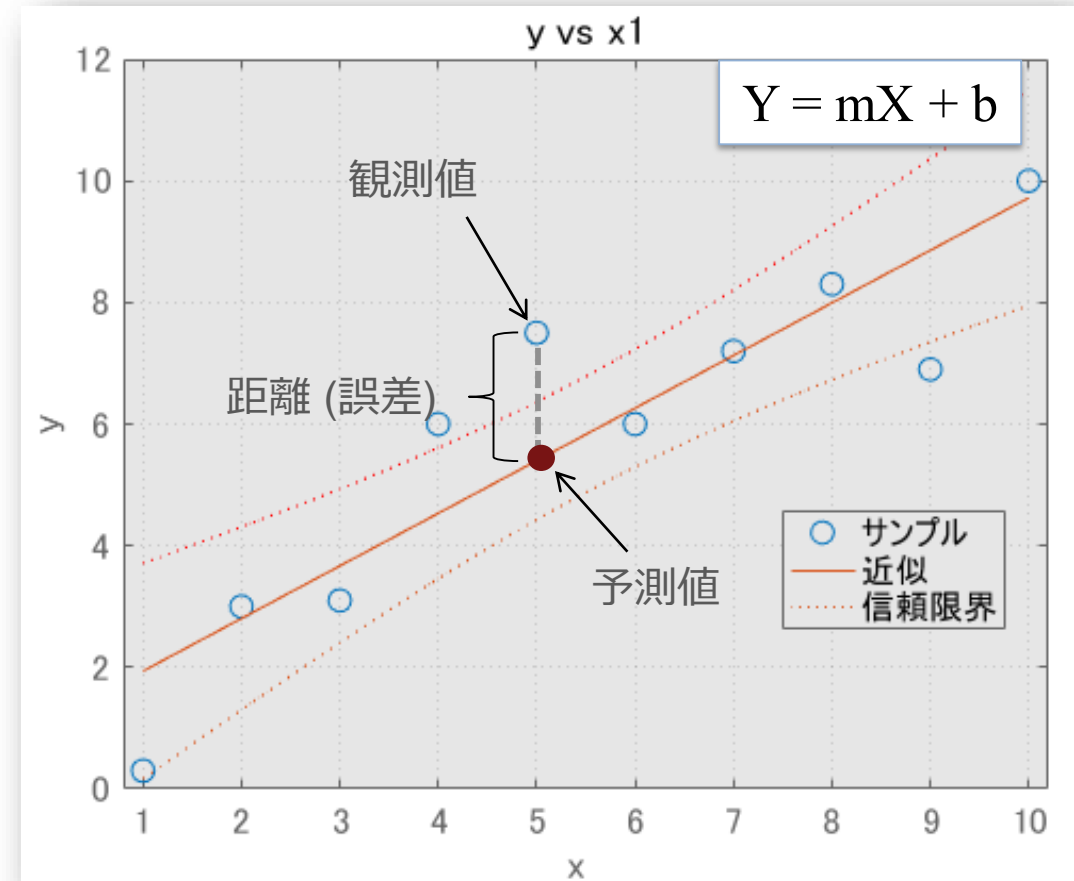
回帰分析とは

- 観測可能なデータから変数間の関係性をモデリング
- 説明変数 X** の関数として **目的変数 Y** を説明する変数の関係性モデルを構築

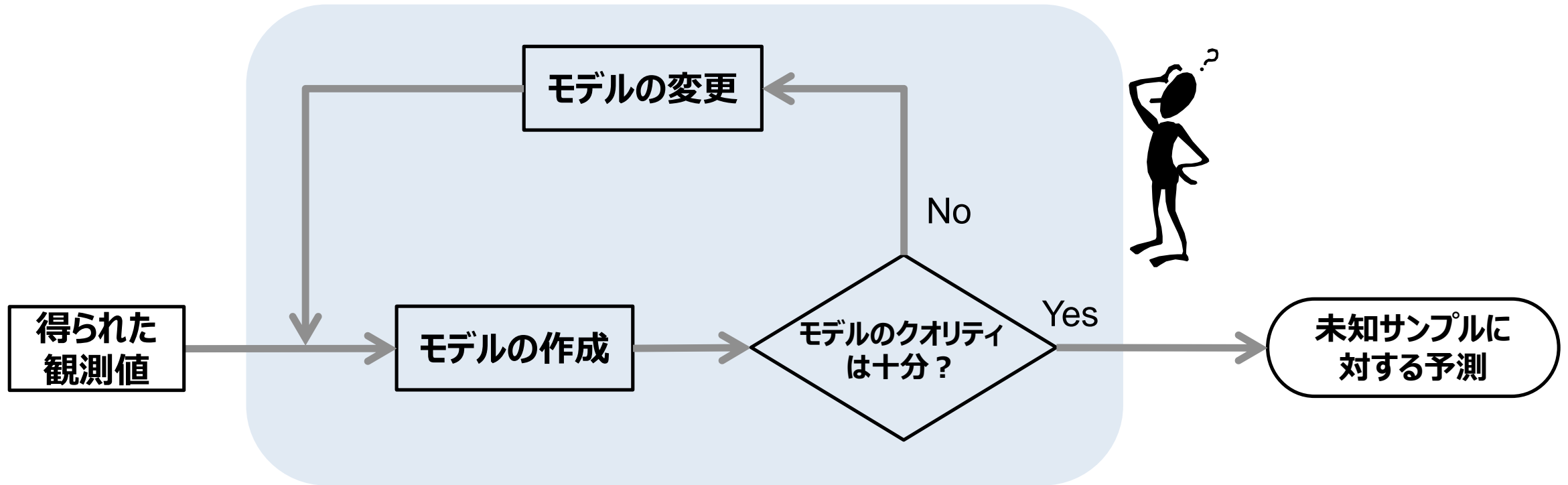


- 観測値と予測値の差を最小化する係数を推定
- 入力 X の未知の出力 Y を予測可能

例: 二乗誤差の和を最小化

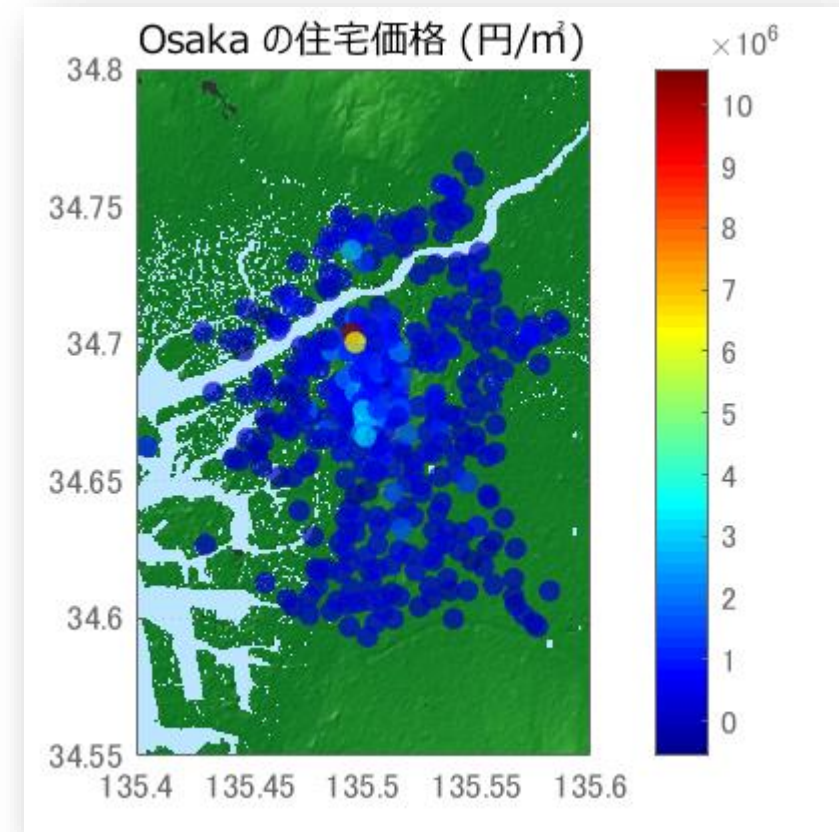


回帰分析の流れ



例: 大阪の住宅価格の推定

- 推定対象 (目的変数)
 - 住宅価格
- 既知の情報 (説明変数)
 - 位置情報
 - 部屋の形状
 - 地積
 - 階数
 - 駅からの距離など
- 目的
 - どの説明変数が住宅価格に関わっているか調べる
 - 正確な住宅価格予測モデルの作成



商業施設の場所の決定、
経済指標として利用

>> HousingPriceEst_Osaka

回帰モデルの推定結果



- 線形回帰モデル:

$$y \sim 1 + x1$$

t値: 説明変数を与える影響

p値: 係数の有意確率
(極端な値を取る確率)

- 推定された係数:

	Estimate	SE	tStat	pValue
(Intercept)	-3749.5	51.154	-73.298	1.188e-21
x1	2.0608	0.025507	80.794	2.5117e-22

観測数: 18、誤差の自由度: 16

二乗平均平方根誤差: 0.561

決定係数: 0.998、自由度調整済み決定係数 0.997

説明変数が、目的変数の変化をどれくらい説明できているか

F 統計量と一定のモデルの比較: 6.53e+03、p 値は 2.51e-22

回帰の種類

■ パラメトリック回帰

例) 線形回帰、ステップワイズ回帰、...

- モデル式を仮定して、データにフィットするようなパラメタを探索

データ = 確定的な成分 + ランダムな誤差

- データの傾向や関数がある程度わかっている場合に有効

■ ノンパラメトリック回帰

- 例) ガウス過程回帰、決定木、...

- 関数の形を定めない

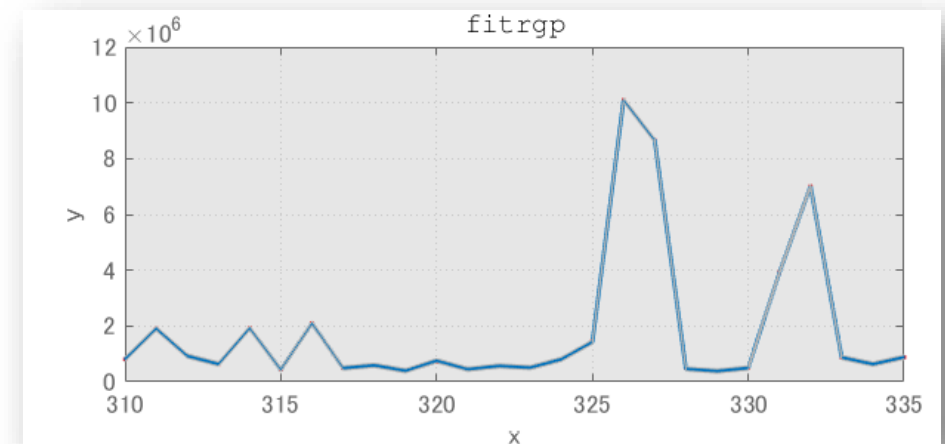
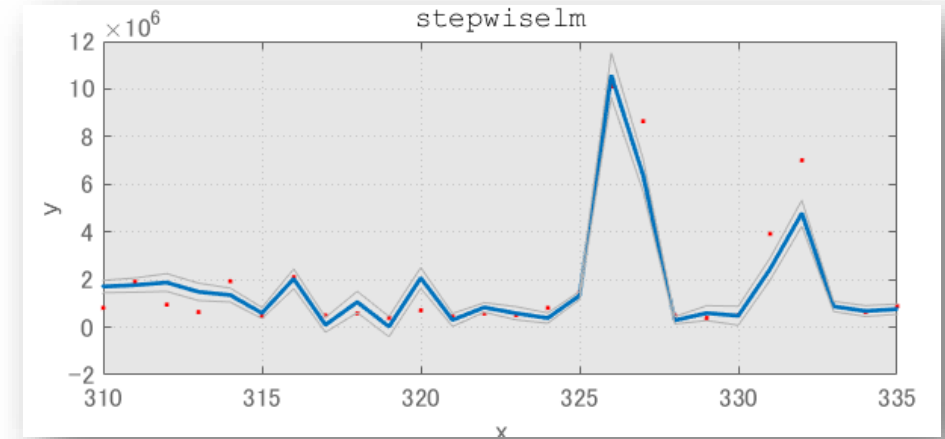
- 事前に関数がわからない場合や**観測可能なサンプル数が少ない場合に有効**

**学習サンプルが少なく
推定結果が良くない場合は？**



ガウス過程回帰

- ガウス過程回帰 (Gaussian Processes, Kriging)
 - ノンパラメトリックな**確率**モデル
 - 訓練データに近ければ分散が小、離れると分散大のガウス分布から確率的に生成されていると仮定 (近傍とのなめらかな遷移を仮定)
- 利点
 - サンプル数が少ない場合にも高い予測精度
 - 途中でサンプルの傾向が変わった場合に対応可能
- 欠点
 - 高次元のデータでは予測精度が高くない



回帰分析のまとめ

モデル		フィッティング関数	クラス名
パラメトリック 回帰	線形回帰	<code>fitlm</code>	<code>LinearModel</code>
	一般化線形回帰	<code>fitglm</code>	<code>GeneralizedLinearModel</code>
	非線形回帰	<code>fitnlm</code>	<code>NonLinearModel</code>
ノンパラメトリック 回帰	サポートベクタ回帰	<code>fitrsvm</code>	<code>RegressionSVM</code>
	ガウス過程回帰	<code>fitrgp</code>	<code>RegressionGP</code>
	回帰木	<code>fitrtree</code>	<code>RegressionTree</code>
	アンサンブル学習 (回帰)	<code>fitensemble</code>	<code>RegressionEnsemble</code>
	ニューラルネットワーク	<code>train</code>	---

本日のトピック

■ 回帰分析

- 回帰分析ワークフロー
- 回帰モデルの見方
- サンプル数が少ないときの回帰分析

■ 分類

- 分類のワークフロー
- コーティングを簡単にする方法
- 機械学習の課題 – ハイパーパラメータの探索

■ 機械学習によるビッグデータ解析

- メモリに収まりきれないビッグデータを扱う場合

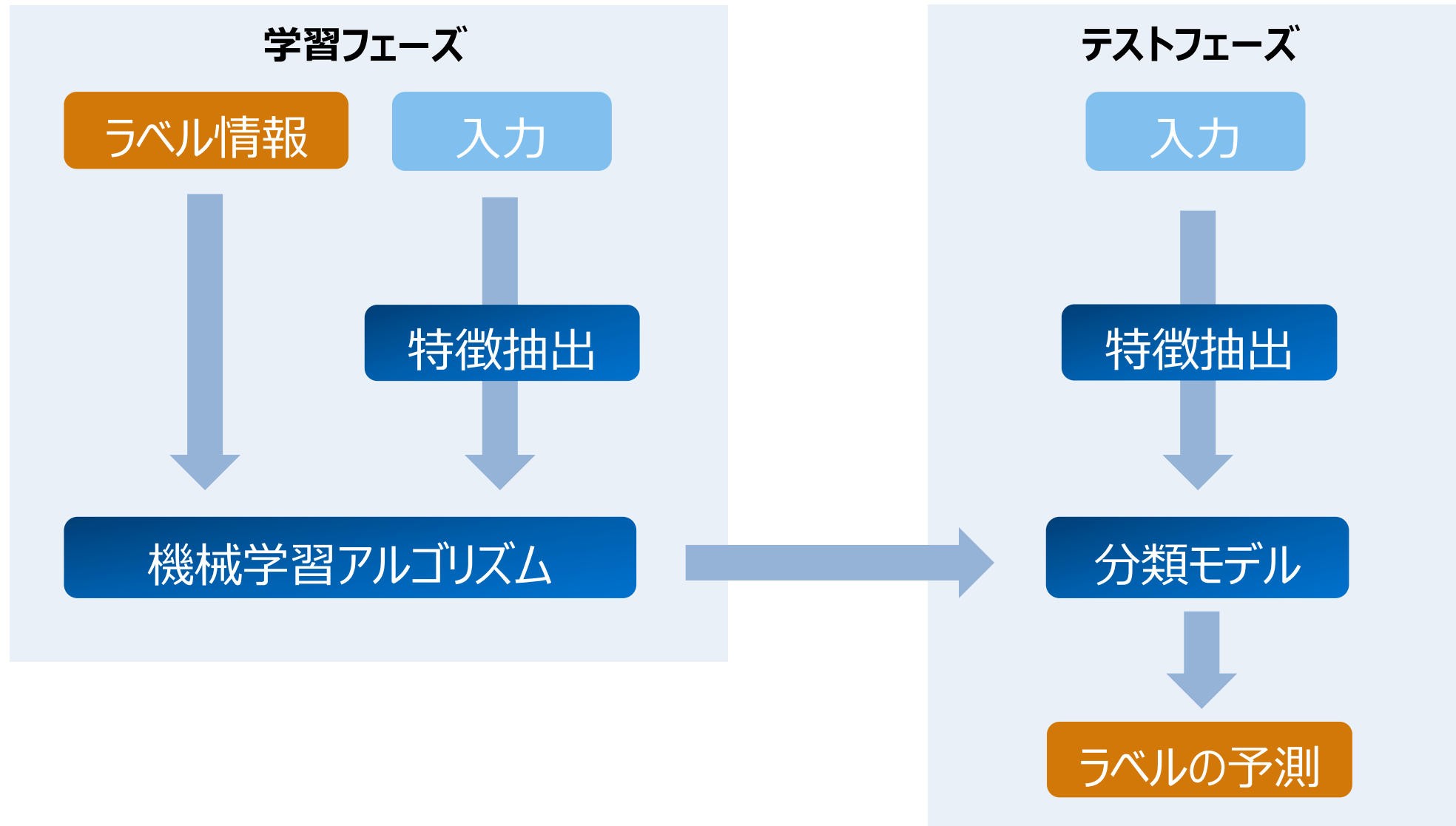
回帰

分類

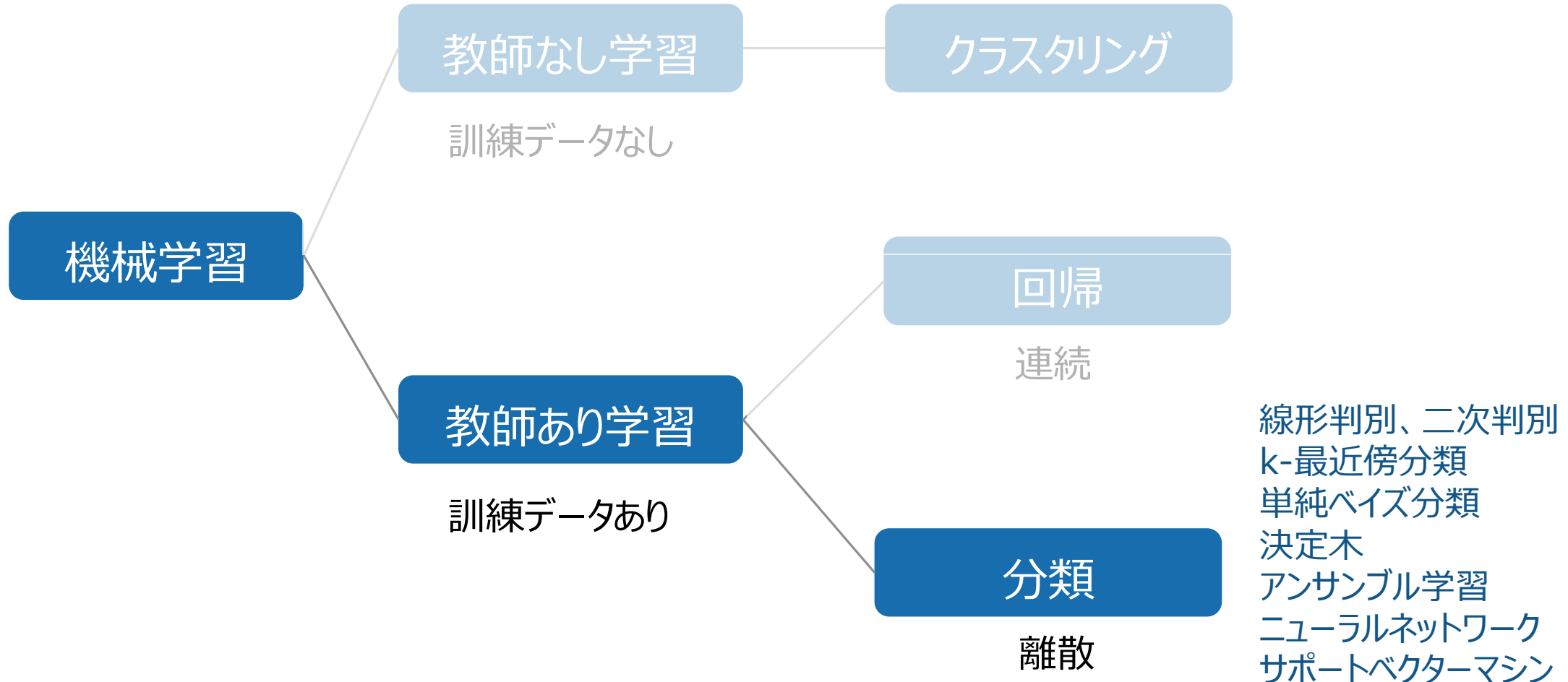
X

ビッグデータ解析

分類の流れ



MATLAB® における機械学習

Statistics and Machine Learning Toolbox
Neural Network Toolbox

アプリを使った機械学習のコード作成

Statistics and Machine Learning Toolbox

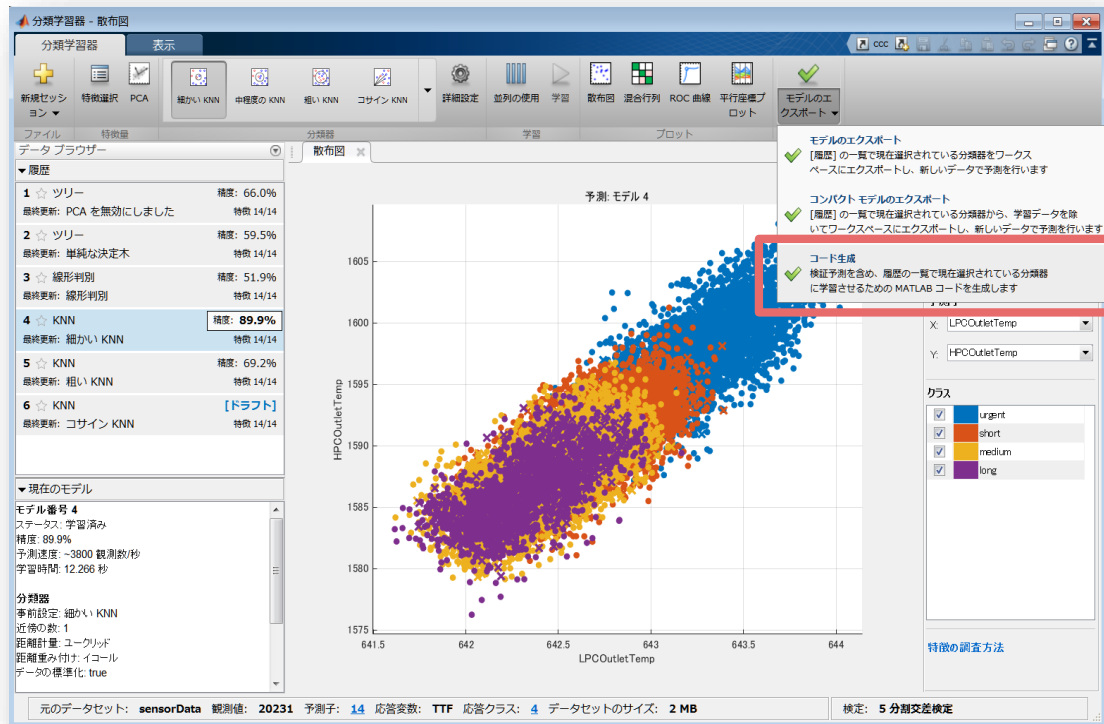
- 分類学習器アプリ
 - データを分類するためのモデル学習 GUI
 - GUI 操作を MATLAB コードとして生成することができる

Parallel Computing Toolbox™
並列モデル学習機能の追加



並列の使用

R2016b



```

52
53 % 分類器の学習
54 % このコードは、すべての分類器オプションを指定し、分類器に:
55 classificationKNN = fitcknn(...
56     predictors, ...
57     response, ...
58     'Distance', 'Euclidean', ...
59     'Exponent', [], ...
60     'NumNeighbors', 1, ...
61     'DistanceWeight', 'Equal', ...
62     'Standardize', true, ...
63     'ClassNames', categorical({'urgent'; 'short'; 'medium'; 'lo
64
65 % 関数 predict で結果の標準化を作成
    
```

MATLAB プログラムの自動生成

例: New York のタクシーチップカテゴリの分類

- 目的
 - チップの多い乗客の傾向を調べる
- 使用するデータ
 - New York のタクシー利用履歴 (乗車の時間、人数、位置など)
 - csv ファイル
- 解析ワークフロー
 - 前処理
 - 解析に使用する特徴の選択
 - チップカテゴリ分類モデルの構築 (**分類学習器アプリ**の使用)
 - チップカテゴリ分類モデルの評価



>> `TaxiTipClassification`

分類器の種類と使い分け

分類器	すべての予測子が数値	すべての予測子がカテゴリカル	一部がカテゴリカル、一部が数値
決定木	あり	あり	あり
判別分析	あり	なし	なし
ロジスティック回帰	あり	あり	あり
SVM	あり	あり	あり
最近傍	ユークリッド距離のみ	ハミング距離のみ	なし
集団	あり	はい。ただし、部分空間判別を除く	はい。ただし、部分空間判別を除く

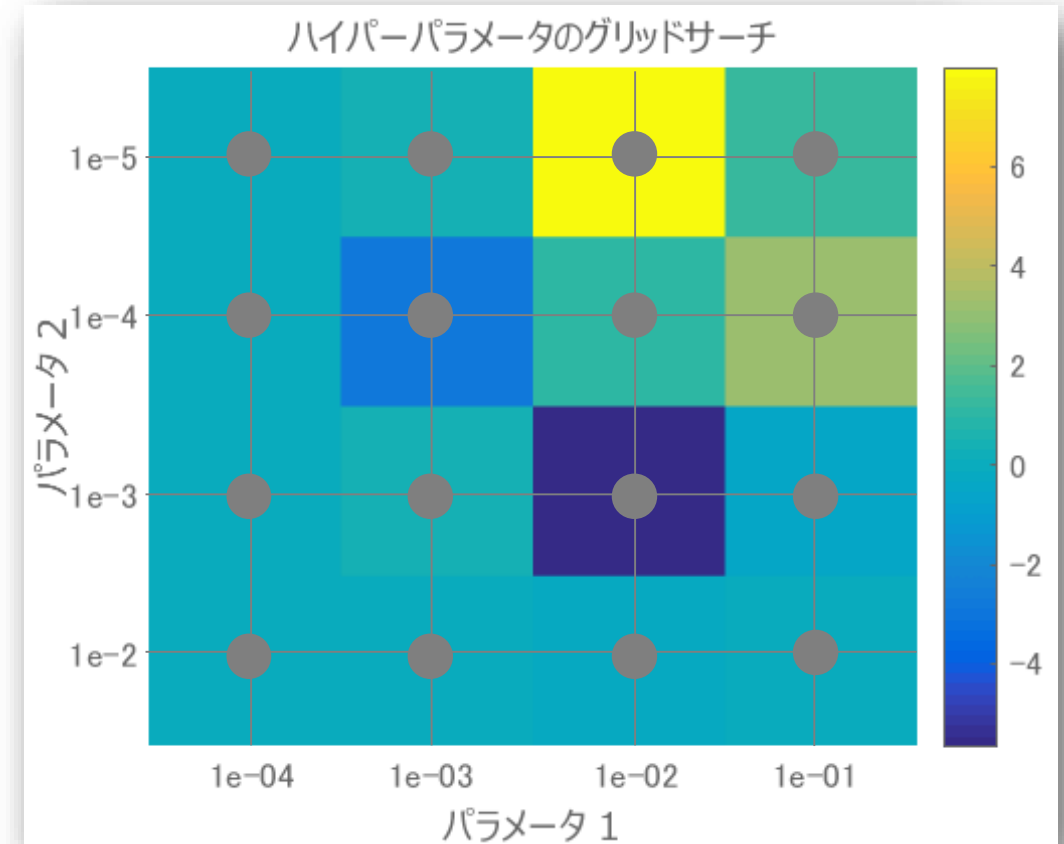
分類器	予測速度	メモリ使用量	解釈可能性
決定木 	高	小	容易
判別分析 	高	線形の場合は小、二次の場合は大	容易
ロジスティック回帰 	高	中	容易
サポートベクターマシン 	線形の場合は中。他の場合は低	線形の場合は中。他のすべて: マルチクラスの場合は中、バイナリの場合は大	線形 SVM の場合は容易。他のすべてのカーネルタイプの場合は困難。
最近傍分類器 	3次元の場合は低。他の場合は中	中	困難
アンサンブル分類器 	アルゴリズムの選択によって高から中	アルゴリズムの選択によって低から高	困難

分類器オプションの選択

<http://jp.mathworks.com/help/stats/choose-a-classifier.html>

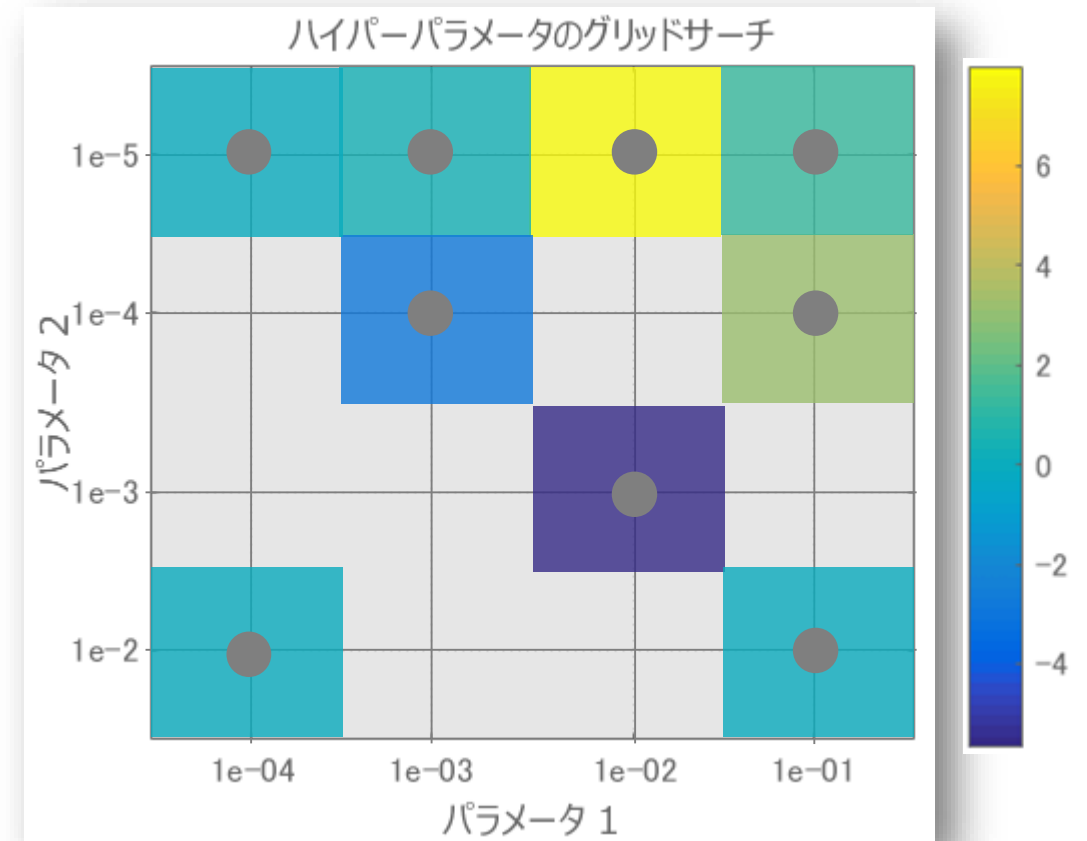
機械学習手法における課題

- 課題
 - ハイパーパラメータの調節
- ハイパーパラメータ
 - データから直接決めることのできないパラメータ
 - ユーザがあらかじめ決めておく必要あり
- **グリッドサーチ** (従来の方法)
 - 格子状の空間で最適なパラメータを探索
 - 課題: ある範囲を総当りするため
計算時間がかかる



機械学習手法における課題: ハイパーパラメータの調節

- **ベイズ最適化** R2016b
 - あるハイパーパラメータでの学習器の精度を目的関数として定義
 - この目的関数を最大化するパラメータを推定
 - ガウス過程回帰でモデル化
 - 精度が上がりやすそうな方向を確率的に推定
 - 条件
 - 低次元データ
 - 目的関数の評価に時間がかかる
 - 低精度
 - 大域的な解を求めたい
 - ハイパーパラメータの決定



ベイズ最適化によるパラメータチューニング

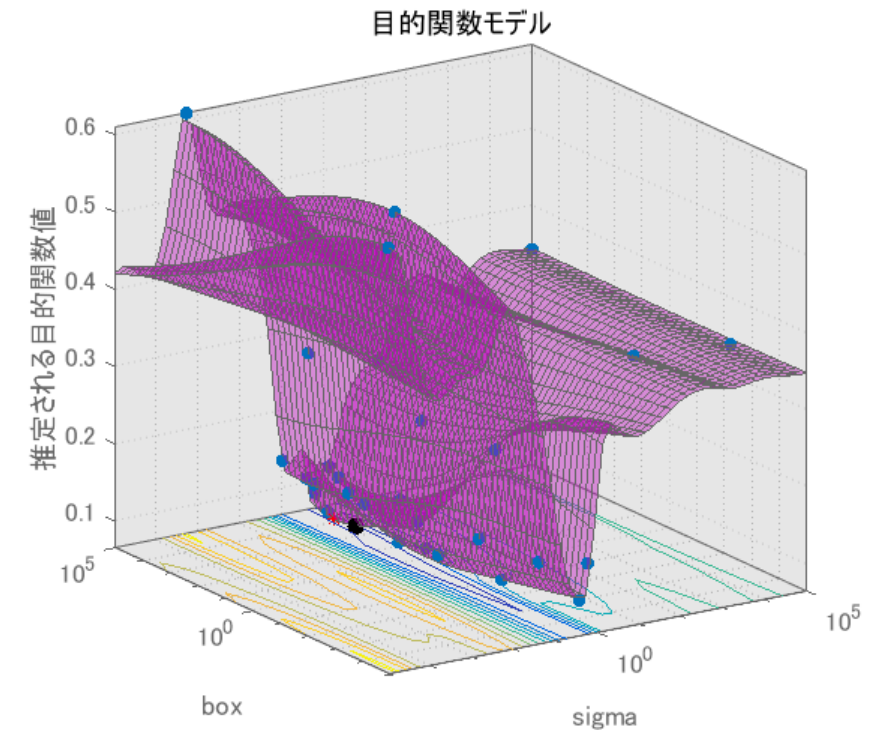
- 機械学習アルゴリズムのハイパーパラメータ推定の自動化
 - “**fit**” 関数の **OptimizeHyperparameters** オプションを追加
パラメータ値固定

```
ctree = fitctree(x,y,'MaxNumSplits', 100);
```

ベイズ最適化によるパラメータ決定

```
ctree_BO = fitctree(x,y,...  
'OptimizeHyperparameters', 'auto');
```

- 目的関数を指定して **bayesopt** 関数を使用
 - 定義した制約内で最適化



ベイズ最適化最終結果

最適化が完了しました。

MaxObjectiveEvaluations の 30 に達しました。

関数評価数の合計: 30

経過時間の合計: 331.7751 秒。

目的関数の評価時間の合計: 312.013

最適な観測実行可能点:

MinLeafSize

154

観測された目的関数値 = 0.083359

推定される目的関数値 = 0.083665

関数の評価時間 = 10.7994

最適な推定実行可能点 (モデルに基づく):

MinLeafSize

168

推定される目的関数値 = 0.083665

推定される関数評価時間 = 10.9553

trainedClassifier_BO =

フィールドをもつ struct:

predictFcn: [function_handle]

RequiredVariables: {1×15 cell}

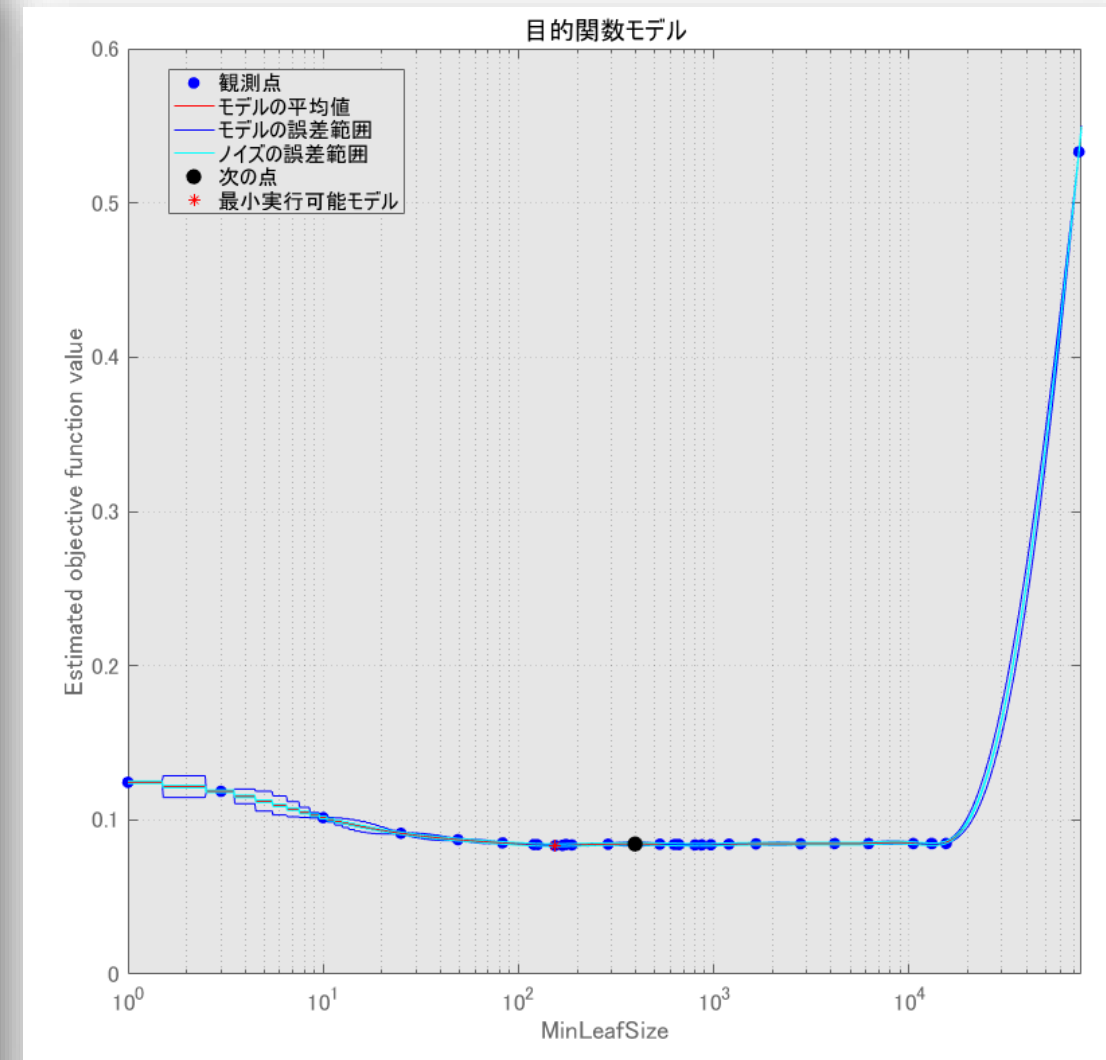
ClassificationTree: [1×1 ClassificationTree]

About: '...'

HowToPredict: '...'

validationAccuracy_BO =

0.9164



本日のトピック

■ 回帰分析

- 回帰分析ワークフロー
- 回帰モデルの見方
- サンプル数が少ないときの回帰分析

■ 分類

- 分類のワークフロー
- コーティングを簡単にする方法
- 機械学習の課題 – ハイパーパラメータの探索

■ **機械学習によるビッグデータ解析**

- メモリに収まりきれないビッグデータを扱う場合

回帰

分類

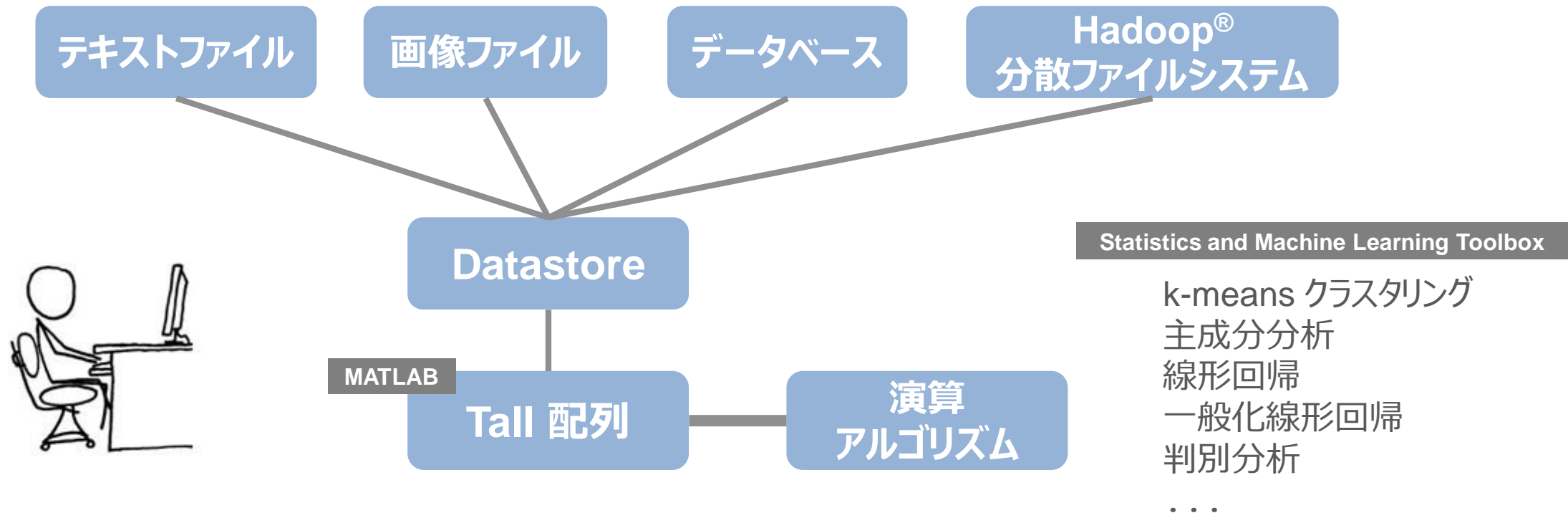
X

ビッグデータ解析

メモリに収まらないデータの扱い

user operation

- **tall**: メモリに収まりきらないデータに対するデータ型



Tall Array 対応関数リスト

https://jp.mathworks.com/help/releases/R2016b/matlab/import_export/functions-that-support-tall-arrays-by-type.html

例: New York のタクシー料金の予測モデルの作成 (回帰分析)

- 使用するデータ
 - 25GB の csv ファイル
- 解析ワークフロー
 - 解析に使用する特徴の選択
 - 前処理
 - データの探索
 - 料金予測モデルの構築
 - 料金モデルの評価

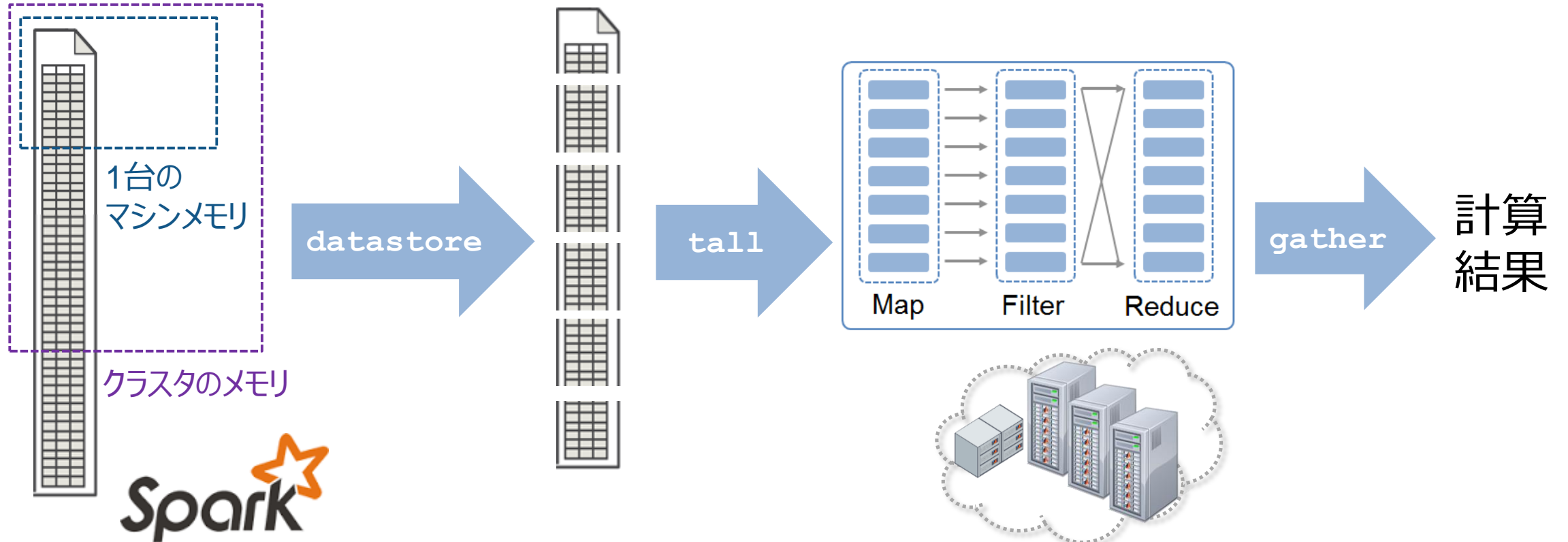


>> TallArrayDemo

メモリに収まらないデータの扱い

Tall Array のしくみ

- **gather** 関数で初めてすべての演算が実行
 - 演算をまとめることでデータへのアクセスを最適化 (遅延評価)



本日のトピック

回帰

分類

X

ビッグデータ解析

■ 回帰分析

- 回帰分析ワークフロー
- 回帰モデルの見方
- サンプル数が少ないときの回帰分析: **ガウス過程回帰 R2015b**

■ 分類

- 分類のワークフロー
- コーティングを簡単にする方法: **分類学習器アプリ R2015a**
- ハイパーパラメターの探索: **ベイズ最適化 R2016b**

■ 機械学習によるビッグデータ解析

- メモリに収まりきれないビッグデータを扱う場合: **Tall 配列 R2016b**

機械学習のトレーニングも
あります！

MATLAB Answers

日本語/英語の Q&A サイト

- MATLAB に関する過去の質問 & 回答が閲覧可能
- MathWorks アカウントがあれば、誰でも投稿できます！
- 日本語/英語両方に対応
- 得意な分野の質問への回答に挑戦してみませんか？



The screenshot shows the MATLAB Answers website interface. At the top, there is a navigation bar with the MathWorks logo and links for products, solutions, academia, support, community (English), events, and company information. Below this is the 'MATLAB Answers' header with a search bar and a dropdown menu. The main content area is divided into two columns. The left column contains filters for language (Japanese and English), status (Answer Accepted and Answered), source (MathWorks Support), and product (MATLAB, Simulink, Communications System). The right column is titled 'Recently Added' and shows a list of questions. The first question is 'スタティックテキストの更新が遅いのはなぜですか?' (Why is the update of static text slow?) with 1 answer. The second question is 'ライセンス マネージャーを再起動したり MATLAB を終了せずに、どのようにして利用可能なキーのプールにツール ボックスのライセンス キーを解放または返却できますか?' (How can I release or return the license key of the tool box to the pool of available keys without restarting the license manager or ending MATLAB?) with 1 answer. A cartoon character with a question mark above its head is standing on the right side of the screenshot.



© 2016 The MathWorks, Inc. MATLAB and Simulink are registered trademarks of The MathWorks, Inc. See www.mathworks.com/trademarks for a list of additional trademarks. Other product or brand names may be trademarks or registered trademarks of their respective holders.