

2023年6月1日
データサイエンティスト集会 in VRC

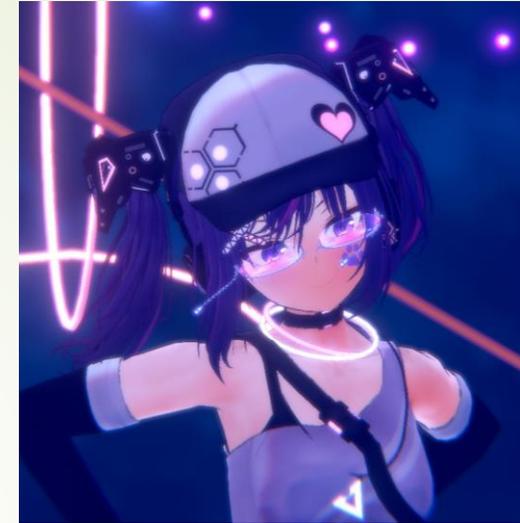
実務で使える Orange Data Miningの便利な機能

1

ぶんちん

自己紹介 ぶんちゃん

- ▶ 複合経営が特徴の企業（製造業）に所属
- ▶ データ分析担当者だったが。。。



Orange Data Mining

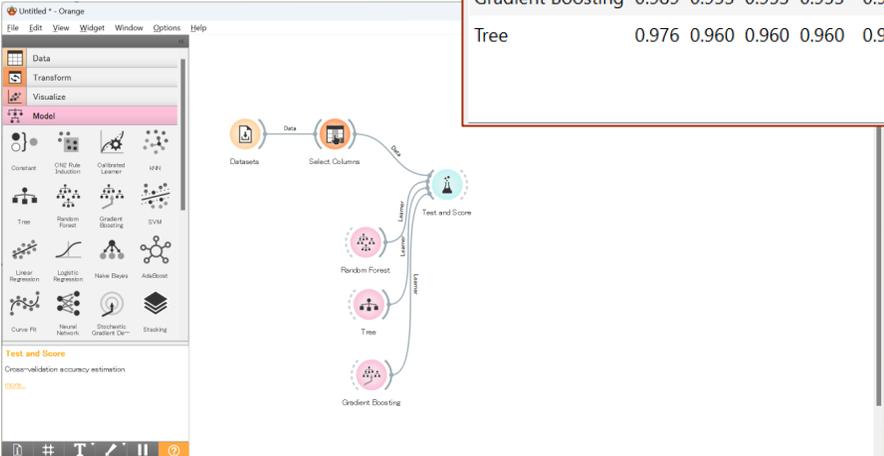
<https://orangedatamining.com/>

- ▶ ビジュアルプログラミング的にデータ分析や機械学習モデル作成・評価が可能
- ▶ 公式HPから入手すれば無料で使用可能

▶ UIが素晴らしい

- ▶ 初心者は勉強に使おう！
- ▶ 専門家は手抜き・教育に使おう！

基本的な使い方は前回紹介
スライド&動画公開してます



The screenshot shows the Orange Data Mining GUI. On the left is a widget palette with categories like Data, Transform, Visualize, and Model. The main workspace contains a workflow: Data -> Select Columns -> Train -> Random Forest, Tree, and Gradient Boosting. A 'Test and Score' widget is connected to the models. An evaluation results window is open, displaying a table of performance metrics for three models.

Evaluation results for target (None, show average over classes)						
Model	AUC	CA	F1	Prec	Recall	MCC
Random Forest	0.988	0.953	0.953	0.953	0.953	0.930
Gradient Boosting	0.989	0.953	0.953	0.953	0.953	0.930
Tree	0.976	0.960	0.960	0.960	0.960	0.940

GUI操作で分析・モデル作成が可能

便利な機能紹介

Orangeには様々な便利な機能が実装されています。

今回はその中から、ちょっと見ただけでは気づきそうにないけど便利なものを中心に紹介します。

- 入門者向け：誰もが使える便利な機能
- 中級者向け：機械学習についてある知識がある人向けの機能
- アドオン：特定の領域に特化した追加機能

入門者向け

- ▶ 基本統計量の一括出力
- ▶ 散布図の便利機能×2

入門者向け：基本統計量の一括出力

基本統計量を表示したいデータと接続するだけ

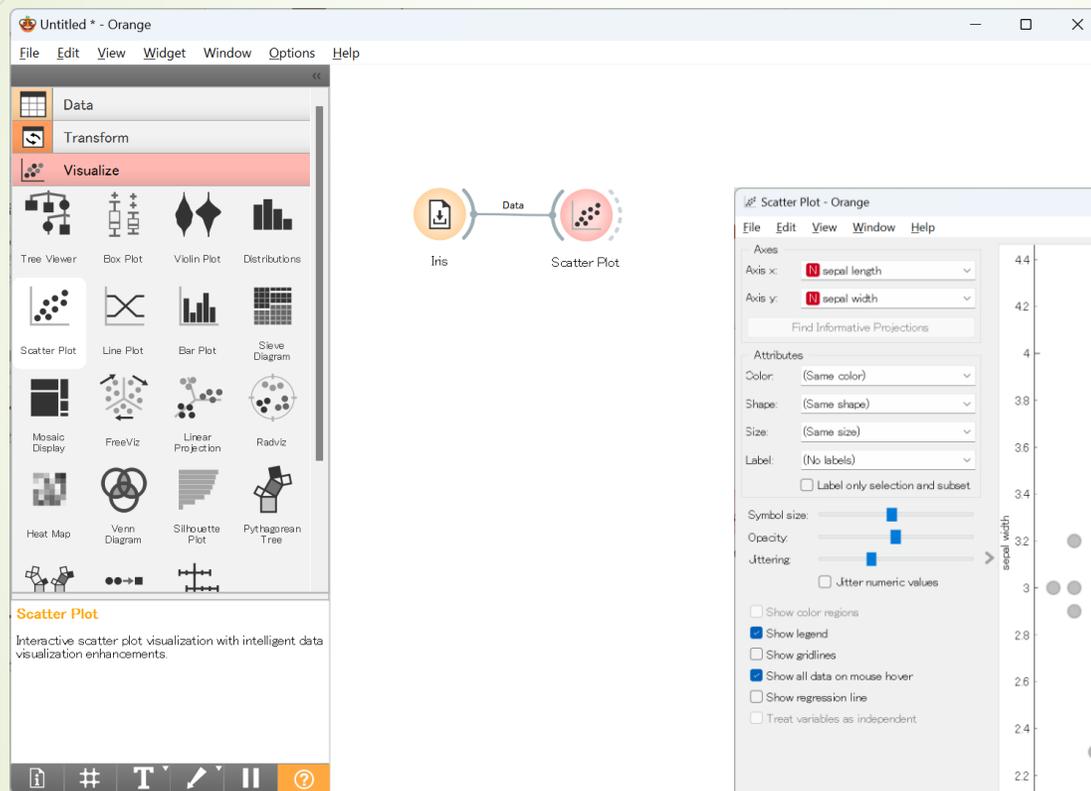
The screenshot shows the Orange software interface. The main window displays a workflow with two widgets: 'Iris' and 'Feature Statistics'. The 'Feature Statistics' widget is selected, and its output window is shown in the foreground. The output window displays a table of basic statistics for the features: sepal length, sepal width, petal length, and petal width. Each feature has a corresponding histogram. The 'iris' feature is also listed at the bottom, with a color selection dropdown set to 'iris'.

Name	Distribution	Mean	Mode	Median	Dispersion	Min.	Max.	Missing	
sepal length		5.843	5.0	5.8	0.141	4.3	7.9	0 (0%)	
sepal width		3.054	3.0	3.0	0.142	2.0	4.4	0 (0%)	
petal length		3.759	1.5	4.350	0.468	1.0	6.9	0 (0%)	
petal width		1.199	0.2	1.3	0.635	0.1	2.5	0 (0%)	
iris							Iris-setosa	1.1	0 (0%)

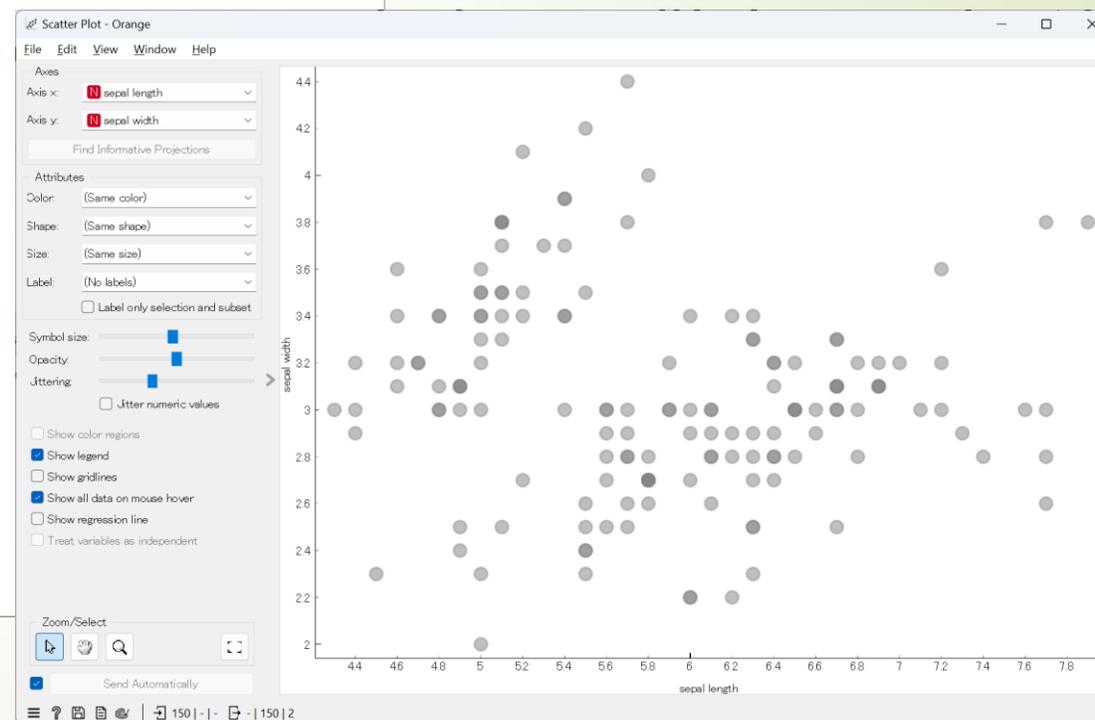
Color: iris ← colorを指定するとグラフに反映される

Send Automatically

入門者向け：散布図の便利機能 1



散布図を表示したいデータと接続、
縦軸と横軸の項目を指定するだけ



これだけだと普通すぎて面白くないですよね？

入門者向け：散布図の便利機能 1

表形式で表示するwidget

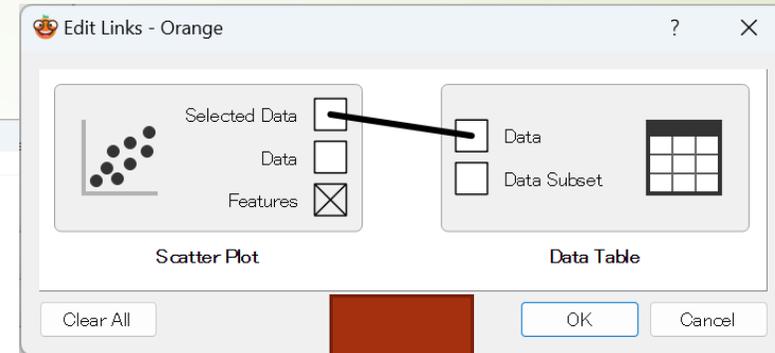
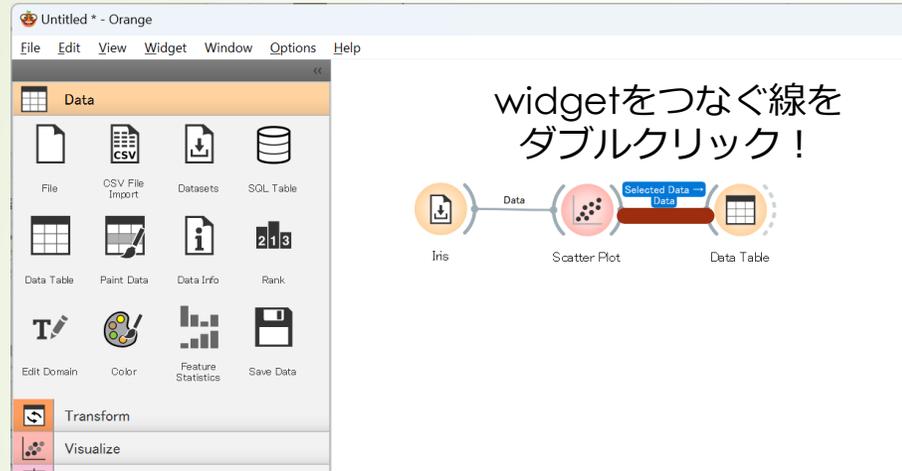
マウスで範囲指定

選択したデータを抽出

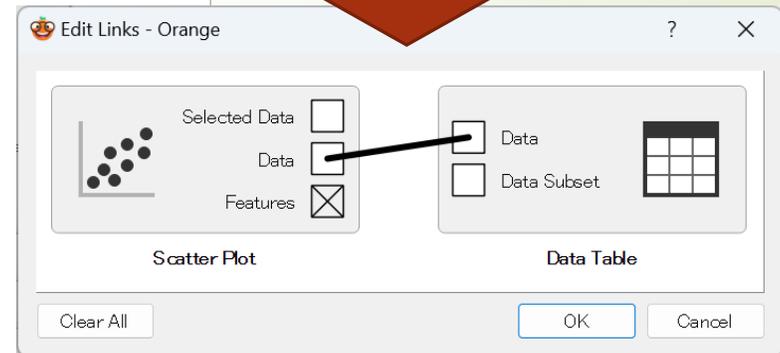
	iris	Group	sepal length	sepal width	petal length	petal width
1	Iris-virginica	G1	7.2	3.6	6.1	2.1
2	Iris-virginica	G1	7.7	3.8	6.7	2.1
3	Iris-virginica	G1	7.9	3.8	6.4	2.1

散布図どころかグラフのwidgeに限らず、
全てのデータ可視化機能からGUIでデータ選択が可能

入門者向け：散布図の便利機能 1



接続を変えると



この機能を使ってフィルターかければ、マウス操作で異常値の除去などが可能即時、他の分析結果に反映できる！

	iris	Selected	sepal length	sepal width	petal length	petal width
110	Iris-virginica	Yes	7.2	3.6	6.1	
118	Iris-virginica	Yes	7.7	3.8	6.7	
132	Iris-virginica	Yes	7.9	3.8	6.4	
1	Iris-setosa	No	5.1	3.5	1.4	
2	Iris-setosa	No	4.9	3.0	1.4	
3	Iris-setosa	No	4.7	3.2	1.3	
4	Iris-setosa	No	4.6	3.1	1.5	
5	Iris-setosa	No	5.0	3.6	1.4	
6	Iris-setosa	No	5.4	3.9	1.7	
7	Iris-setosa	No	4.6	3.4	1.4	
8	Iris-setosa	No	5.0	3.4	1.5	

選択有無のフラグ情報を取得可能

注意：恣意的なデータ選択はやめよう

入門者向け：散布図の便利機能 2

散布図を表示したいデータと接続、縦軸と横軸の項目を指定するだけ

① 目的変数を設定 →

② 押す

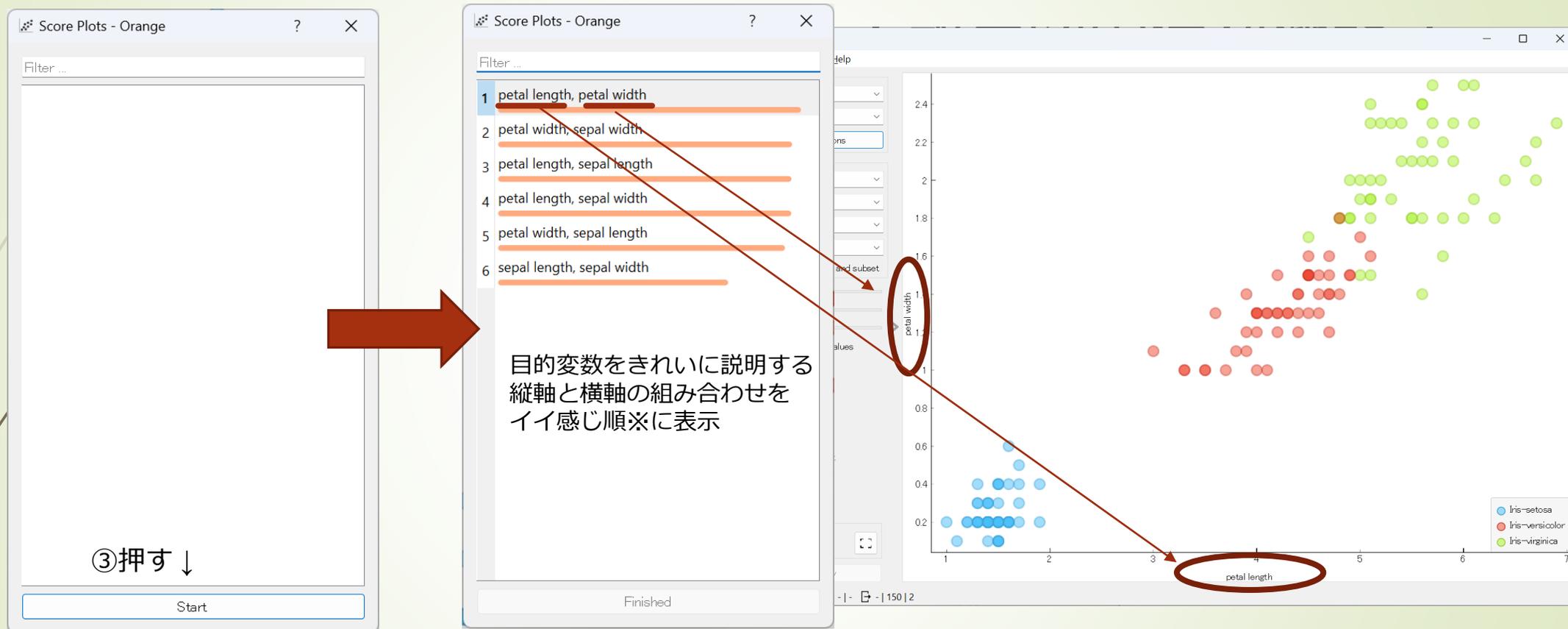
The screenshot shows the Orange data mining software interface. On the left, the 'Visualize' widget shelf contains various visualization options, with 'Scatter Plot' selected. The main workspace shows a workflow: 'Iris' data source connected to a 'Scatter Plot' widget. A configuration window for the 'Scatter Plot' widget is open, showing the following settings:

- Axes: Axis x: sepal length, Axis y: sepal width
- Attributes: Color: iris, Shape: (Same shape), Size: (Same size), Label: (No labels)
- Symbol size: [slider]
- Opacity: [slider]
- Jittering: [checkbox] Jitter numeric values
- Options:
 - Show color regions: [checkbox]
 - Show legend: [checkbox]
 - Show gridlines: [checkbox]
 - Show all data on mouse hover: [checkbox]
 - Show regression line: [checkbox]
 - Treat variables as independent: [checkbox]
- Zoom/Select: [checkbox] Send Automatically

The scatter plot itself displays data points colored by the 'iris' attribute, with 'sepal length' on the x-axis and 'sepal width' on the y-axis. A legend in the bottom right corner identifies the species: Iris-setosa (blue), Iris-versicolor (red), and Iris-virginica (green).

これだけだと普通すぎて面白くないですね？

入門者向け：散布図の便利機能 2



※イイ感じ順 の具体的な内容
knn(k=10)で全ての特徴量の組み合わせでモデル作成・評価
精度（オレンジ色のバー）の良い順に表示する

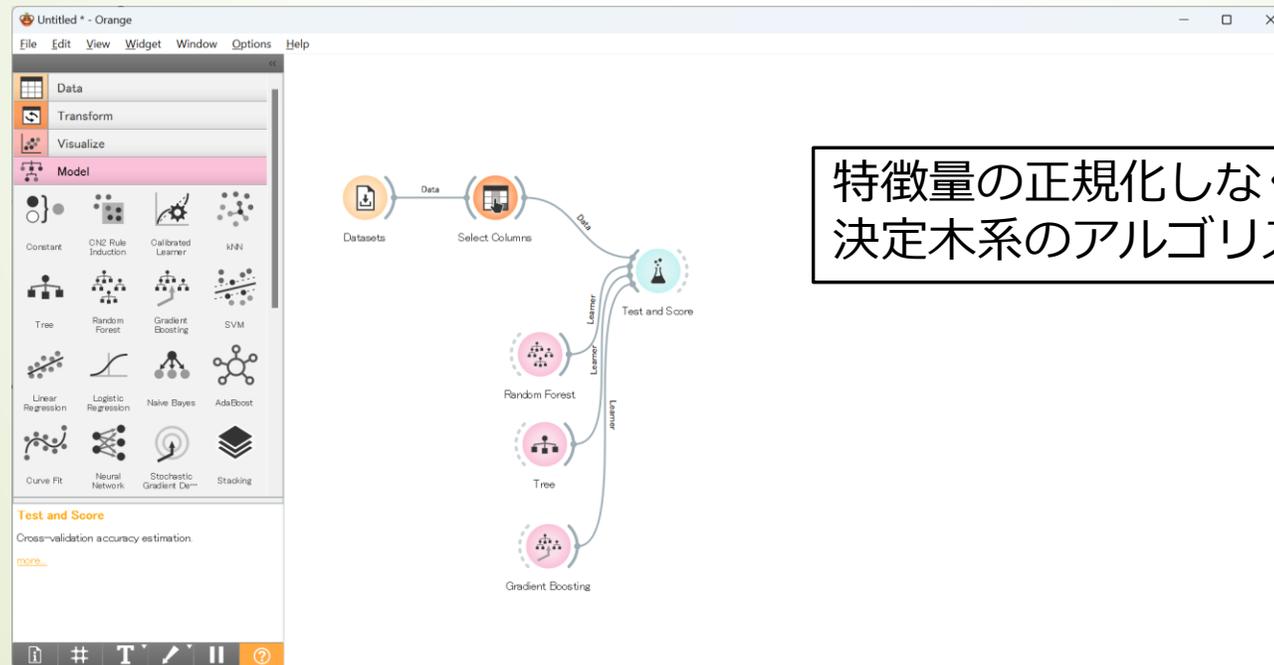
カテゴリ変数や相関係数には表れない相関を扱える
非専門家向けのデータ可視化に便利！

中級者向け

- ▶ モデル作成のデータの前処理
- ▶ 異常データの除去（異常検知）

中級者向け：モデル作成のデータ前処理

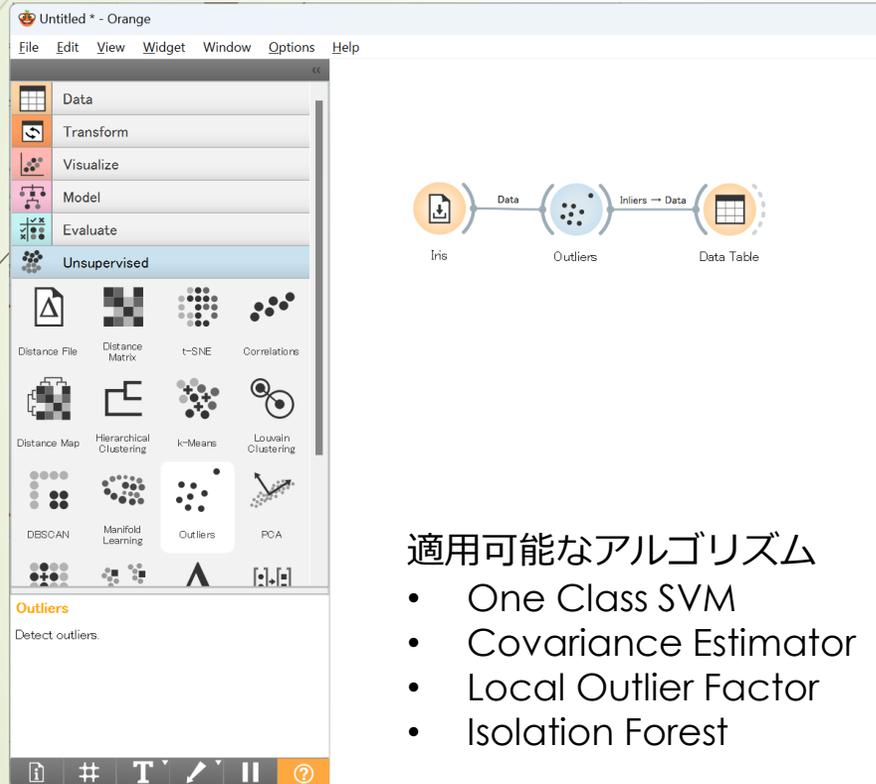
- ▶ 前回の資料、説明をシンプルにするため、適切に機械学習モデルを作成するために恣意的なことをしていました。



特徴量の正規化しなくても問題ない
決定木系のアルゴリズムを選択

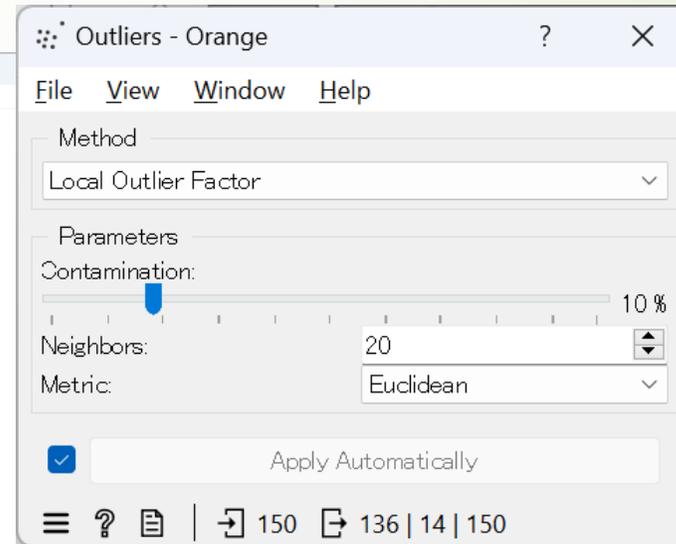
正規化をはじめ、便利なデータ前処理機能を説明します

中級者向け：異常データの除去（異常検知）



適用可能なアルゴリズム

- One Class SVM
- Covariance Estimator
- Local Outlier Factor
- Isolation Forest



Outliers - Orange

File View Window Help

Method
Local Outlier Factor

Parameters
Contamination: 10 %
Neighbors: 20
Metric: Euclidean

Apply Automatically

150 136 | 14 | 150

異常検知アルゴリズムを使い、
一定比率の異常値を簡単に除去可能

前述のwidgetの接続を変えれば、
逆に異常データの抽出も可能

アドオン

- ad-on（追加機能）の導入方法
- 特徴量重要度の算出

アドオン : ad-on (追加機能) の導入方法

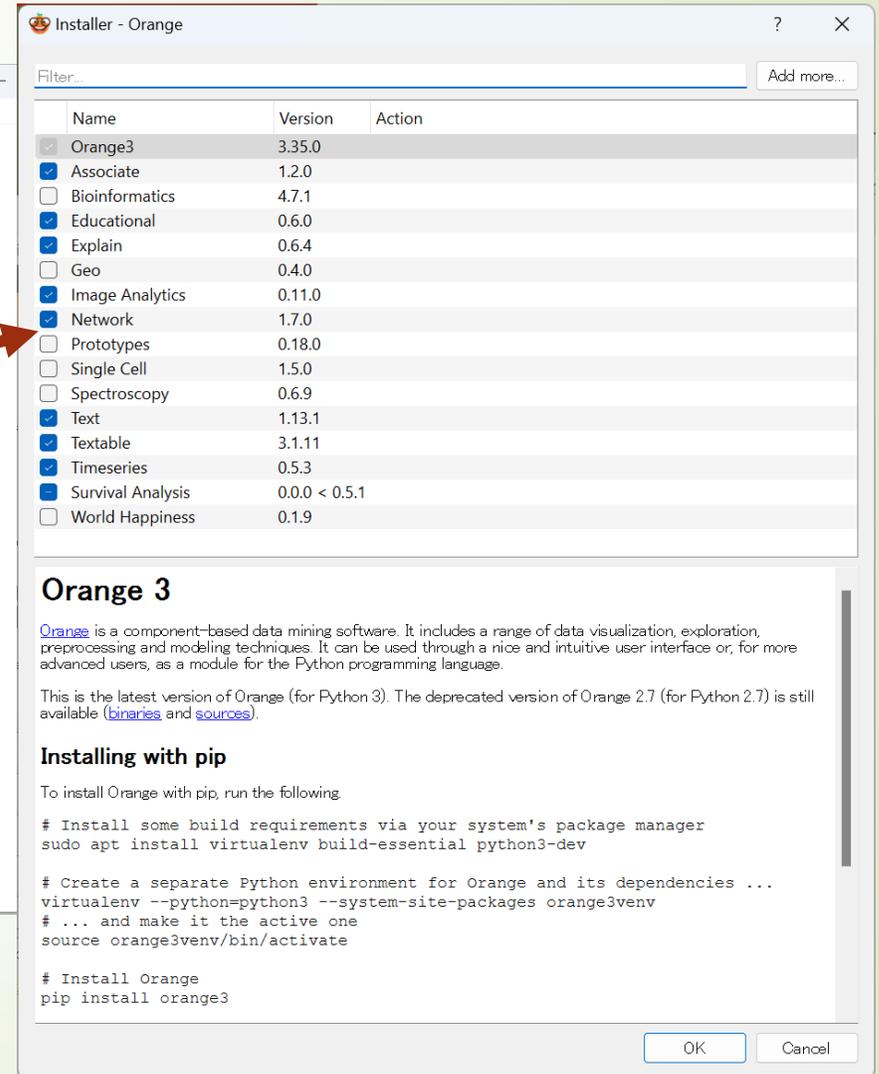


導入したいアドオンにチェック

例えば

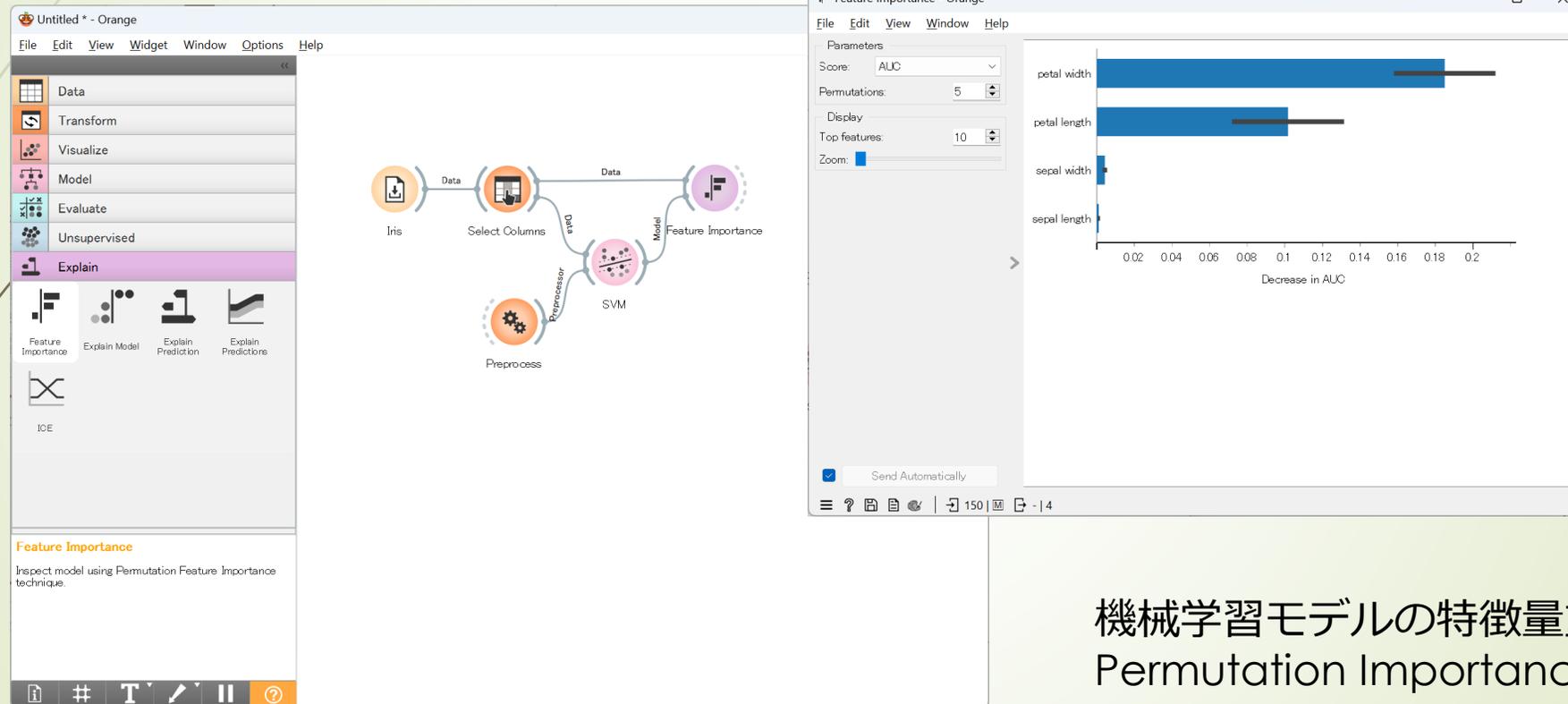
- 機械学習モデルの説明
 - 画像処理
 - 自然言語
 - 時系列分析
 - 生存分析
 - ネットワーク分析 など
- 様々な領域の手法に対応

必要に応じて導入してください
あまり導入しすぎると起動が重くなるので注意



アドオン：特徴量重要度の算出

Explainのad-on



機械学習モデルの特徴量重要度を
Permutation Importanceで評価

補足) Permutation Importanceを使ってモデルがどの特徴量から学習したかを定量化する
<https://www.datarobot.com/jp/blog/permutation-importance/>

ご清聴、ありがとうございました。

他にも話したいネタがたくさんあります

- ▶ 超初心者向け機械学習の考え方
- ▶ 組織の基礎レベル向上 ノーコード分析の紹介
- ▶ データ分析プロジェクトの進め方ネタ
- ▶ あまり知られていない良書紹介 など



Twitter : @bunnchinn3

今後もLTでいろんなお話をしていきたいです。

どれにするかtwitterでアンケートを考えているので、投票してもらえると嬉しいです。

詳しい内容については、個別に対応するのでお気軽にお声がけください。