

Introduction to tsfresh

AKLab, 22.03.2024, ABI House

Andreas W. Kempa-Liehr

Department of Engineering Science and Biomedical Engineering



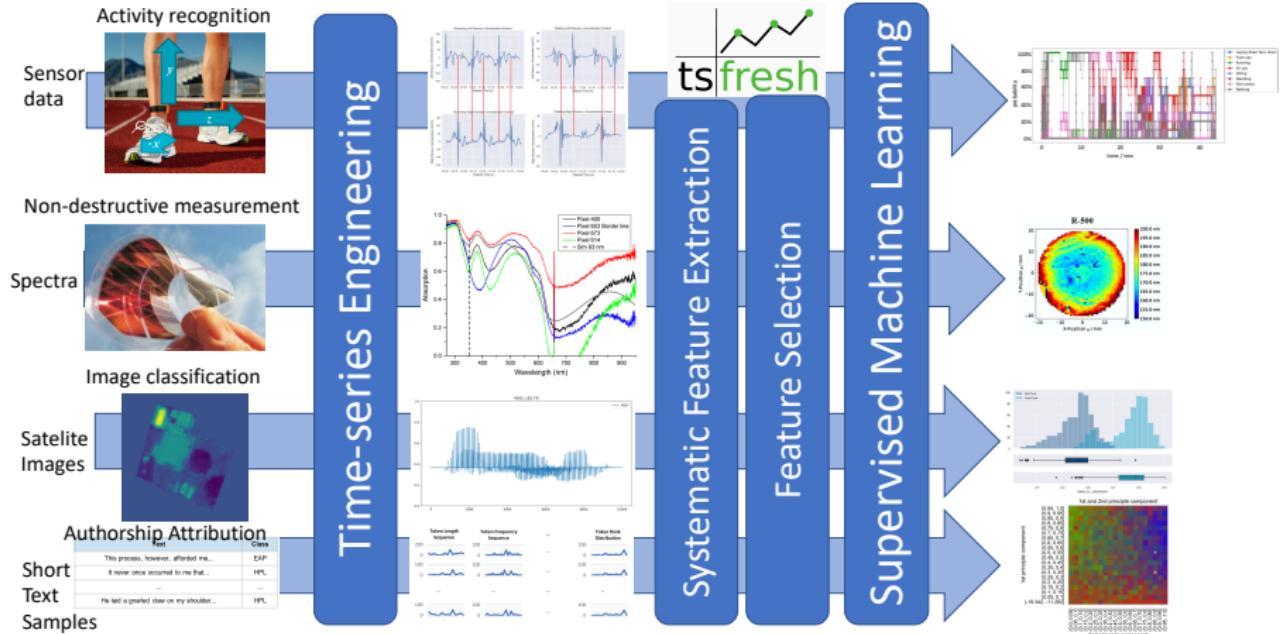
ENGINEERING

22.03.2024 – AKLab

Outline

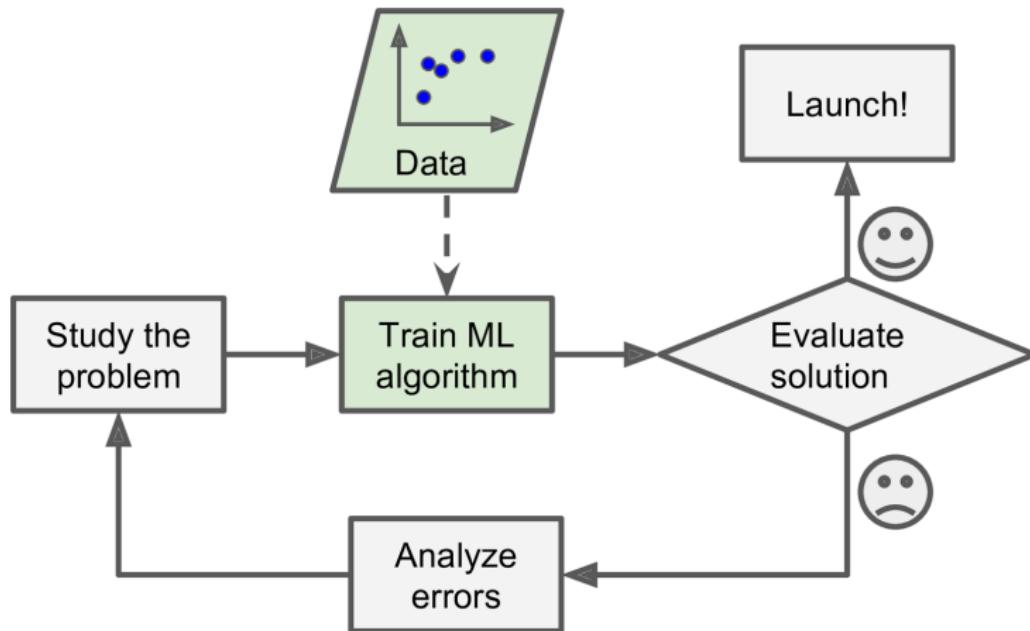
1. Overview
2. Machine Learning explained
3. tsfresh
4. Logistic Regression (Decision model for signal classification)
5. Random Forest

Applications of tsfresh



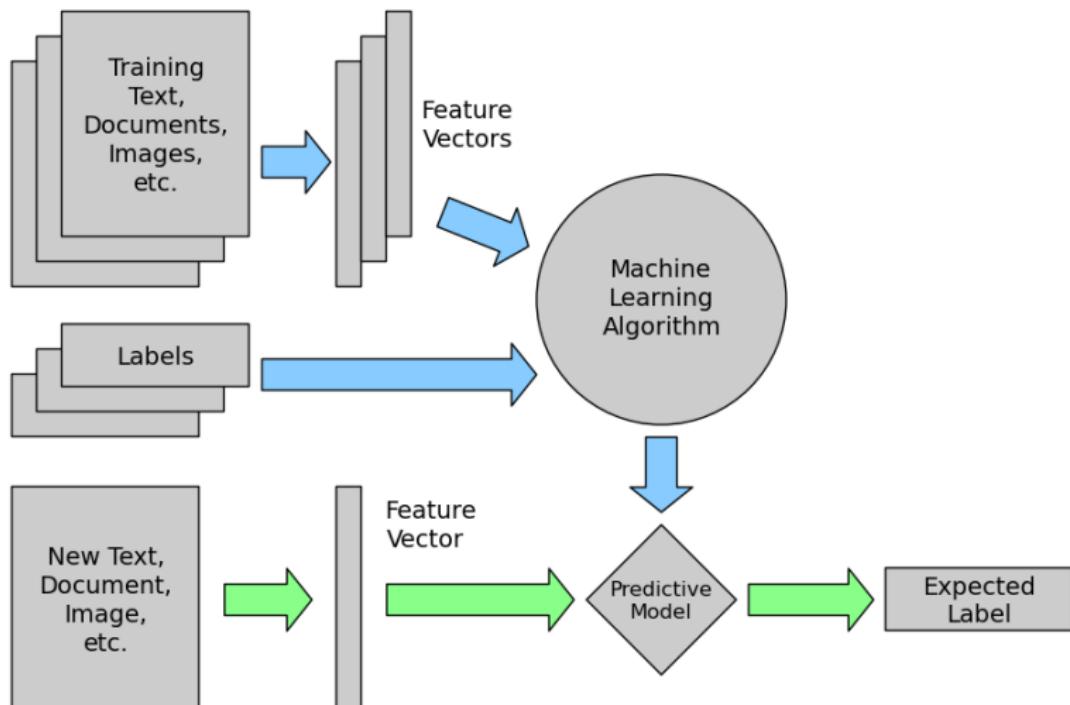
Machine learning

Training a ML algorithm is also called *fitting*.



Aurélien Géron. *Hands-On Machine Learning with Scikit-Learn and TensorFlow*. 4th release. Sebastopol, CA, United States: O'Reilly Media, 2017, Fig. 1-2

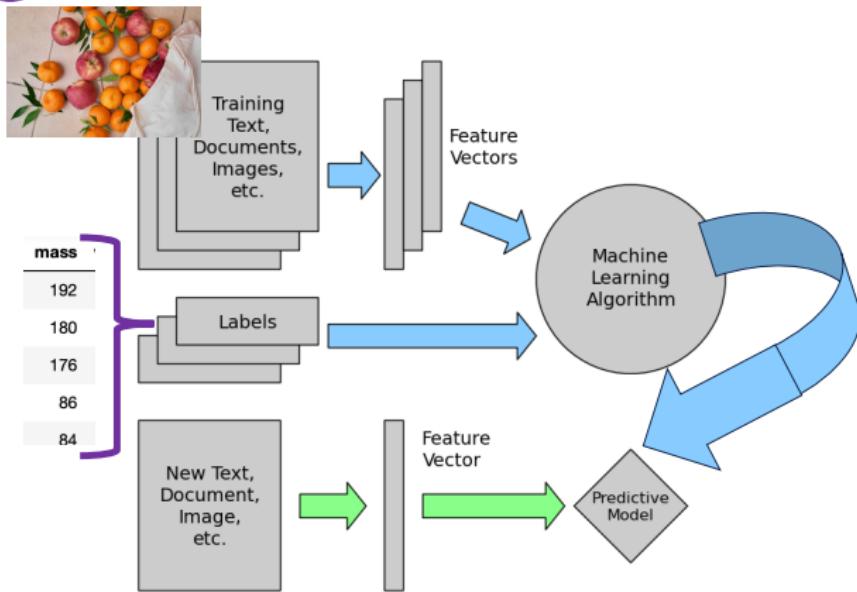
Flowchart of supervised learning



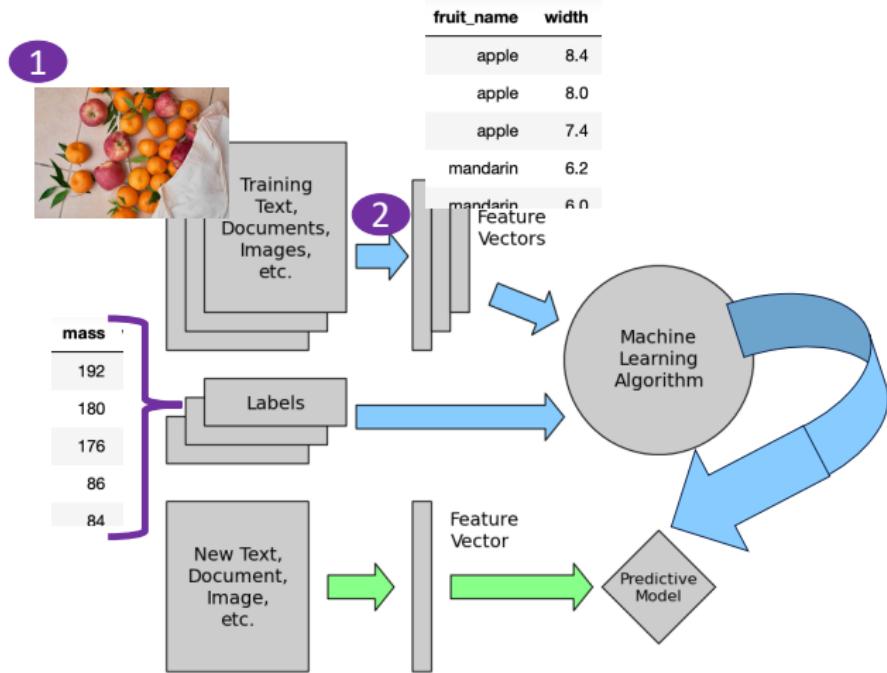
Jakob VanderPlas. *Astronomy with scikit-learn*. Release Scipy2012. Seattle: University of Washington, 2012. URL: http://astroML.github.com/sklearn_tutorial

Machine learning in action - (1) Data collection

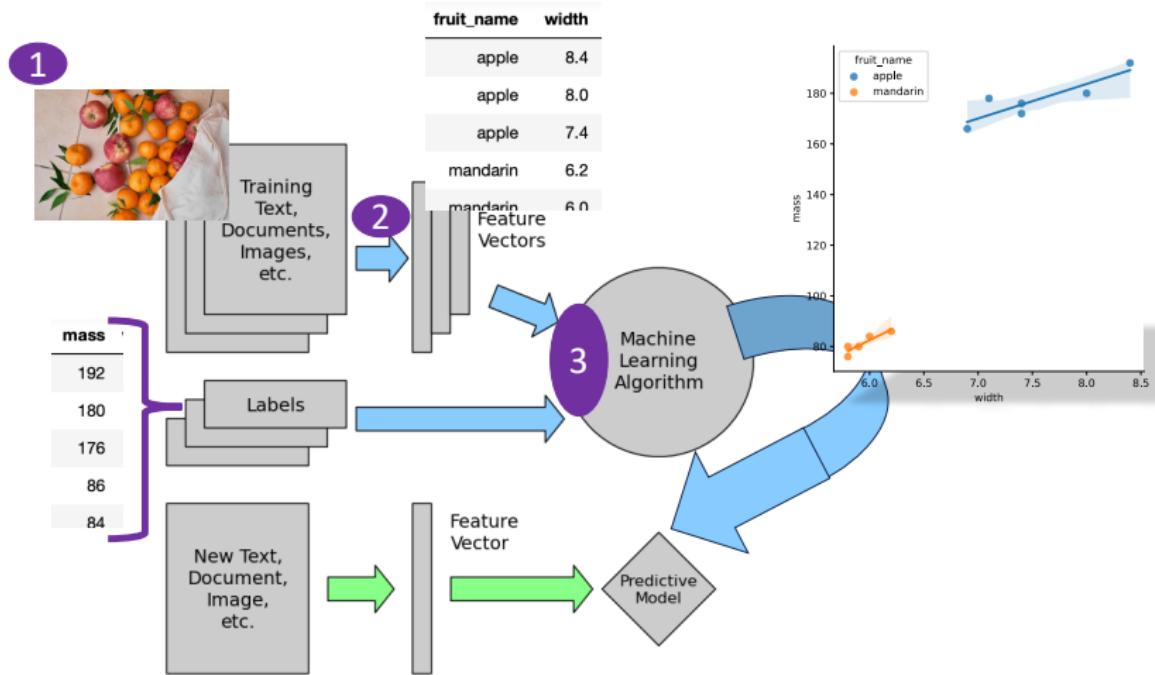
1



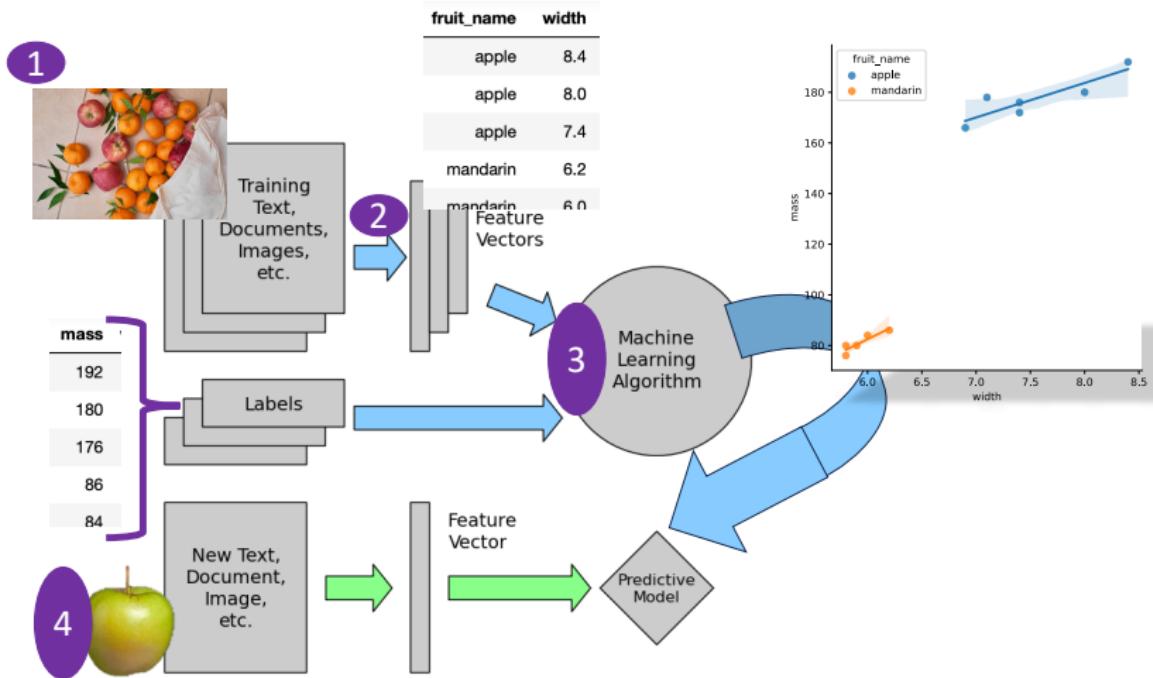
Machine learning in action - (2) Feature extraction



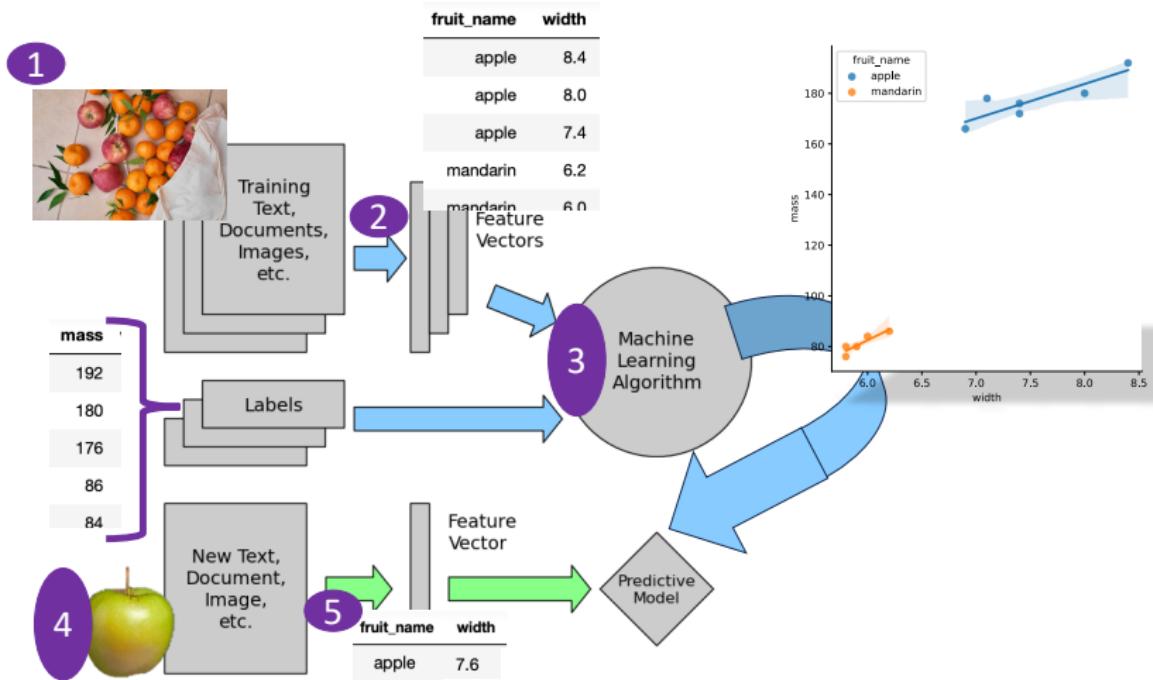
Machine learning in action - (3) Model fitting/training



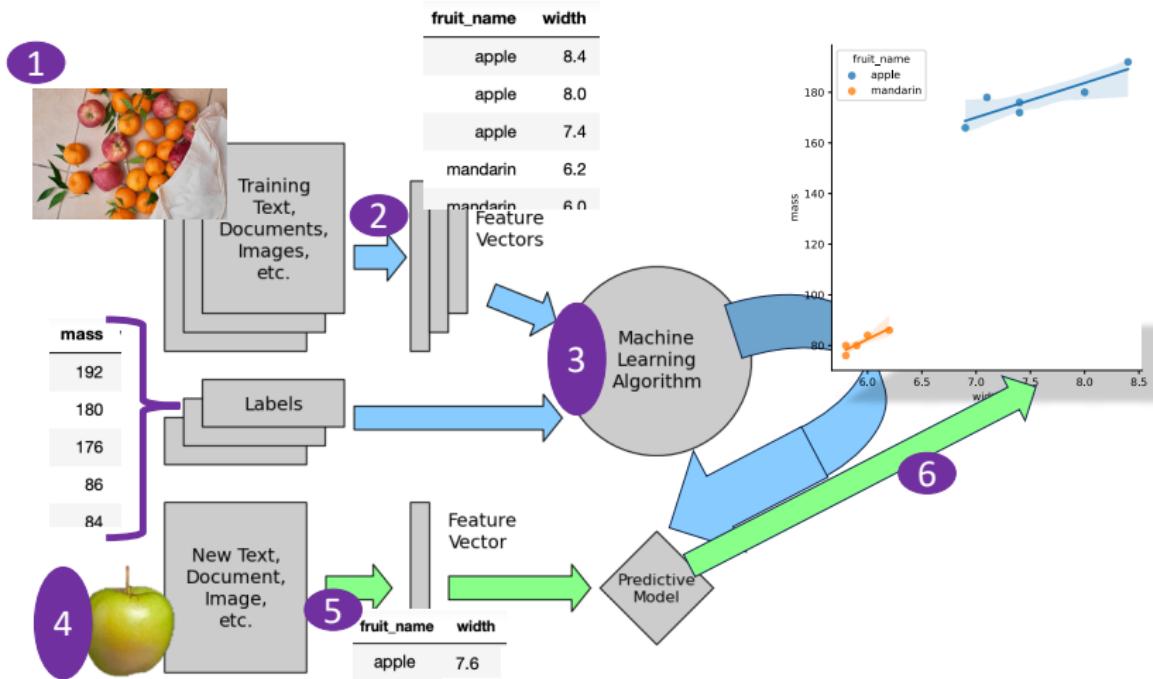
Machine learning in action - (4) New observation



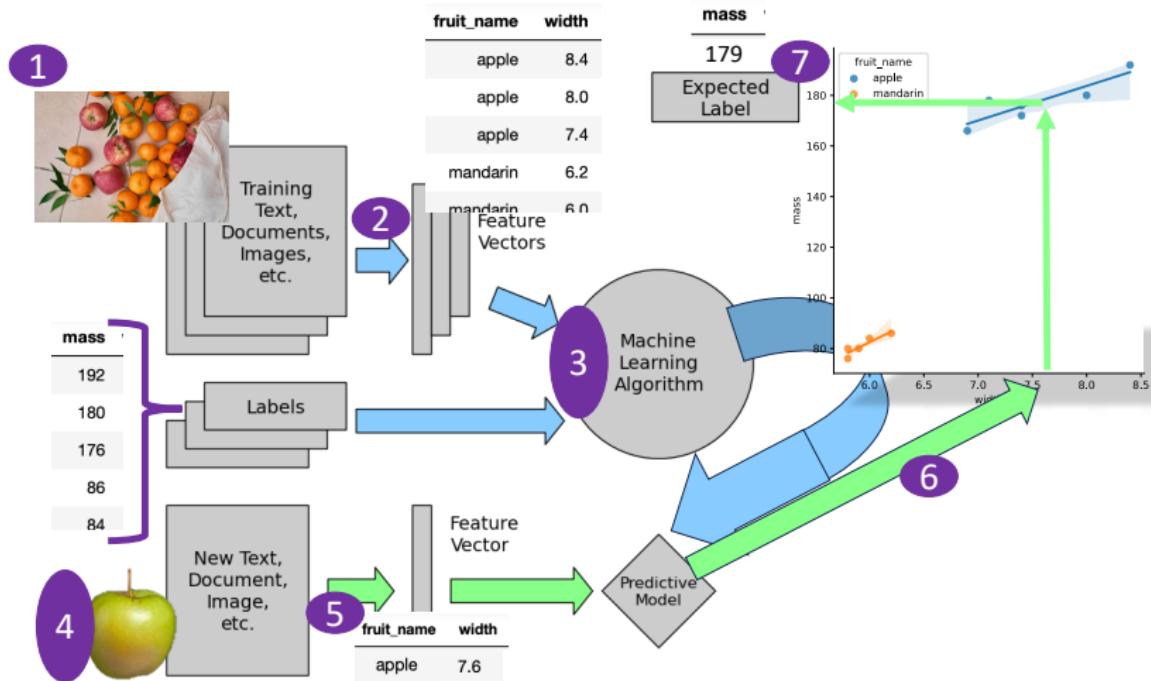
Machine learning in action - (5) Feature extraction



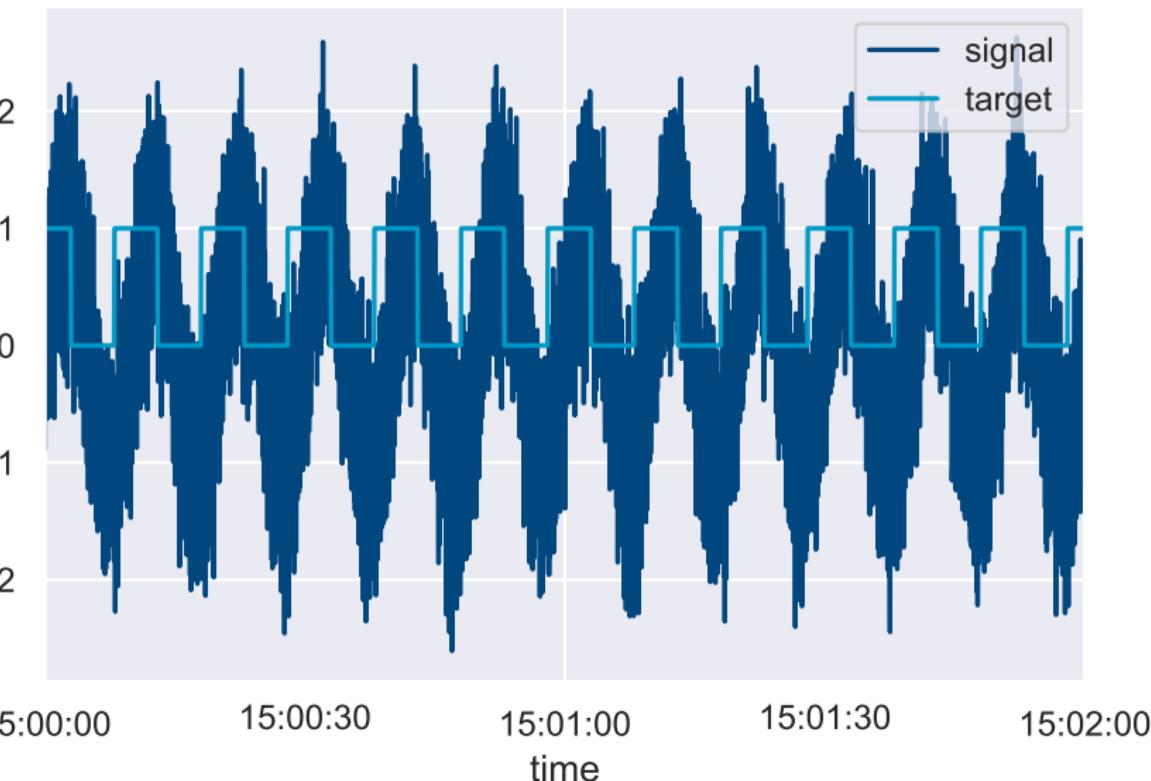
Machine learning in action - (6) Use model



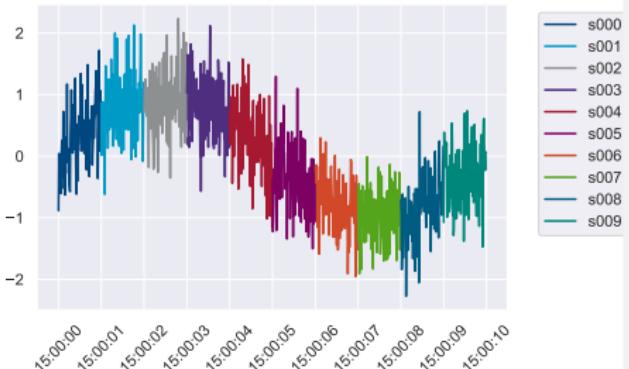
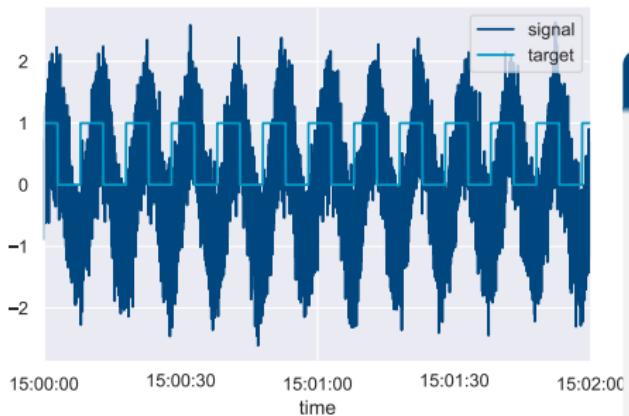
Machine learning in action - (7) Prediction



Noisy sine example



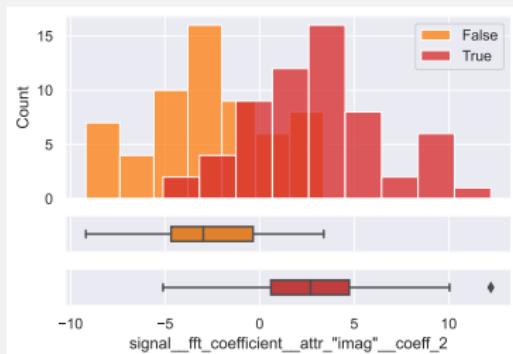
Signal Processing Example: Noisy Sine



$T = 100$ $N = 120$ $K = 778$

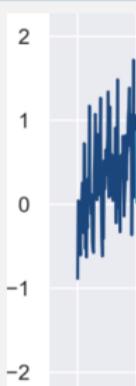
Systematic feature engineering X_ϕ

$$\begin{pmatrix} \phi_1(\vec{x}_1) & \dots & \phi_k(\vec{x}_1) & \dots & \phi_K(\vec{x}_1) \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \phi_1(\vec{x}_i) & \dots & \phi_k(\vec{x}_i) & \dots & \phi_K(\vec{x}_i) \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \phi_1(\vec{x}_N) & \dots & \phi_k(\vec{x}_N) & \dots & \phi_K(\vec{x}_N) \end{pmatrix}$$



Signal classification - A special type of signal operator

Signal classification



input signal $\vec{x} \rightarrow$ Signal Operator $\rightarrow y = (\mathbb{P}[y = \text{True} | \vec{x}] \geq \frac{1}{2})$

Logistic regression

- Let $X_\phi \in \mathbb{R}^{N \times \kappa}$ be the feature matrix of $1 \leq \kappa \leq K$ statistically significant time-series features, and $\phi_k(\vec{x}_1), \dots, \phi_k(\vec{x}_N)$ its k^{th} feature column.
- Let y_1, \dots, y_N be the labels of signals $\vec{x}_1, \dots, \vec{x}_N$.
- Probability $\mathbb{P}[y = \text{True} | \vec{x}]$ that target y is True given observed signal \vec{x} is predicted as

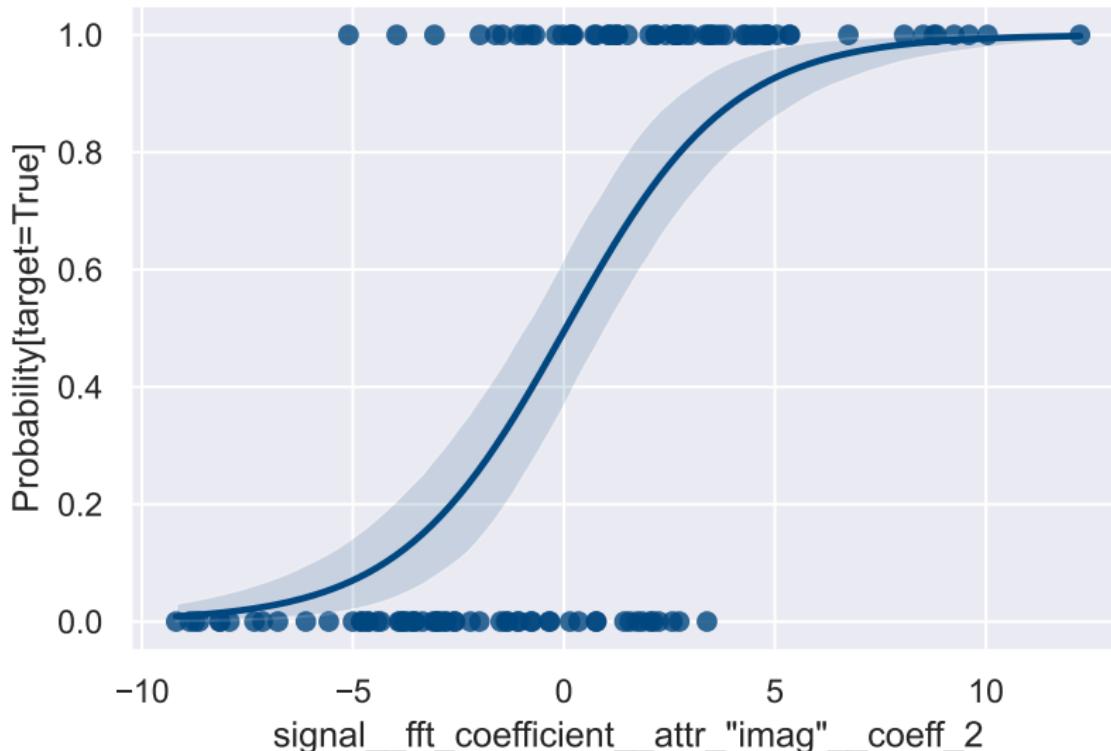
$$\mathbb{P}[y = \text{True} | \mathbf{x}] = \left(1 + e^{-\left(\hat{w}_0 + \sum_{k=1}^K \hat{w}_k \phi_k(\mathbf{x}) \right)} \right)^{-1}$$

$$\hat{w}_0, \hat{w}_1, \dots, \hat{w}_k = \arg \min_{w_0, \dots, w_k} \left[\sum_{i=1}^N -\log \mathbb{P}[y = y_i | \mathbf{x}_i] + \frac{1}{C} \sum_{k=0}^K |w_k| \right]$$

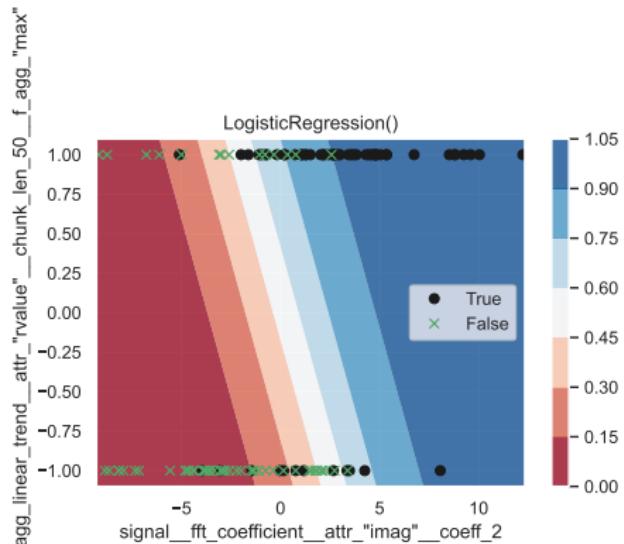
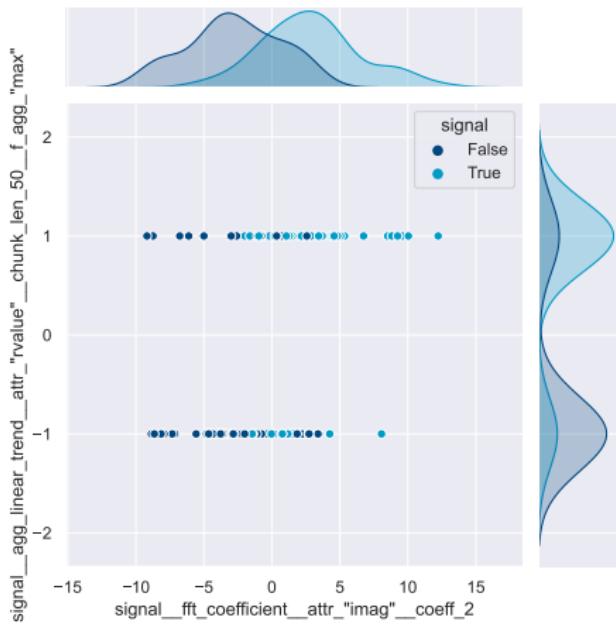
with C being the inverse regularisation strength.

Logistic regression with one time-series feature

$$C = 1.0, w_0 = -0.01387571, w_1 = 0.50363426$$

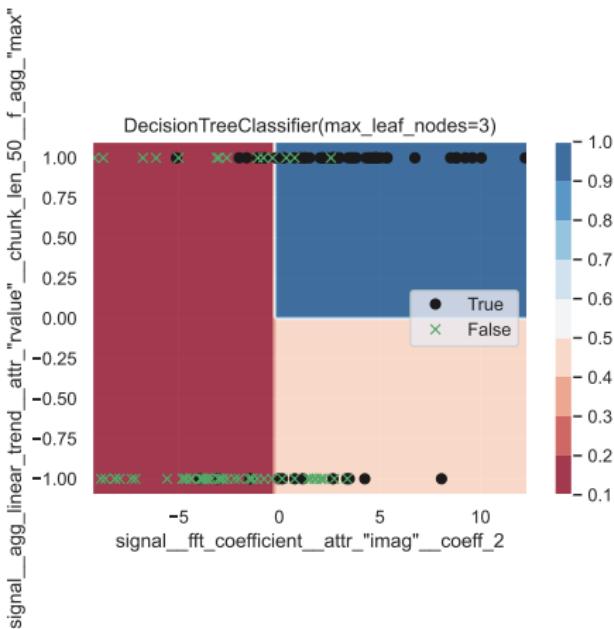
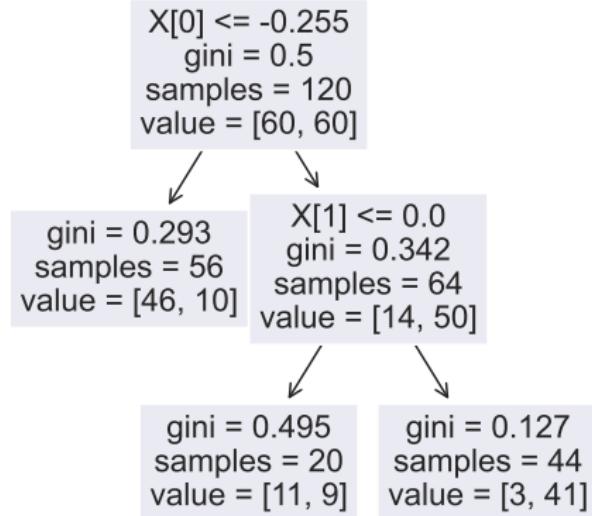


Logistic regression with two time-series features



Decision tree with two time-series features

```
X[0] = signal_fft_coefficient_attr_"imag"__coeff_2'
X[1] = signal_agg_linear_trend_attr_"rvalue"__chunk_len_50__f_agg_"max"
```



$$\text{Gini impurity bottom left node } G = 1 - \left(\frac{11}{20}\right)^2 - \left(\frac{9}{20}\right)^2 = 0.495$$

Random Forests

Definition

Random forests or random decision forests are an ensemble learning method for classification, regression and other tasks that operates by constructing a multitude of decision trees at training time and outputting the class that is the mode of the classes (classification) or mean prediction (regression) of the individual trees.

Random decision forests correct for decision trees' habit of overfitting to their training set.

How does it work?

The Random Forest algorithm introduces extra randomness. Instead of searching for the very best feature when splitting a node, it search for the best feature among a random subset of features.

Random forest with two time-series features

