# tsfresh Documentation

*Release 0.11.1*

**Sep 07, 2018**

# Contents

This is the documentation of **tsfresh**.

tsfresh is a python package. It automatically calculates a large number of time series characteristics, the so called features. Further the package contains methods to evaluate the explaining power and importance of such characteristics for regression or classification tasks.
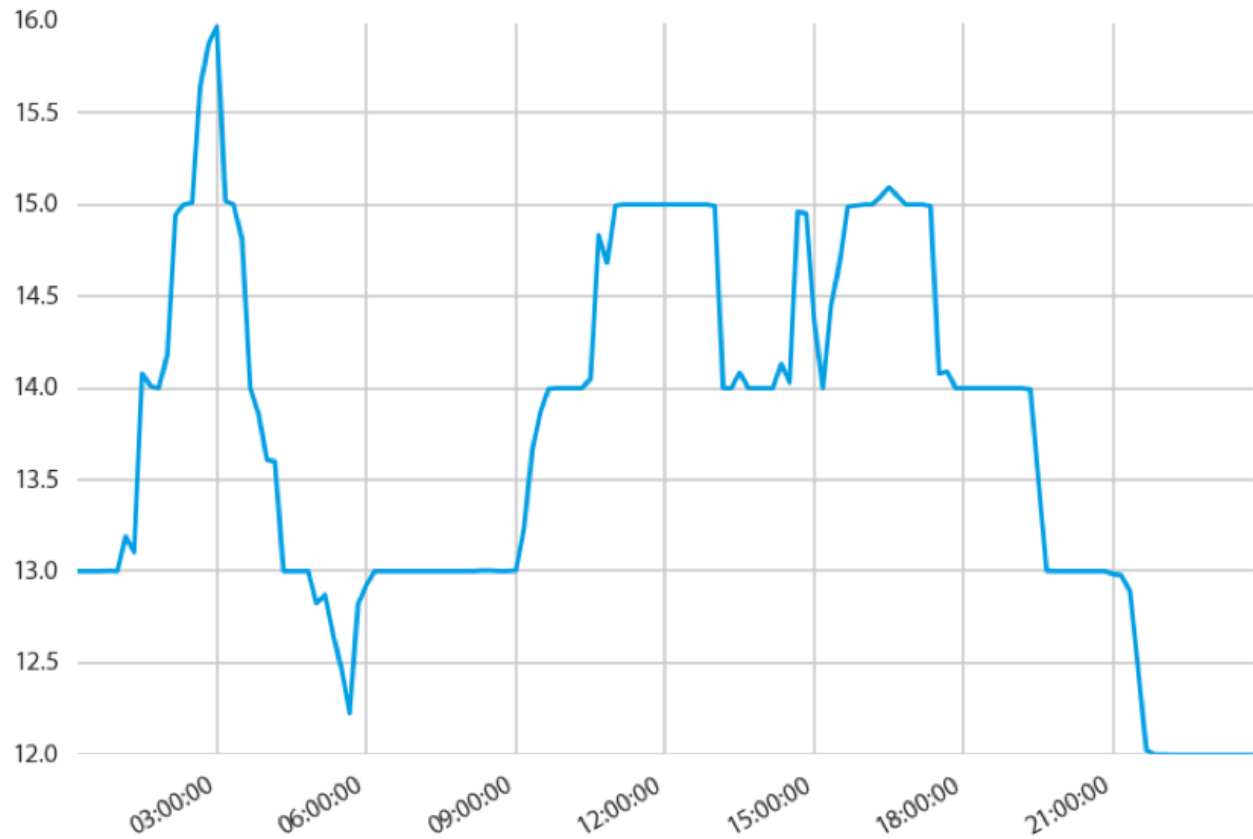
You can jump right into the package by looking into our *Quick Start*.

# Contents

The following chapters will explain the tsfresh package in detail:
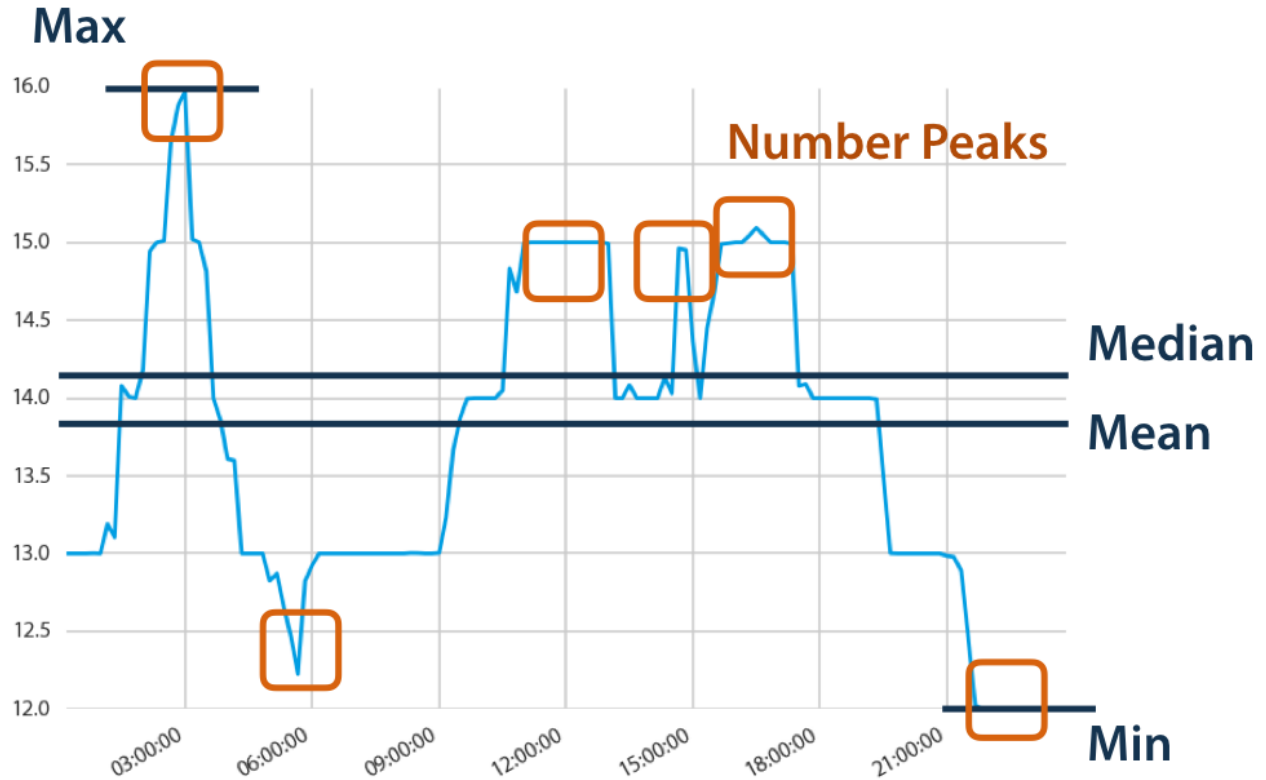
## 1.1 Introduction

### 1.1.1 Why do you need such a module?

tsfresh is used to to extract characteristics from time series. Let's assume you recorded the ambient temperature around your computer over one day as the following time series:

Now you want to calculate different characteristics such as the maximal or minimal temperature, the average temperature or the number of temporary temperature peaks:

Without tsfresh, you would have to calculate all those characteristics by hand. With tsfresh this process is automated and all those features can be calculated automatically.

Further tsfresh is compatible with pythons `pandas` and `scikit-learn` APIs, two important packages for Data Science endeavours in python.

### 1.1.2 What to do with these features?

The extracted features can be used to describe or cluster time series based on the extracted characteristics. Further, they can be used to build models that perform classification/regression tasks on the time series. Often the features give new insights into time series and their dynamics.

The tsfresh package has been used successfully in projects involving

- the prediction of the life span of machines
- the prediction of the quality of steel billets during a continuous casting process

### 1.1.3 What not to do with tsfresh?

Currently, tsfresh is not suitable

- for usage with streaming data
- for batch processing over a distributed architecture, where different time series are fragmented over different computational units
- to train models on the features (we do not want to reinvent the wheel, check out the python package scikit-learn for example)

However, some of these use cases could be implemented, if you have an application in mind, open an issue at https://github.com/blue-yonder/tsfresh/issues, or feel free to contact us.

### 1.1.4 What else is out there?

There is a matlab package called hctsa which can be used to automatically extract features from time series. It is also possible to use hctsa from within python by means of the pyopy package.

## 1.2 Quick Start

### 1.2.1 Install tsfresh

As the compiled tsfresh package is hosted on the Python Package Index (PyPI) you can easily install it with pip

```
pip install tsfresh
```

### 1.2.2 Dive in

Before boring yourself by reading the docs in detail, you can dive right into tsfresh with the following example:

We are given a data set containing robot failures as discussed in[1]. Each robot records time series from six different sensors. For each sample denoted by a different id we are going to classify if the robot reports a failure or not. From a machine learning point of view, our goal is to classify each group of time series.

To start, we load the data into python

```
from tsfresh.examples.robot_execution_failures import download_robot_execution_
↪failures, \
    load_robot_execution_failures
download_robot_execution_failures()
timeseries, y = load_robot_execution_failures()
```

and end up with a pandas.DataFrame *timeseries* having the following shape

```
print(timeseries.head())
```

|   | id | time | F_x | F_y | F_z | T_x | T_y | T_z |
|---|----|------|-----|-----|-----|-----|-----|-----|
| 0 | 1  | 0    | -1  | -1  | 63  | -3  | -1  | 0   |
| 1 | 1  | 1    | 0   | 0   | 62  | -3  | -1  | 0   |
| 2 | 1  | 2    | -1  | -1  | 61  | -3  | 0   | 0   |
| 3 | 1  | 3    | -1  | -1  | 63  | -2  | -1  | 0   |
| 4 | 1  | 4    | -1  | -1  | 63  | -3  | -1  | 0   |
| … | …  | …    | …   | …   | …   | …   | …   | …   |

The first column is the DataFrame index and has no meaning here. There are six different time series (a-f) for the different sensors. The different robots are denoted by the ids column.

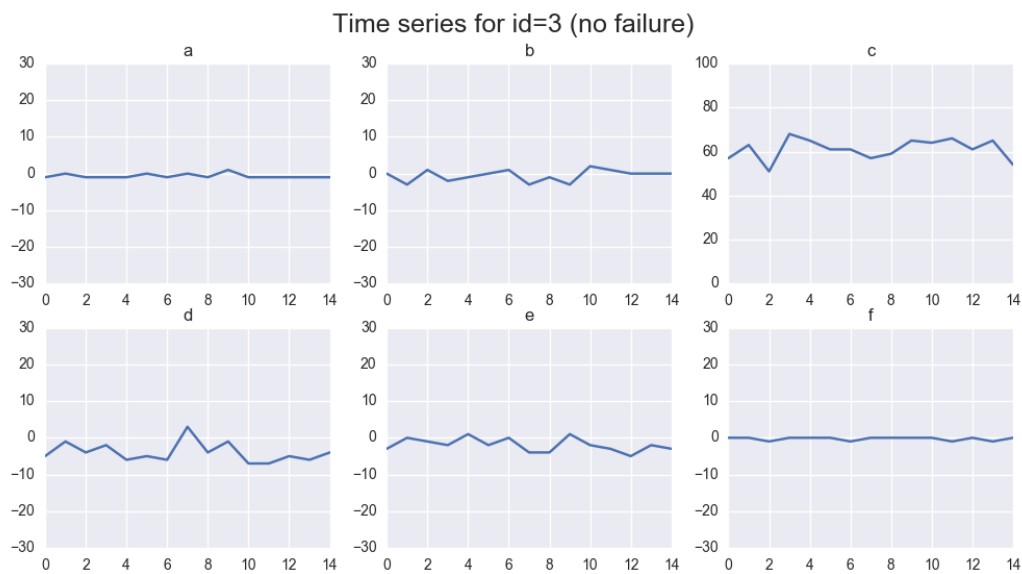On the other hand, `y` contains the information which robot *id* reported a failure and which not:

---

[1] http://archive.ics.uci.edu/ml/datasets/Robot+Execution+Failures

| 1 | 0 |
|---|---|
| 2 | 0 |
| 3 | 0 |
| 4 | 0 |
| 5 | 0 |
| ... | ... |

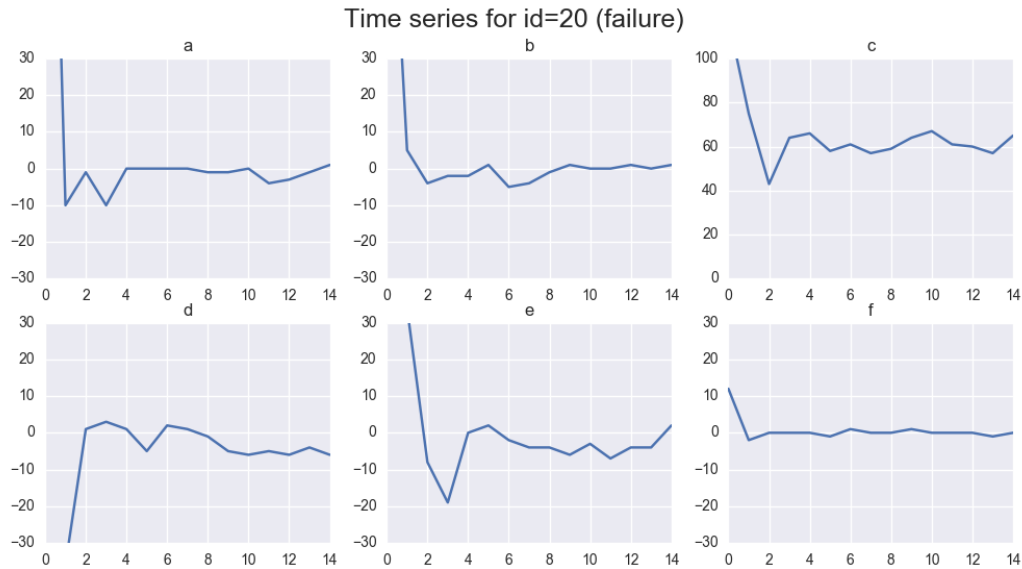Here, for the samples with ids 1 to 5 no failure was reported.

In the following we illustrate the time series of the sample id 3 reporting no failure:

```python
import matplotlib.pyplot as plt
timeseries[timeseries['id'] == 3].plot(subplots=True, sharex=True, figsize=(10,10))
plt.show()
```



And for id 20 reporting a failure:

```python
timeseries[timeseries['id'] == 21].plot(subplots=True, sharex=True, figsize=(10,10))
plt.show()
```

Time series for id=20 (failure)

You can already see some differences by eye - but for successful machine learning we have to put these differences into numbers.

For this, tsfresh comes into place. It allows us to automatically extract over 1200 features from those six different time series for each robot.

For extracting all features, we do:

```
from tsfresh import extract_features
extracted_features = extract_features(timeseries, column_id="id", column_sort="time")
```

You end up with a DataFrame *extracted_features* with all more than 1200 different extracted features. We will now remove all NaN values (that were created by feature calculators, than can not be used on the given data, e.g. because it has too low statistics) and select only the relevant features next:

```
from tsfresh import select_features
from tsfresh.utilities.dataframe_functions import impute

impute(extracted_features)
features_filtered = select_features(extracted_features, y)
```

Only around 300 features were classified as relevant enough.

Further, you can even perform the extraction, imputing and filtering at the same time with the tsfresh. extract_relevant_features() function:

```
from tsfresh import extract_relevant_features

features_filtered_direct = extract_relevant_features(timeseries, y,
                                                     column_id='id', column_sort='time
↪')
```

You can now use the features contained in the DataFrame *features_filtered* (which is equal to *features_filtered_direct*) in conjunction with *y* to train your classification model. Please see the *robot_failure_example.ipynb* Jupyter Notebook in the folder named notebook for this. In this notebook a RandomForestClassifier is trained on the extracted features.

References

# 1.3 tsfresh

## 1.3.1 tsfresh package

**Subpackages**

**tsfresh.convenience package**

**Submodules**

**tsfresh.convenience.relevant_extraction module**

tsfresh.convenience.relevant_extraction.**extract_relevant_features**(*timeseries_container,*
*y, X=None,*
*de-*
*fault_fc_parameters=None,*
*kind_to_fc_parameters=None,*
*col-*
*umn_id=None,*
*col-*
*umn_sort=None,*
*col-*
*umn_kind=None,*
*col-*
*umn_value=None,*
*show_warnings=False,*
*dis-*
*able_progressbar=False,*
*pro-*
*file=False,*
*profil-*
*ing_filename='profile.txt',*
*profil-*
*ing_sorting='cumulative',*
*test_for_binary_target_binary_featu*
*test_for_binary_target_real_feature=*
*test_for_real_target_binary_feature=*
*test_for_real_target_real_feature='k*
*fdr_level=0.05,*
*hypothe-*
*ses_independent=False,*
*n_jobs=2,*
*chunk-*
*size=None,*
*ml_task='auto')*

High level convenience function to extract time series features from *timeseries_container*. Then return feature
matrix *X* possibly augmented with relevant features with respect to target vector *y*.

For more details see the documentation of `extract_features()` and `select_features()`.

### Examples

```
>>> from tsfresh.examples import load_robot_execution_failures
>>> from tsfresh import extract_relevant_features
>>> df, y = load_robot_execution_failures()
>>> X = extract_relevant_features(df, y, column_id='id', column_sort='time')
```

Parameters

- **timeseries_container** – The pandas.DataFrame with the time series to compute the features for, or a dictionary of pandas.DataFrames. See *extract_features()*.

- **X** (*pandas.DataFrame*) – A DataFrame containing additional features

- **y** (*pandas.Series*) – The target vector

- **default_fc_parameters** (*dict*) – mapping from feature calculator names to parameters. Only those names which are keys in this dict will be calculated. See the class:*ComprehensiveFCParameters* for more information.

- **kind_to_fc_parameters** (*dict*) – mapping from kind names to objects of the same type as the ones for default_fc_parameters. If you put a kind as a key here, the fc_parameters object (which is the value), will be used instead of the default_fc_parameters.

- **column_id** (*str*) – The name of the id column to group by.

- **column_sort** (*str*) – The name of the sort column.

- **column_kind** (*str*) – The name of the column keeping record on the kind of the value.

- **column_value** (*str*) – The name for the column keeping the value itself.

- **chunksize** (*None or int*) – The size of one chunk that is submitted to the worker process for the parallelisation. Where one chunk is defined as a singular time series for one id and one kind. If you set the chunksize to 10, then it means that one task is to calculate all features for 10 time series. If it is set it to None, depending on distributor, heuristics are used to find the optimal chunksize. If you get out of memory exceptions, you can try it with the dask distributor and a smaller chunksize.

- **n_jobs** (*int*) – The number of processes to use for parallelization. If zero, no parallelization is used.

- **disable_progressbar** (*bool*) – Do not show a progressbar while doing the calculation.

- **profile** (*bool*) – Turn on profiling during feature extraction

- **profiling_sorting** (*basestring*) – How to sort the profiling results (see the documentation of the profiling package for more information)

- **profiling_filename** (*basestring*) – Where to save the profiling results.

- **test_for_binary_target_binary_feature** (*str*) – Which test to be used for binary target, binary feature (currently unused)

- **test_for_binary_target_real_feature** (*str*) – Which test to be used for binary target, real feature

- **test_for_real_target_binary_feature** (*str*) – Which test to be used for real target, binary feature (currently unused)

- **test_for_real_target_real_feature** (*str*) – Which test to be used for real target, real feature (currently unused)

- **fdr_level** (*float*) – The FDR level that should be respected, this is the theoretical expected percentage of irrelevant features among all created features.

- **hypotheses_independent** (*bool*) – Can the significance of the features be assumed to be independent? Normally, this should be set to False as the features are never independent (e.g. mean and median)

- **ml_task** (*str*) – The intended machine learning task. Either *'classification'*, *'regression'* or *'auto'*. Defaults to *'auto'*, meaning the intended task is inferred from *y*. If *y* has a boolean, integer or object dtype, the task is assumend to be classification, else regression.

> **Param** show_warnings: Show warnings during the feature extraction (needed for debugging of calculators).

> **Returns** Feature matrix X, possibly extended with relevant time series features.

## Module contents

The *convenience* submodule contains methods that allow the user to extract and filter features conveniently.

## tsfresh.examples package

## Submodules

## tsfresh.examples.driftbif_simulation module

tsfresh.examples.driftbif_simulation.**load_driftbif**(*n*, *l*, *m=2*, *classification=True*, *kappa_3=0.3*, *seed=False*)

> Simulates n time-series with l time steps each for the m-dimensional velocity of a dissipative soliton

> classification=True: target 0 means tau<=1/0.3, Dissipative Soliton with Brownian motion (purely noise driven) target 1 means tau> 1/0.3, Dissipative Soliton with Active Brownian motion (intrinsiv velocity with overlaid noise)

> classification=False: target is bifurcation parameter tau

> **Parameters**
>
> - **n** (*int*) – number of samples
> - **l** (*int*) – length of the time series
> - **m** (*int*) – number of spatial dimensions (default m=2) the dissipative soliton is propagating in
> - **classification** (*bool*) – distinguish between classification (default True) and regression target
> - **kappa_3** (*float*) – inverse bifurcation parameter (default 0.3)
> - **seed** (*float*) – random seed (default False)

> **Returns** X, y. Time series container and target vector

> **Rtype X** pandas.DataFrame

> **Rtype y** pandas.DataFrame

`tsfresh.examples.driftbif_simulation.`**`sample_tau`**(*n=10*, *kappa_3=0.3*, *ratio=0.5*, *rel_increase=0.15*)

> Return list of control parameters
>
> > **Parameters**
> >
> > - **n** (*int*) – number of samples
> >
> > - **kappa_3** (*float*) – inverse bifurcation point
> >
> > - **ratio** (*float*) – ratio (default 0.5) of samples before and beyond drift-bifurcation
> >
> > - **rel_increase** (*float*) – relative increase from bifurcation point
> >
> > **Returns** tau. List of sampled bifurcation parameter
> >
> > **Rtype tau** list

**class** `tsfresh.examples.driftbif_simulation.`**`velocity`**(*tau=3.8*, *kappa_3=0.3*, *Q=1950.0*, *R=0.0003*, *delta_t=0.05*, *seed=None*)

> Bases: `object`
>
> Simulates the velocity of a dissipative soliton (kind of self organized particle)[6]. The equilibrium velocity without noise R=0 for $tau>1.0/kappa\_3$ is $kappa\_3 sqrt\{(tau - 1.0/kappa\_3)/Q\}$. Before the drift-bifurcation $tau le 1.0/kappa\_3$ the velocity is zero.

### References

```
>>> ds = velocity(tau=3.5) # Dissipative soliton with equilibrium velocity 1.5e-3
>>> print(ds.label) # Discriminating before or beyond Drift-Bifurcation
1
>>> print(ds.deterministic) # Equilibrium velocity
0.0015191090506254991
>>> v = ds.simulate(20000) #Simulated velocity as a time series with 20000 time
↪steps being disturbed by Gaussian white noise
```

**`simulate`**(*N*, *v0=array([0., 0.])*)

> > **Parameters**
> >
> > - **N** (*int*) – number of time steps
> >
> > - **v0** (*ndarray*) – initial velocity vector
> >
> > **Returns** time series of velocity vectors with shape (N, v0.shape[0])
> >
> > **Return type** ndarray

### tsfresh.examples.har_dataset module

This module implements functions to download and load the Human Activity Recognition dataset[4]. A description of the data set can be found in[5].

---

[6] Andreas Kempa-Liehr (2013, p. 159-170) Dynamics of Dissipative Soliton Dissipative Solitons in Reaction Diffusion Systems. Springer: Berlin

[4] https://archive.ics.uci.edu/ml/datasets/Human+Activity+Recognition+Using+Smartphones

[5] Davide Anguita, Alessandro Ghio, Luca Oneto, Xavier Parra and Jorge L. Reyes-Ortiz. (2013) A Public Domain Dataset for Human Activity Recognition Using Smartphones. 21th European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning, ESANN 2013. Bruges, Belgium 24-26 April 2013.

### References

`tsfresh.examples.har_dataset.`**`download_har_dataset`**`()`
> Download human activity recognition dataset from UCI ML Repository and store it at /tsfresh/notebooks/data.

#### Examples

```
>>> from tsfresh.examples import har_dataset
>>> har_dataset.download_har_dataset()
```

`tsfresh.examples.har_dataset.`**`load_har_classes`**`()`

`tsfresh.examples.har_dataset.`**`load_har_dataset`**`()`

### tsfresh.examples.robot_execution_failures module

This module implements functions to download the Robot Execution Failures LP1 Data Set[1],[2],[3] and load it as as DataFrame.

*Important:* You need to download the data set yourself, either manually or via the function *download_robot_execution_failures()*

### References

`tsfresh.examples.robot_execution_failures.`**`download_robot_execution_failures`**`()`
> Download the Robot Execution Failures LP1 Data Set[#1] from the UCI Machine Learning Repository [#2] and store it locally.

> > **Returns**

#### Examples

```
>>> from tsfresh.examples import download_robot_execution_failures
>>> download_robot_execution_failures()
```

`tsfresh.examples.robot_execution_failures.`**`load_robot_execution_failures`**`(`*multiclass=False*`)`
> Load the Robot Execution Failures LP1 Data Set[1]. The Time series are passed as a flat DataFrame.

#### Examples

```
>>> from tsfresh.examples import load_robot_execution_failures
>>> df, y = load_robot_execution_failures()
>>> print(df.shape)
(1320, 8)
```

---

[1] https://archive.ics.uci.edu/ml/datasets/Robot+Execution+Failures

[2] Lichman, M. (2013). UCI Machine Learning Repository [https://archive.ics.uci.edu/ml]. Irvine, CA: University of California, School of Information and Computer Science.

[3] Camarinha-Matos, L.M., L. Seabra Lopes, and J. Barata (1996). Integration and Learning in Supervision of Flexible Assembly Systems. "IEEE Transactions on Robotics and Automation", 12 (2), 202-219

---

> **Parameters multiclass** (*bool*) – If True, return all target labels. The default returns only "normal" vs all other labels.
>
> **Returns** time series data as `pandas.DataFrame` and target vector as `pandas.Series`
>
> **Return type** tuple

### tsfresh.examples.test_tsfresh_baseline_dataset module

This module implements a function to download a json timeseries data set that is utilised by tests/baseline/tsfresh_features_test.py to test calculated feature names and their calculated values are consistent with the known baseline.

`tsfresh.examples.test_tsfresh_baseline_dataset.`**`download_json_dataset`**`()`
> Download the tests baseline timeseries json data set and store it at ts-fresh/examples/data/test_tsfresh_baseline_dataset/data.json.

#### Examples

```
>>> from tsfresh.examples import test_tsfresh_baseline_dataset
>>> download_json_dataset()
```

### Module contents

Module with exemplary data sets to play around with.

See for eample the *Quick Start* section on how to use them.

### tsfresh.feature_extraction package

### Submodules

### tsfresh.feature_extraction.extraction module

This module contains the main function to interact with tsfresh: extract features

`tsfresh.feature_extraction.extraction.`**`extract_features`**`(`*timeseries_container*, *default_fc_parameters=None*, *kind_to_fc_parameters=None*, *column_id=None*, *column_sort=None*, *column_kind=None*, *column_value=None*, *chunksize=None*, *n_jobs=2*, *show_warnings=False*, *disable_progressbar=False*, *impute_function=None*, *profile=False*, *profiling_filename='profile.txt'*, *profiling_sorting='cumulative'*, *distributor=None*`)`

Extract features from

- a `pandas.DataFrame` containing the different time series

or

- a dictionary of `pandas.DataFrame` each containing one type of time series

In both cases a `pandas.DataFrame` with the calculated features will be returned.

For a list of all the calculated time series features, please see the `ComprehensiveFCParameters` class, which is used to control which features with which parameters are calculated.

For a detailed explanation of the different parameters and data formats please see *Data Formats*.

### Examples

```
>>> from tsfresh.examples import load_robot_execution_failures
>>> from tsfresh import extract_features
>>> df, _ = load_robot_execution_failures()
>>> X = extract_features(df, column_id='id', column_sort='time')
```

Parameters

- **timeseries_container** (*pandas.DataFrame or dict*) – The pandas.DataFrame with the time series to compute the features for, or a dictionary of pandas.DataFrames.

- **default_fc_parameters** (*dict*) – mapping from feature calculator names to parameters. Only those names which are keys in this dict will be calculated. See the class:*ComprehensiveFCParameters* for more information.

- **kind_to_fc_parameters** (*dict*) – mapping from kind names to objects of the same type as the ones for default_fc_parameters. If you put a kind as a key here, the fc_parameters object (which is the value), will be used instead of the default_fc_parameters.

- **column_id** (*str*) – The name of the id column to group by.

- **column_sort** (*str*) – The name of the sort column.

- **column_kind** (*str*) – The name of the column keeping record on the kind of the value.

- **column_value** (*str*) – The name for the column keeping the value itself.

- **n_jobs** (*int*) – The number of processes to use for parallelization. If zero, no parallelization is used.

- **chunksize** (*None or int*) – The size of one chunk that is submitted to the worker process for the parallelisation. Where one chunk is defined as a singular time series for one id and one kind. If you set the chunksize to 10, then it means that one task is to calculate all features for 10 time series. If it is set it to None, depending on distributor, heuristics are used to find the optimal chunksize. If you get out of memory exceptions, you can try it with the dask distributor and a smaller chunksize.

- **disable_progressbar** (*bool*) – Do not show a progressbar while doing the calculation.

- **impute_function** (*None or callable*) – None, if no imputing should happen or the function to call for imputing.

- **profile** (*bool*) – Turn on profiling during feature extraction

- **profiling_sorting** (*basestring*) – How to sort the profiling results (see the documentation of the profiling package for more information)

- **profiling_filename** (*basestring*) – Where to save the profiling results.

- **distributor** (*class*) – Advanced parameter: set this to a class name that you want to use as a distributor. See the utilities/distribution.py for more information. Leave to None, if you want TSFresh to choose the best distributor.

**Param** show_warnings: Show warnings during the feature extraction (needed for debugging of calculators).

**Returns** The (maybe imputed) DataFrame containing extracted features.

**Return type** pandas.DataFrame

## tsfresh.feature_extraction.feature_calculators module

This module contains the feature calculators that take time series as input and calculate the values of the feature. There are two types of features:

1. feature calculators which calculate a single number (simple)

2. feature calculators which calculate a bunch of features for a list of parameters at once, to use e.g. cached results (combiner). They return a list of (key, value) pairs for each input parameter.

They are specified using the "fctype" parameter of each feature calculator, which is added using the set_property function. Only functions in this python module, which have a parameter called "fctype" are seen by tsfresh as a feature calculator. Others will not be calculated.

tsfresh.feature_extraction.feature_calculators.**abs_energy**(*x*)
    Returns the absolute energy of the time series which is the sum over the squared values

$$E = \sum_{i=1,\dots,n} x_i^2$$

    **Parameters x** (*pandas.Series*) – the time series to calculate the feature of

    **Returns** the value of this feature

    **Return type** float

*This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**absolute_sum_of_changes**(*x*)
    Returns the sum over the absolute value of consecutive changes in the series x

$$\sum_{i=1,\dots,n-1} \mid x_{i+1} - x_i \mid$$

    **Parameters x** (*pandas.Series*) – the time series to calculate the feature of

    **Returns** the value of this feature

    **Return type** float

*This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**agg_autocorrelation**(*x*, *param*)
    Calculates the value of an aggregation function f_agg (e.g. var or mean) of the autocorrelation (Compare to

http://en.wikipedia.org/wiki/Autocorrelation#Estimation), taken over different all possible lags (1 to length of x)

$$\frac{1}{n-1} \sum_{l=1,\dots,n} \frac{1}{(n-l)\sigma^2} \sum_{t=1}^{n-l} (X_t - \mu)(X_{t+l} - \mu)$$

where $n$ is the length of the time series $X_i$, $\sigma^2$ its variance and $\mu$ its mean.

> **Parameters**
>
> - **x** (*pandas.Series*) – the time series to calculate the feature of
> - **param** (*list*) – contains dictionaries {"attr": x} with x str, name of a numpy function (e.g. mean, var, std, median), the name of the aggregator function that is applied to the autocorrelations
>
> **Returns** the value of this feature
>
> **Return type** float

*This function is of type: combiner*

tsfresh.feature_extraction.feature_calculators.**agg_linear_trend**(*x*, *param*)
Calculates a linear least-squares regression for values of the time series that were aggregated over chunks versus the sequence from 0 up to the number of chunks minus one.

This feature assumes the signal to be uniformly sampled. It will not use the time stamps to fit the model.

The parameters attr controls which of the characteristics are returned. Possible extracted attributes are "pvalue", "rvalue", "intercept", "slope", "stderr", see the documentation of linregress for more information.

The chunksize is regulated by "chunk_len". It specifies how many time series values are in each chunk.

Further, the aggregation function is controlled by "f_agg", which can use "max", "min" or , "mean", "median"

> **Parameters**
>
> - **x** (*pandas.Series*) – the time series to calculate the feature of
> - **param** (*list*) – contains dictionaries {"attr": x, "chunk_len": l, "f_agg": f} with x, f an string and l an int
>
> **Returns** the different feature values
>
> **Return type** pandas.Series

*This function is of type: combiner*

tsfresh.feature_extraction.feature_calculators.**approximate_entropy**(*x*, *m*, *r*)
Implements a vectorized Approximate entropy algorithm.

> https://en.wikipedia.org/wiki/Approximate_entropy

For short time-series this method is highly dependent on the parameters, but should be stable for N > 2000, see:

> Yentes et al. (2012) - *The Appropriate Use of Approximate Entropy and Sample Entropy with Short Data Sets*

Other shortcomings and alternatives discussed in:

> Richman & Moorman (2000) - *Physiological time-series analysis using approximate entropy and sample entropy*

> **Parameters**
>
> - **x** (*pandas.Series*) – the time series to calculate the feature of

- **m** (*int*) – Length of compared run of data

- **r** (*float*) – Filtering level, must be positive

**Returns** Approximate entropy

**Return type** float

*This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**ar_coefficient**(*x*, *param*)

This feature calculator fits the unconditional maximum likelihood of an autoregressive AR(k) process. The k parameter is the maximum lag of the process

$$X_t = \varphi_0 + \sum_{i=1}^{k} \varphi_i X_{t-i} + \varepsilon_t$$

For the configurations from param which should contain the maxlag "k" and such an AR process is calculated. Then the coefficients $\varphi_i$ whose index $i$ contained from "coeff" are returned.

**Parameters**

- **x** (*pandas.Series*) – the time series to calculate the feature of

- **param** (*list*) – contains dictionaries {"coeff": x, "k": y} with x,y int

**Return x** the different feature values

**Return type** pandas.Series

*This function is of type: combiner*

tsfresh.feature_extraction.feature_calculators.**augmented_dickey_fuller**(*x*,

*param*)

The Augmented Dickey-Fuller test is a hypothesis test which checks whether a unit root is present in a time series sample. This feature calculator returns the value of the respective test statistic.

See the statsmodels implementation for references and more details.

**Parameters**

- **x** (*pandas.Series*) – the time series to calculate the feature of

- **param** (*list*) – contains dictionaries {"attr": x} with x str, either "teststat", "pvalue" or "usedlag"

**Returns** the value of this feature

**Return type** float

*This function is of type: combiner*

tsfresh.feature_extraction.feature_calculators.**autocorrelation**(*x*, *lag*)

Calculates the autocorrelation of the specified lag, according to the formula [1]

$$\frac{1}{(n-l)\sigma^2} \sum_{t=1}^{n-l} (X_t - \mu)(X_{t+l} - \mu)$$

where $n$ is the length of the time series $X_i$, $\sigma^2$ its variance and $\mu$ its mean. $l$ denotes the lag.

**References**

[1] https://en.wikipedia.org/wiki/Autocorrelation#Estimation

> **Parameters**
>
> - **x** (*pandas.Series*) – the time series to calculate the feature of
>
> - **lag** (*int*) – the lag
>
> **Returns** the value of this feature
>
> **Return type** float

*This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**binned_entropy**(*x*, *max_bins*)
First bins the values of x into max_bins equidistant bins. Then calculates the value of

$$- \sum_{k=0}^{min(max\_bins, len(x))} p_k log(p_k) \cdot \mathbf{1}_{(p_k > 0)}$$

where $p_k$ is the percentage of samples in bin $k$.

> **Parameters**
>
> - **x** (*pandas.Series*) – the time series to calculate the feature of
>
> - **max_bins** (*int*) – the maximal number of bins
>
> **Returns** the value of this feature
>
> **Return type** float

*This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**c3**(*x*, *lag*)
This function calculates the value of

$$\frac{1}{n - 2lag} \sum_{i=0}^{n-2lag} x_{i+2 \cdot lag}^2 \cdot x_{i+lag} \cdot x_i$$

which is

$$\mathbb{E}[L^2(X)^2 \cdot L(X) \cdot X]$$

where $\mathbb{E}$ is the mean and $L$ is the lag operator. It was proposed in [1] as a measure of non linearity in the time series.

**References**

[1] Schreiber, T. and Schmitz, A. (1997).
Discrimination power of measures for nonlinearity in a time series
PHYSICAL REVIEW E, VOLUME 55, NUMBER 5

> **Parameters**
>
> - **x** (*pandas.Series*) – the time series to calculate the feature of
>
> - **lag** (*int*) – the lag that should be used in the calculation of the feature

> **Returns** the value of this feature
>
> **Return type** float

*This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**change_quantiles**(*x*, *ql*, *qh*, *isabs*,
*f_agg*)

> First fixes a corridor given by the quantiles ql and qh of the distribution of x. Then calculates the average, absolute value of consecutive changes of the series x inside this corridor.
>
> Think about selecting a corridor on the y-Axis and only calculating the mean of the absolute change of the time series inside this corridor.
>
> > **Parameters**
> >
> > - **x** (*pandas.Series*) – the time series to calculate the feature of
> > - **ql** (*float*) – the lower quantile of the corridor
> > - **qh** (*float*) – the higher quantile of the corridor
> > - **isabs** (*bool*) – should the absolute differences be taken?
> > - **f_agg** (*str, name of a numpy function (e.g. mean, var, std, median)*) – the aggregator function that is applied to the differences in the bin
>
> **Returns** the value of this feature
>
> **Return type** float

*This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**cid_ce**(*x*, *normalize*)

> This function calculator is an estimate for a time series complexity [1] (A more complex time series has more peaks, valleys etc.). It calculates the value of
>
> $$\sqrt{\sum_{i=0}^{n-2lag} (x_i - x_{i+1})^2}$$

### References

[1] Batista, Gustavo EAPA, et al (2014).
CID: an efficient complexity-invariant distance for time series.
Data Mining and Knowledge Difscovery 28.3 (2014): 634-669.

> > **Parameters**
> >
> > - **x** (*pandas.Series*) – the time series to calculate the feature of
> > - **normalize** (*bool*) – should the time series be z-transformed?
>
> **Returns** the value of this feature
>
> **Return type** float

*This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**count_above_mean**(*x*)

> Returns the number of values in x that are higher than the mean of x

> **Parameters x** (*pandas.Series*) – the time series to calculate the feature of
>
> **Returns** the value of this feature
>
> **Return type** float

*This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**count_below_mean**(*x*)

Returns the number of values in x that are lower than the mean of x

> **Parameters x** (*pandas.Series*) – the time series to calculate the feature of
>
> **Returns** the value of this feature
>
> **Return type** float

*This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**cwt_coefficients**(*x*, *param*)

Calculates a Continuous wavelet transform for the Ricker wavelet, also known as the "Mexican hat wavelet" which is defined by

$$\frac{2}{\sqrt{3a}\pi^{\frac{1}{4}}}(1 - \frac{x^2}{a^2})exp(-\frac{x^2}{2a^2})$$

where $a$ is the width parameter of the wavelet function.

This feature calculator takes three different parameter: widths, coeff and w. The feature calculater takes all the different widths arrays and then calculates the cwt one time for each different width array. Then the values for the different coefficient for coeff and width w are returned. (For each dic in param one feature is returned)

> **Parameters**
>
> - **x** (*pandas.Series*) – the time series to calculate the feature of
> - **param** (*list*) – contains dictionaries {"widths":x, "coeff": y, "w": z} with x array of int and y,z int
>
> **Returns** the different feature values
>
> **Return type** pandas.Series

*This function is of type: combiner*

tsfresh.feature_extraction.feature_calculators.**energy_ratio_by_chunks**(*x*, *param*)

Calculates the sum of squares of chunk i out of N chunks expressed as a ratio with the sum of squares over the whole series

Takes as input parameters the number num_segments of segments to divide the series into and segment_focus which is the segment number (starting at zero) to return a feature on.

Note that the answer for num_segments=1 is a trivial "1" but we handle this scenario in case somebody calls it. Sum of the ratios should be 1.0.

Returns an error for N <= 0

> **Parameters**
>
> - **x** (*pandas.Series*) – the time series to calculate the feature of
> - **param** – contains dictionaries {"num_segments": N, "segment_focus": i} with N, i both ints
>
> **Returns** the feature values

> **Return type** list of tuples (index, data)

*This function is of type: combiner*

tsfresh.feature_extraction.feature_calculators.**fft_aggregated**(*x*, *param*)
> Returns the spectral centroid (mean), variance, skew, and kurtosis of the absolute fourier transform spectrum.
>
> > **Parameters**
> >
> > - **x** (*pandas.Series*) – the time series to calculate the feature of
> >
> > - **param** (*list*) – contains dictionaries {"aggtype": s} where s str and in ["centroid", "variance", "skew", "kurtosis"]
> >
> > **Returns** the different feature values
> >
> > **Return type** pandas.Series

*This function is of type: combiner*

tsfresh.feature_extraction.feature_calculators.**fft_coefficient**(*x*, *param*)
> Calculates the fourier coefficients of the one-dimensional discrete Fourier Transform for real input by fast fourier transformation algorithm

$$A_k = \sum_{m=0}^{n-1} a_m \exp\left\{-2\pi i \frac{mk}{n}\right\}, \qquad k = 0, \dots, n-1.$$

> The resulting coefficients will be complex, this feature calculator can return the real part (attr=="real"), the imaginary part (attr=="imag), the absolute value (attr=""abs) and the angle in degrees (attr=="angle).
>
> > **Parameters**
> >
> > - **x** (*pandas.Series*) – the time series to calculate the feature of
> >
> > - **param** (*list*) – contains dictionaries {"coeff": x, "attr": s} with x int and x >= 0, s str and in ["real", "imag", "abs", "angle"]
> >
> > **Returns** the different feature values
> >
> > **Return type** pandas.Series

*This function is of type: combiner*

tsfresh.feature_extraction.feature_calculators.**first_location_of_maximum**(*x*)
> Returns the first location of the maximum value of x. The position is calculated relatively to the length of x.
>
> > **Parameters** **x** (*pandas.Series*) – the time series to calculate the feature of
> >
> > **Returns** the value of this feature
> >
> > **Return type** float

*This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**first_location_of_minimum**(*x*)
> Returns the first location of the minimal value of x. The position is calculated relatively to the length of x.
>
> > **Parameters** **x** (*pandas.Series*) – the time series to calculate the feature of
> >
> > **Returns** the value of this feature
> >
> > **Return type** float

*This function is of type: simple*

`tsfresh.feature_extraction.feature_calculators.`**`friedrich_coefficients`**(*x*,

*param*)

Coefficients of polynomial $h(x)$, which has been fitted to the deterministic dynamics of Langevin model

$$\dot{x}(t) = h(x(t)) + \mathcal{N}(0, R)$$

as described by [1].

For short time-series this method is highly dependent on the parameters.

### References

[1] Friedrich et al. (2000): Physics Letters A 271, p. 217-222
*Extracting model equations from experimental data*

> **Parameters**
>
> - **x** (*pandas.Series*) – the time series to calculate the feature of
>
> - **c** (*str*) – the time series name
>
> - **param** (*list*) – contains dictionaries {"coeff": x} with x int and x >= 0
>
> **Returns**  the different feature values
>
> **Return type**  pandas.Series

*This function is of type: combiner*

`tsfresh.feature_extraction.feature_calculators.`**`has_duplicate`**(*x*)

Checks if any value in x occurs more than once

> **Parameters x** (*pandas.Series*) – the time series to calculate the feature of
>
> **Returns**  the value of this feature
>
> **Return type**  bool

*This function is of type: simple*

`tsfresh.feature_extraction.feature_calculators.`**`has_duplicate_max`**(*x*)

Checks if the maximum value of x is observed more than once

> **Parameters x** (*pandas.Series*) – the time series to calculate the feature of
>
> **Returns**  the value of this feature
>
> **Return type**  bool

*This function is of type: simple*

`tsfresh.feature_extraction.feature_calculators.`**`has_duplicate_min`**(*x*)

Checks if the minimal value of x is observed more than once

> **Parameters x** (*pandas.Series*) – the time series to calculate the feature of
>
> **Returns**  the value of this feature
>
> **Return type**  bool

*This function is of type: simple*

`tsfresh.feature_extraction.feature_calculators.`**`index_mass_quantile`**(*x, param*)

Those apply features calculate the relative index i where q% of the mass of the time series x lie left of i. For example for q = 50% this feature calculator will return the mass center of the time series

**Parameters**

- **x** (*pandas.Series*) – the time series to calculate the feature of

- **param** (*list*) – contains dictionaries {"q": x} with x float

**Returns** the different feature values

**Return type** pandas.Series

*This function is of type: combiner*

`tsfresh.feature_extraction.feature_calculators.`**`kurtosis`**(*x*)

Returns the kurtosis of x (calculated with the adjusted Fisher-Pearson standardized moment coefficient G2).

**Parameters x** (*pandas.Series*) – the time series to calculate the feature of

**Returns** the value of this feature

**Return type** float

*This function is of type: simple*

`tsfresh.feature_extraction.feature_calculators.`**`large_standard_deviation`**(*x,*
*r*)

Boolean variable denoting if the standard dev of x is higher than 'r' times the range = difference between max and min of x. Hence it checks if

$$std(x) > r * (max(X) - min(X))$$

According to a rule of the thumb, the standard deviation should be a forth of the range of the values.

**Parameters**

- **x** (*pandas.Series*) – the time series to calculate the feature of

- **r** (*float*) – the percentage of the range to compare with

**Returns** the value of this feature

**Return type** bool

*This function is of type: simple*

`tsfresh.feature_extraction.feature_calculators.`**`last_location_of_maximum`**(*x*)

Returns the relative last location of the maximum value of x. The position is calculated relatively to the length of x.

**Parameters x** (*pandas.Series*) – the time series to calculate the feature of

**Returns** the value of this feature

**Return type** float

*This function is of type: simple*

`tsfresh.feature_extraction.feature_calculators.`**`last_location_of_minimum`**(*x*)

Returns the last location of the minimal value of x. The position is calculated relatively to the length of x.

**Parameters x** (*pandas.Series*) – the time series to calculate the feature of

**Returns** the value of this feature

**Return type** float

*This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**length**(*x*)

Returns the length of x

> **Parameters** **x** (*pandas.Series*) – the time series to calculate the feature of
>
> **Returns** the value of this feature
>
> **Return type** int

*This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**linear_trend**(*x*, *param*)

Calculate a linear least-squares regression for the values of the time series versus the sequence from 0 to length of the time series minus one. This feature assumes the signal to be uniformly sampled. It will not use the time stamps to fit the model. The parameters control which of the characteristics are returned.

Possible extracted attributes are "pvalue", "rvalue", "intercept", "slope", "stderr", see the documentation of linregress for more information.

> **Parameters**
>
> - **x** (*pandas.Series*) – the time series to calculate the feature of
>
> - **param** (*list*) – contains dictionaries {"attr": x} with x an string, the attribute name of the regression model
>
> **Returns** the different feature values
>
> **Return type** pandas.Series

*This function is of type: combiner*

tsfresh.feature_extraction.feature_calculators.**longest_strike_above_mean**(*x*)

Returns the length of the longest consecutive subsequence in x that is bigger than the mean of x

> **Parameters** **x** (*pandas.Series*) – the time series to calculate the feature of
>
> **Returns** the value of this feature
>
> **Return type** float

*This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**longest_strike_below_mean**(*x*)

Returns the length of the longest consecutive subsequence in x that is smaller than the mean of x

> **Parameters** **x** (*pandas.Series*) – the time series to calculate the feature of
>
> **Returns** the value of this feature
>
> **Return type** float

*This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**max_langevin_fixed_point**(*x*, *r*, *m*)

Largest fixed point of dynamics :math:argmax_x {h(x)=0}' estimated from polynomial $h(x)$, which has been fitted to the deterministic dynamics of Langevin model

$$\dot{(x)}(t) = h(x(t)) + R(N)(0,1)$$

as described by

Friedrich et al. (2000): Physics Letters A 271, p. 217-222 *Extracting model equations from experimental data*

For short time-series this method is highly dependent on the parameters.

> **Parameters**
>
> > • **x** (*pandas.Series*) – the time series to calculate the feature of
> >
> > • **m** (*int*) – order of polynom to fit for estimating fixed points of dynamics
> >
> > • **r** (*float*) – number of quantils to use for averaging
>
> **Returns** Largest fixed point of deterministic dynamics
>
> **Return type** float

*This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**maximum**(*x*)
> Calculates the highest value of the time series x.
>
> > **Parameters x** (*pandas.Series*) – the time series to calculate the feature of
> >
> > **Returns** the value of this feature
> >
> > **Return type** float

*This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**mean**(*x*)
> Returns the mean of x
>
> > **Parameters x** (*pandas.Series*) – the time series to calculate the feature of
> >
> > **Returns** the value of this feature
> >
> > **Return type** float

*This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**mean_abs_change**(*x*)
> Returns the mean over the absolute differences between subsequent time series values which is
>
> $$\frac{1}{n} \sum_{i=1,\dots,n-1} |x_{i+1} - x_i|$$
>
> > **Parameters x** (*pandas.Series*) – the time series to calculate the feature of
> >
> > **Returns** the value of this feature
> >
> > **Return type** float

*This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**mean_change**(*x*)
> Returns the mean over the absolute differences between subsequent time series values which is
>
> $$\frac{1}{n} \sum_{i=1,\dots,n-1} x_{i+1} - x_i$$
>
> > **Parameters x** (*pandas.Series*) – the time series to calculate the feature of
> >
> > **Returns** the value of this feature
> >
> > **Return type** float

*This function is of type: simple*

`tsfresh.feature_extraction.feature_calculators.`**`mean_second_derivative_central`**`(x)`

Returns the mean value of a central approximation of the second derivative

$$\frac{1}{n} \sum_{i=1,\ldots,n-1} \frac{1}{2}(x_{i+2} - 2 \cdot x_{i+1} + x_i)$$

> **Parameters** **x** (`pandas.Series`) – the time series to calculate the feature of
>
> **Returns** the value of this feature
>
> **Return type** float

*This function is of type: simple*

`tsfresh.feature_extraction.feature_calculators.`**`median`**`(x)`

Returns the median of x

> **Parameters** **x** (`pandas.Series`) – the time series to calculate the feature of
>
> **Returns** the value of this feature
>
> **Return type** float

*This function is of type: simple*

`tsfresh.feature_extraction.feature_calculators.`**`minimum`**`(x)`

Calculates the lowest value of the time series x.

> **Parameters** **x** (`pandas.Series`) – the time series to calculate the feature of
>
> **Returns** the value of this feature
>
> **Return type** float

*This function is of type: simple*

`tsfresh.feature_extraction.feature_calculators.`**`number_crossing_m`**`(x, m)`

Calculates the number of crossings of x on m. A crossing is defined as two sequential values where the first value is lower than m and the next is greater, or vice-versa. If you set m to zero, you will get the number of zero crossings.

> **Parameters**
>
> - **x** (`pandas.Series`) – the time series to calculate the feature of
> - **m** (`float`) – the threshold for the crossing
>
> **Returns** the value of this feature
>
> **Return type** int

*This function is of type: simple*

`tsfresh.feature_extraction.feature_calculators.`**`number_cwt_peaks`**`(x, n)`

This feature calculator searches for different peaks in x. To do so, x is smoothed by a ricker wavelet and for widths ranging from 1 to n. This feature calculator returns the number of peaks that occur at enough width scales and with sufficiently high Signal-to-Noise-Ratio (SNR)

> **Parameters**
>
> - **x** (`pandas.Series`) – the time series to calculate the feature of
> - **n** (`int`) – maximum width to consider
>
> **Returns** the value of this feature
>
> **Return type** int

*This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**number_peaks**(*x, n*)

> Calculates the number of peaks of at least support n in the time series x. A peak of support n is defined as a subsequence of x where a value occurs, which is bigger than its n neighbours to the left and to the right.
>
> Hence in the sequence

```
>>> x = [3, 0, 0, 4, 0, 0, 13]
```

> 4 is a peak of support 1 and 2 because in the subsequences

```
>>> [0, 4, 0]
>>> [0, 0, 4, 0, 0]
```

> 4 is still the highest value. Here, 4 is not a peak of support 3 because 13 is the 3th neighbour to the right of 4 and its bigger than 4.
>
> > **Parameters**
> >
> > > - **x** (*pandas.Series*) – the time series to calculate the feature of
> > >
> > > - **n** (*int*) – the support of the peak
> >
> > **Returns**  the value of this feature
> >
> > **Return type**  float

*This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**partial_autocorrelation**(*x,*
> > > > > > > > > > > > > > > > > > > > > > > > > > > > > > > > > > > > > *param*)

> Calculates the value of the partial autocorrelation function at the given lag. The lag $k$ partial autocorrelation of a time series $\{x_t, t = 1 \ldots T\}$ equals the partial correlation of $x_t$ and $x_{t-k}$, adjusted for the intermediate variables $\{x_{t-1}, \ldots, x_{t-k+1}\}$ ([1]). Following [2], it can be defined as
>
> $$\alpha_k = \frac{Cov(x_t, x_{t-k}|x_{t-1}, \ldots, x_{t-k+1})}{\sqrt{Var(x_t|x_{t-1}, \ldots, x_{t-k+1})Var(x_{t-k}|x_{t-1}, \ldots, x_{t-k+1})}}$$
>
> with (a) $x_t = f(x_{t-1}, \ldots, x_{t-k+1})$ and (b) $x_{t-k} = f(x_{t-1}, \ldots, x_{t-k+1})$ being AR(k-1) models that can be fitted by OLS. Be aware that in (a), the regression is done on past values to predict $x_t$ whereas in (b), future values are used to calculate the past value $x_{t-k}$. It is said in [1] that "for an AR(p), the partial autocorrelations [ $\alpha_k$ ] will be nonzero for *k<=p* and zero for *k>p*." With this property, it is used to determine the lag of an AR-Process.

### References

[1] Box, G. E., Jenkins, G. M., Reinsel, G. C., & Ljung, G. M. (2015).

Time series analysis: forecasting and control. John Wiley & Sons.

[2] https://onlinecourses.science.psu.edu/stat510/node/62

> > **Parameters**
> >
> > > - **x** (*pandas.Series*) – the time series to calculate the feature of
> > >
> > > - **param** (*list*) – contains dictionaries {"lag": val} with int val indicating the lag to be returned
> >
> > **Returns**  the value of this feature

**Return type** float

*This function is of type: combiner*

tsfresh.feature_extraction.feature_calculators.**percentage_of_reoccurring_datapoints_to_all_**
Returns the percentage of unique values, that are present in the time series more than once.

len(different values occurring more than once) / len(different values)

This means the percentage is normalized to the number of unique values, in contrast to the percentage_of_reoccurring_values_to_all_values.

**Parameters** **x** (*pandas.Series*) – the time series to calculate the feature of

**Returns** the value of this feature

**Return type** float

*This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**percentage_of_reoccurring_values_to_all_valu**
Returns the ratio of unique values, that are present in the time series more than once.

# of data points occurring more than once / # of all data points

This means the ratio is normalized to the number of data points in the time series, in contrast to the percentage_of_reoccurring_datapoints_to_all_datapoints.

**Parameters** **x** (*pandas.Series*) – the time series to calculate the feature of

**Returns** the value of this feature

**Return type** float

*This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**quantile**(*x*, *q*)
Calculates the q quantile of x. This is the value of x greater than q% of the ordered values from x.

**Parameters**

- **x** (*pandas.Series*) – the time series to calculate the feature of

- **q** (*float*) – the quantile to calculate

**Returns** the value of this feature

**Return type** float

*This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**range_count**(*x*, *min*, *max*)
Count observed values within the interval [min, max).

**Parameters**

- **x** (*pandas.Series*) – the time series to calculate the feature of

- **min** (*int or float*) – the inclusive lower bound of the range

- **max** (*int or float*) – the exclusive upper bound of the range

**Returns** the count of values within the range

**Return type** int

*This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**ratio_beyond_r_sigma**(*x*, *r*)

    Ratio of values that are more than r*std(x) (so r sigma) away from the mean of x.

        **Parameters** **x** (*iterable*) – the time series to calculate the feature of

        **Returns** the value of this feature

        **Return type** float

    *This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**ratio_value_number_to_time_series_length**(*x*)

    Returns a factor which is 1 if all values in the time series occur only once, and below one if this is not the case. In principle, it just returns

        # unique values / # values

        **Parameters** **x** (*pandas.Series*) – the time series to calculate the feature of

        **Returns** the value of this feature

        **Return type** float

    *This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**sample_entropy**(*x*)

    Calculate and return sample entropy of x.

### References

    [1] http://en.wikipedia.org/wiki/Sample_Entropy

    [2] https://www.ncbi.nlm.nih.gov/pubmed/10843903?dopt=Abstract

        **Parameters** **x** (*pandas.Series*) – the time series to calculate the feature of

        **Returns** the value of this feature

        **Return type** float

    *This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**set_property**(*key*, *value*)

    This method returns a decorator that sets the property key of the function to value

tsfresh.feature_extraction.feature_calculators.**skewness**(*x*)

    Returns the sample skewness of x (calculated with the adjusted Fisher-Pearson standardized moment coefficient G1).

        **Parameters** **x** (*pandas.Series*) – the time series to calculate the feature of

        **Returns** the value of this feature

        **Return type** float

    *This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**spkt_welch_density**(*x*, *param*)

    This feature calculator estimates the cross power spectral density of the time series x at different frequencies. To do so, the time series is first shifted from the time domain to the frequency domain.

    The feature calculators returns the power spectrum of the different frequencies.

**Parameters**

- **x** (*pandas.Series*) – the time series to calculate the feature of
- **param** (*list*) – contains dictionaries {"coeff": x} with x int

**Returns** the different feature values

**Return type** pandas.Series

*This function is of type: combiner*

tsfresh.feature_extraction.feature_calculators.**standard_deviation**(*x*)

Returns the standard deviation of x

    **Parameters** **x** (*pandas.Series*) – the time series to calculate the feature of

    **Returns** the value of this feature

    **Return type** float

*This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**sum_of_reoccurring_data_points**(*x*)

Returns the sum of all data points, that are present in the time series more than once.

    **Parameters** **x** (*pandas.Series*) – the time series to calculate the feature of

    **Returns** the value of this feature

    **Return type** float

*This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**sum_of_reoccurring_values**(*x*)

Returns the sum of all values, that are present in the time series more than once.

    **Parameters** **x** (*pandas.Series*) – the time series to calculate the feature of

    **Returns** the value of this feature

    **Return type** float

*This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**sum_values**(*x*)

Calculates the sum over the time series values

    **Parameters** **x** (*pandas.Series*) – the time series to calculate the feature of

    **Returns** the value of this feature

    **Return type** bool

*This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**symmetry_looking**(*x*, *param*)

Boolean variable denoting if the distribution of x *looks symmetric*. This is the case if

$$|mean(X) - median(X)| < r * (max(X) - min(X))$$

**Parameters**

- **x** (*pandas.Series*) – the time series to calculate the feature of
- **r** (*float*) – the percentage of the range to compare with

**Returns** the value of this feature

> **Return type** bool

*This function is of type: combiner*

tsfresh.feature_extraction.feature_calculators.**time_reversal_asymmetry_statistic**(*x*,
*lag*)

> This function calculates the value of

$$\frac{1}{n-2lag} \sum_{i=0}^{n-2lag} x_{i+2\cdot lag}^2 \cdot x_{i+lag} - x_{i+lag} \cdot x_i^2$$

> which is

$$\mathbb{E}[L^2(X)^2 \cdot L(X) - L(X) \cdot X^2]$$

> where $\mathbb{E}$ is the mean and $L$ is the lag operator. It was proposed in [1] as a promising feature to extract from time series.

### References

[1] Fulcher, B.D., Jones, N.S. (2014).
Highly comparative feature-based time-series classification.
Knowledge and Data Engineering, IEEE Transactions on 26, 3026–3037.

> **Parameters**
>
> - **x** (*pandas.Series*) – the time series to calculate the feature of
> - **lag** (*int*) – the lag that should be used in the calculation of the feature
>
> **Returns** the value of this feature
>
> **Return type** float

*This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**value_count**(*x*, *value*)
Count occurrences of *value* in time series x.

> **Parameters**
>
> - **x** (*pandas.Series*) – the time series to calculate the feature of
> - **value** (*int or float*) – the value to be counted
>
> **Returns** the count
>
> **Return type** int

*This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**variance**(*x*)
Returns the variance of x

> **Parameters** **x** (*pandas.Series*) – the time series to calculate the feature of
>
> **Returns** the value of this feature
>
> **Return type** float

*This function is of type: simple*

tsfresh.feature_extraction.feature_calculators.**variance_larger_than_standard_deviation**(*x*)

Boolean variable denoting if the variance of x is greater than its standard deviation. Is equal to variance of x being larger than 1

> **Parameters x** (*pandas.Series*) – the time series to calculate the feature of
>
> **Returns** the value of this feature
>
> **Return type** bool

*This function is of type: simple*

## tsfresh.feature_extraction.settings module

This file contains methods/objects for controlling which features will be extracted when calling extract_features. For the naming of the features, see *Feature Calculation*.

**class** tsfresh.feature_extraction.settings.**ComprehensiveFCParameters**

Bases: `dict`

Create a new ComprehensiveFCParameters instance. You have to pass this instance to the extract_feature instance.

It is basically a dictionary (and also based on one), which is a mapping from string (the same names that are in the feature_calculators.py file) to a list of dictionary of parameters, which should be used when the function with this name is called.

Only those strings (function names), that are keys in this dictionary, will be later used to extract features - so whenever you delete a key from this dict, you disable the calculation of this feature.

You can use the settings object with

```
>>> from tsfresh.feature_extraction import extract_features,
↪ComprehensiveFCParameters
>>> extract_features(df, default_fc_parameters=ComprehensiveFCParameters())
```

to extract all features (which is the default nevertheless) or you change the ComprehensiveFCParameters object to other types (see below).

**class** tsfresh.feature_extraction.settings.**EfficientFCParameters**

Bases: *tsfresh.feature_extraction.settings.ComprehensiveFCParameters*

This class is a child class of the ComprehensiveFCParameters class and has the same functionality as its base class.

The only difference is, that the features with high computational costs are not calculated. Those are denoted by the attribute "high_comp_cost".

You should use this object when calling the extract function, like so:

```
>>> from tsfresh.feature_extraction import extract_features, EfficientFCParameters
>>> extract_features(df, default_fc_parameters=EfficientFCParameters())
```

**class** tsfresh.feature_extraction.settings.**MinimalFCParameters**

Bases: *tsfresh.feature_extraction.settings.ComprehensiveFCParameters*

This class is a child class of the ComprehensiveFCParameters class and has the same functionality as its base class. The only difference is, that most of the feature calculators are disabled and only a small subset of calculators will be calculated at all. Those are donated by an attribute called "minimal".

Use this class for quick tests of your setup before calculating all features which could take some time depending of your data set size.

You should use this object when calling the extract function, like so:

```
>>> from tsfresh.feature_extraction import extract_features, MinimalFCParameters
>>> extract_features(df, default_fc_parameters=MinimalFCParameters())
```

tsfresh.feature_extraction.settings.**from_columns**(*columns*, *columns_to_ignore=[]*)

Creates a mapping from kind names to fc_parameters objects (which are itself mappings from feature calculators to settings) to extract only the features contained in the columns. To do so, for every feature name in columns this method

1. split the column name into col, feature, params part

2. decide which feature we are dealing with (aggregate with/without params or apply)

3. add it to the new name_to_function dict

4. set up the params

> **Parameters**
>
> * **columns** (*list of str*) – containing the feature names
>
> * **columns_to_ignore** (*list of str*) – columns which do not contain tsfresh feature names
>
> **Returns** The kind_to_fc_parameters object ready to be used in the extract_features function.
>
> **Return type** dict

## Module contents

The *tsfresh.feature_extraction* module contains methods to extract the features from the time series

## tsfresh.feature_selection package

## Submodules

## tsfresh.feature_selection.benjamini_hochberg_test module

tsfresh.feature_selection.benjamini_hochberg_test.**benjamini_hochberg_test**(*df_pvalues*, *hypotheses_independent*, *fdr_level*)

This is an implementation of the benjamini hochberg procedure[1] that determines if the null hypothesis for a given feature can be rejected. For this the test regards the features' p-values and controls the global false discovery rate, which is the ratio of false rejections by all rejections:

$$FDR = \mathbb{E}\left[\frac{|\text{false rejections}|}{|\text{all rejections}|}\right]$$

---

[1] Benjamini, Yoav and Yekutieli, Daniel (2001). The control of the false discovery rate in multiple testing under dependency. Annals of statistics, 1165–1188

> **Parameters**
>
> - **df_pvalues** (*pandas.DataFrame*) – This DataFrame should contain the p_values of the different hypotheses in a column named "p_values".
>
> - **hypotheses_independent** (*bool*) – Can the significance of the features be assumed to be independent? Normally, this should be set to False as the features are never independent (e.g. mean and median)
>
> - **fdr_level** (*float*) – The FDR level that should be respected, this is the theoretical expected percentage of irrelevant features among all created features.
>
> **Returns** The same DataFrame as the input, but with an added boolean column "relevant" denoting if the null hypotheses has been rejected for a given feature.
>
> **Return type** pandas.DataFrame

## tsfresh.feature_selection.relevance module

Contains a feature selection method that evaluates the importance of the different extracted features. To do so, for every feature the influence on the target is evaluated by an univariate tests and the p-Value is calculated. The methods that calculate the p-values are called feature selectors.

Afterwards the Benjamini Hochberg procedure which is a multiple testing procedure decides which features to keep and which to cut off (solely based on the p-values).

tsfresh.feature_selection.relevance.**calculate_relevance_table**(*X*, *y*, *ml_task='auto'*, *n_jobs=2*, *chunksize=None*, *test_for_binary_target_binary_feature='fi...* *test_for_binary_target_real_feature='man...* *test_for_real_target_binary_feature='man...* *test_for_real_target_real_feature='kendal...* *fdr_level=0.05*, *hypotheses_independent=False*)

Calculate the relevance table for the features contained in feature matrix *X* with respect to target vector *y*. The relevance table is calculated for the intended machine learning task *ml_task*.

To accomplish this for each feature from the input pandas.DataFrame an univariate feature significance test is conducted. Those tests generate p values that are then evaluated by the Benjamini Hochberg procedure to decide which features to keep and which to delete.

We are testing

> $H_0$ = the Feature is not relevant and should not be added

against

> $H_1$ = the Feature is relevant and should be kept

or in other words

> $H_0$ = Target and Feature are independent / the Feature has no influence on the target
>
> $H_1$ = Target and Feature are associated / dependent

When the target is binary this becomes

$$H_0 = (F_{\text{target}=1} = F_{\text{target}=0})$$

$$H_1 = (F_{\text{target}=1} \neq F_{\text{target}=0})$$

Where $F$ is the distribution of the target.

In the same way we can state the hypothesis when the feature is binary

$$H_0 = (T_{\text{feature}=1} = T_{\text{feature}=0})$$

$$H_1 = (T_{\text{feature}=1} \neq T_{\text{feature}=0})$$

Here $T$ is the distribution of the target.

TODO: And for real valued?

> **Parameters**
>
> - **X** (`pandas.DataFrame`) – Feature matrix in the format mentioned before which will be reduced to only the relevant features. It can contain both binary or real-valued features at the same time.
>
> - **y** (`pandas.Series or numpy.ndarray`) – Target vector which is needed to test which features are relevant. Can be binary or real-valued.
>
> - **ml_task** (`str`) – The intended machine learning task. Either *'classification'*, *'regression'* or *'auto'*. Defaults to *'auto'*, meaning the intended task is inferred from *y*. If *y* has a boolean, integer or object dtype, the task is assumend to be classification, else regression.
>
> - **test_for_binary_target_binary_feature** (`str`) – Which test to be used for binary target, binary feature (currently unused)
>
> - **test_for_binary_target_real_feature** (`str`) – Which test to be used for binary target, real feature
>
> - **test_for_real_target_binary_feature** (`str`) – Which test to be used for real target, binary feature (currently unused)
>
> - **test_for_real_target_real_feature** (`str`) – Which test to be used for real target, real feature (currently unused)
>
> - **fdr_level** (`float`) – The FDR level that should be respected, this is the theoretical expected percentage of irrelevant features among all created features.
>
> - **hypotheses_independent** (`bool`) – Can the significance of the features be assumed to be independent? Normally, this should be set to False as the features are never independent (e.g. mean and median)
>
> - **n_jobs** (`int`) – Number of processes to use during the p-value calculation
>
> - **chunksize** (`None or int`) – The size of one chunk that is submitted to the worker process for the parallelisation. Where one chunk is defined as a singular time series for one id and one kind. If you set the chunksize to 10, then it means that one task is to calculate all features for 10 time series. If it is set it to None, depending on distributor, heuristics are used to find the optimal chunksize. If you get out of memory exceptions, you can try it with the dask distributor and a smaller chunksize.
>
> **Returns** A pandas.DataFrame with each column of the input DataFrame X as index with information on the significance of this particular feature. The DataFrame has the columns "Feature", "type" (binary, real or const), "p_value" (the significance of this feature as a p-value, lower means more significant) "relevant" (True if the Benjamini Hochberg procedure rejected the null hypothesis [the feature is not relevant] for this feature)
>
> **Return type** pandas.DataFrame

`tsfresh.feature_selection.relevance.`**`combine_relevance_tables`**(*relevance_tables*)

    Create a combined relevance table out of a list of relevance tables, aggregating the p-values and the relevances.

        **Parameters** **`relevance_tables`**(`List[pd.DataFrame]`) – A list of relevance tables

        **Returns** The combined relevance table

        **Return type** pandas.DataFrame

`tsfresh.feature_selection.relevance.`**`get_feature_type`**(*feature_column*)

    For a given feature, determine if it is real, binary or constant. Here binary means that only two unique values occur in the feature.

        **Parameters** **`feature_column`**(*`pandas.Series`*) – The feature column

        **Returns** 'constant', 'binary' or 'real'

`tsfresh.feature_selection.relevance.`**`infer_ml_task`**(*y*)

    Infer the machine learning task to select for. The result will be either *'regression'* or *'classification'*. If the target vector only consists of integer typed values or objects, we assume the task is *'classification'*. Else *'regression'*.

        **Parameters** **`y`**(*`pandas.Series`*) – The target vector y.

        **Returns** 'classification' or 'regression'

        **Return type** str

## tsfresh.feature_selection.selection module

This module contains the filtering process for the extracted features. The filtering procedure can also be used on other features that are not based on time series.

`tsfresh.feature_selection.selection.`**`select_features`**(*X,* *y,* *test_for_binary_target_binary_feature='fisher',* *test_for_binary_target_real_feature='mann',* *test_for_real_target_binary_feature='mann',* *test_for_real_target_real_feature='kendall',* *fdr_level=0.05,* *hypotheses_independent=False,* *n_jobs=2,* *chunksize=None,* *ml_task='auto'*)

    Check the significance of all features (columns) of feature matrix X and return a possibly reduced feature matrix only containing relevant features.

    The feature matrix must be a pandas.DataFrame in the format:

| index | feature_1 | feature_2 | ... | feature_N |
|-------|-----------|-----------|-----|-----------|
| A | ... | ... | ... | ... |
| B | ... | ... | ... | ... |
| ... | ... | ... | ... | ... |
| ... | ... | ... | ... | ... |
| ... | ... | ... | ... | ... |

    Each column will be handled as a feature and tested for its significance to the target.

    The target vector must be a pandas.Series or numpy.array in the form

| index | target |
|-------|--------|
| A | ... |
| B | ... |
| . | ... |
| . | ... |

and must contain all id's that are in the feature matrix. If y is a numpy.array without index, it is assumed that y has the same order and length than X and the rows correspond to each other.

## Examples

```
>>> from tsfresh.examples import load_robot_execution_failures
>>> from tsfresh import extract_features, select_features
>>> df, y = load_robot_execution_failures()
>>> X_extracted = extract_features(df, column_id='id', column_sort='time')
>>> X_selected = select_features(X_extracted, y)
```

**Parameters**

- **X** (*pandas.DataFrame*) – Feature matrix in the format mentioned before which will be reduced to only the relevant features. It can contain both binary or real-valued features at the same time.

- **y** (*pandas.Series or numpy.ndarray*) – Target vector which is needed to test which features are relevant. Can be binary or real-valued.

- **test_for_binary_target_binary_feature** (*str*) – Which test to be used for binary target, binary feature (currently unused)

- **test_for_binary_target_real_feature** (*str*) – Which test to be used for binary target, real feature

- **test_for_real_target_binary_feature** (*str*) – Which test to be used for real target, binary feature (currently unused)

- **test_for_real_target_real_feature** (*str*) – Which test to be used for real target, real feature (currently unused)

- **fdr_level** (*float*) – The FDR level that should be respected, this is the theoretical expected percentage of irrelevant features among all created features.

- **hypotheses_independent** (*bool*) – Can the significance of the features be assumed to be independent? Normally, this should be set to False as the features are never independent (e.g. mean and median)

- **n_jobs** (*int*) – Number of processes to use during the p-value calculation

- **chunksize** (*None or int*) – The size of one chunk that is submitted to the worker process for the parallelisation. Where one chunk is defined as a singular time series for one id and one kind. If you set the chunksize to 10, then it means that one task is to calculate all features for 10 time series. If it is set it to None, depending on distributor, heuristics are used to find the optimal chunksize. If you get out of memory exceptions, you can try it with the dask distributor and a smaller chunksize.

- **ml_task** (*str*) – The intended machine learning task. Either *'classification'*, *'regression'* or *'auto'*. Defaults to *'auto'*, meaning the intended task is inferred from *y*. If *y* has a boolean, integer or object dtype, the task is assumend to be classification, else regression.

**Returns** The same DataFrame as X, but possibly with reduced number of columns ( = features).

**Return type** pandas.DataFrame

**Raises** `ValueError` when the target vector does not fit to the feature matrix or *ml_task* is not one of *'auto'*, *'classification'* or *'regression'*.

### tsfresh.feature_selection.significance_tests module

Contains the methods from the following paper about the FRESH algorithm[2]

Fresh is based on hypothesis tests that individually check the significance of every generated feature on the target. It makes sure that only features are kept, that are relevant for the regression or classification task at hand. FRESH decide between four settings depending if the features and target are binary or not.

The four functions are named

1. *target_binary_feature_binary_test()*: Target and feature are both binary

2. *target_binary_feature_real_test()*: Target is binary and feature real

3. *target_real_feature_binary_test()*: Target is real and the feature is binary

4. *target_real_feature_real_test()*: Target and feature are both real

### References

tsfresh.feature_selection.significance_tests.**target_binary_feature_binary_test**(*x*, *y*)

Calculate the feature significance of a binary feature to a binary target as a p-value. Use the two-sided univariate fisher test from `fisher_exact()` for this.

> **Parameters**
>
> - **x** (*pandas.Series*) – the binary feature vector
>
> - **y** (*pandas.Series*) – the binary target vector
>
> **Returns** the p-value of the feature significance test. Lower p-values indicate a higher feature significance
>
> **Return type** float
>
> **Raise** `ValueError` if the target or the feature is not binary.

tsfresh.feature_selection.significance_tests.**target_binary_feature_real_test**(*x*, *y*, *test*)

Calculate the feature significance of a real-valued feature to a binary target as a p-value. Use either the *Mann-Whitney U* or *Kolmogorov Smirnov* from `mannwhitneyu()` or `ks_2samp()` for this.

> **Parameters**
>
> - **x** (*pandas.Series*) – the real-valued feature vector
>
> - **y** (*pandas.Series*) – the binary target vector
>
> - **test** (*str*) – The significance test to be used. Either `'mann'` for the Mann-Whitney-U test or `'smir'` for the Kolmogorov-Smirnov test

---

[2] Christ, M., Kempa-Liehr, A.W. and Feindt, M. (2016). Distributed and parallel time series feature extraction for industrial big data applications. ArXiv e-prints: 1610.07717 https://arxiv.org/abs/1610.07717

> **Returns** the p-value of the feature significance test. Lower p-values indicate a higher feature significance
>
> **Return type** float
>
> **Raise** `ValueError` if the target is not binary.

`tsfresh.feature_selection.significance_tests.`**`target_real_feature_binary_test`**(*x*, *y*)

> Calculate the feature significance of a binary feature to a real-valued target as a p-value. Use the *Kolmogorov-Smirnov* test from from `ks_2samp()` for this.
>
> **Parameters**
>
> - **x** (*pandas.Series*) – the binary feature vector
>
> - **y** (*pandas.Series*) – the real-valued target vector
>
> **Returns** the p-value of the feature significance test. Lower p-values indicate a higher feature significance.
>
> **Return type** float
>
> **Raise** `ValueError` if the feature is not binary.

`tsfresh.feature_selection.significance_tests.`**`target_real_feature_real_test`**(*x*, *y*)

> Calculate the feature significance of a real-valued feature to a real-valued target as a p-value. Use *Kendall's tau* from `kendalltau()` for this.
>
> **Parameters**
>
> - **x** (*pandas.Series*) – the real-valued feature vector
>
> - **y** (*pandas.Series*) – the real-valued target vector
>
> **Returns** the p-value of the feature significance test. Lower p-values indicate a higher feature significance.
>
> **Return type** float

### Module contents

The *feature_selection* module contains feature selection algorithms. Those methods were suited to pick the best explaining features out of a massive amount of features. Often the features have to be picked in situations where one has more features than samples. Traditional feature selection methods can be not suitable for such situations which is why we propose a p-value based approach that inspects the significance of the features individually to avoid overfitting and spurious correlations.

### tsfresh.scripts package

### Submodules

### tsfresh.scripts.run_tsfresh module

This script can be run with:

```
python run_tsfresh.py path_to_your_csv.csv
```

A corresponding csv containing time series features will be saved as features_path_to_your_csv.csv

There are a few limitations though

- Currently this only samples to first 50 values.

- Your csv must be space delimited.

- Output is saved as path_to_your_csv.features.csv

`tsfresh.scripts.run_tsfresh.`**`main`**(*console_args=None*)

## tsfresh.scripts.test_timing module

`tsfresh.scripts.test_timing.`**`plot_results`**()

`tsfresh.scripts.test_timing.`**`test_timing`**()

`tsfresh.scripts.test_timing.`**`test_with_length`**(*length*, *df*)

## Module contents

## tsfresh.transformers package

## Submodules

## tsfresh.transformers.feature_augmenter module

**class** `tsfresh.transformers.feature_augmenter.`**`FeatureAugmenter`**(*default_fc_parameters=None*, *kind_to_fc_parameters=None*, *column_id=None*, *column_sort=None*, *column_kind=None*, *column_value=None*, *timeseries_container=None*, *chunksize=None*, *n_jobs=2*, *show_warnings=False*, *disable_progressbar=False*, *impute_function=None*, *profile=False*, *profiling_filename='profile.txt'*, *profiling_sorting='cumulative'*)

Bases: `sklearn.base.BaseEstimator`, `sklearn.base.TransformerMixin`

Sklearn-compatible estimator, for calculating and adding many features calculated from a given time series to the data. It is basically a wrapper around `extract_features()`.

The features include basic ones like min, max or median, and advanced features like fourier transformations or statistical tests. For a list of all possible features, see the module *feature_calculators*. The column name of each added feature contains the name of the function of that module, which was used for the calculation.

For this estimator, two datasets play a crucial role:

1. the time series container with the timeseries data. This container (for the format see *Data Formats*) contains the data which is used for calculating the features. It must be groupable by ids which are used to identify which feature should be attached to which row in the second dataframe:

2. the input data, where the features will be added to.

Imagine the following situation: You want to classify 10 different financial shares and you have their development in the last year as a time series. You would then start by creating features from the metainformation of the shares, e.g. how long they were on the market etc. and filling up a table - the features of one stock in one row.

```
>>> df = pandas.DataFrame()
>>> # Fill in the information of the stocks
>>> df["started_since_days"] = 0 # add a feature
```

You can then extract all the features from the time development of the shares, by using this estimator:

```
>>> time_series = read_in_timeseries() # get the development of the shares
>>> from tsfresh.transformers import FeatureAugmenter
>>> augmenter = FeatureAugmenter()
>>> augmenter.set_timeseries_container(time_series)
>>> df_with_time_series_features = augmenter.transform(df)
```

The settings for the feature calculation can be controlled with the settings object. If you pass `None`, the default settings are used. Please refer to *ComprehensiveFCParameters* for more information.

This estimator does not select the relevant features, but calculates and adds all of them to the DataFrame. See the *RelevantFeatureAugmenter* for calculating and selecting features.

For a description what the parameters column_id, column_sort, column_kind and column_value mean, please see *extraction*.

**fit** (*X=None*, *y=None*)
    The fit function is not needed for this estimator. It just does nothing and is here for compatibility reasons.

> **Parameters**
>
> > • **X** (*Any*) – Unneeded.
> >
> > • **y** (*Any*) – Unneeded.
>
> **Returns** The estimator instance itself
>
> **Return type** *FeatureAugmenter*

**set_timeseries_container** (*timeseries_container*)
    Set the timeseries, with which the features will be calculated. For a format of the time series container, please refer to *extraction*. The timeseries must contain the same indices as the later DataFrame, to which the features will be added (the one you will pass to *transform()*). You can call this function as often as you like, to change the timeseries later (e.g. if you want to extract for different ids).

> **Parameters timeseries_container** (*pandas.DataFrame or dict*) – The timeseries as a pandas.DataFrame or a dict. See *extraction* for the format.
>
> **Returns** None
>
> **Return type** None

**transform**(*X*)

    Add the features calculated using the timeseries_container and add them to the corresponding rows in the input pandas.DataFrame X.

    To save some computing time, you should only include those time serieses in the container, that you need. You can set the timeseries container with the method *set_timeseries_container()*.

        **Parameters X** (*pandas.DataFrame*) – the DataFrame to which the calculated timeseries features will be added. This is *not* the dataframe with the timeseries itself.

        **Returns** The input DataFrame, but with added features.

        **Return type** pandas.DataFrame

## tsfresh.transformers.feature_selector module

**class** tsfresh.transformers.feature_selector.**FeatureSelector**(*test_for_binary_target_binary_feature='fishe test_for_binary_target_real_feature='mann', test_for_real_target_binary_feature='mann', test_for_real_target_real_feature='kendall', fdr_level=0.05, hypotheses_independent=False, n_jobs=2, chunksize=None, ml_task='auto'*)

    Bases: sklearn.base.BaseEstimator, sklearn.base.TransformerMixin

Sklearn-compatible estimator, for reducing the number of features in a dataset to only those, that are relevant and significant to a given target. It is basically a wrapper around check_fs_sig_bh().

The check is done by testing the hypothesis

    $H_0$ = the Feature is not relevant and can not be added'

against

    $H_1$ = the Feature is relevant and should be kept

using several statistical tests (depending on whether the feature or/and the target is binary or not). Using the Benjamini Hochberg procedure, only features in $H_0$ are rejected.

This estimator - as most of the sklearn estimators - works in a two step procedure. First, it is fitted on training data, where the target is known:

```
>>> import pandas as pd
>>> X_train, y_train = pd.DataFrame(), pd.Series() # fill in with your features
↪and target
>>> from tsfresh.transformers import FeatureSelector
>>> selector = FeatureSelector()
>>> selector.fit(X_train, y_train)
```

In this example the list of relevant features is empty:

```
>>> selector.relevant_features
>>> []
```

The same holds for the feature importance:

```
>>> selector.feature_importances_
>>> array([], dtype=float64)
```

The estimator keeps track on those features, that were relevant in the training step. If you apply the estimator after the training, it will delete all other features in the testing data sample:

```
>>> X_test = pd.DataFrame()
>>> X_selected = selector.transform(X_test)
```

After that, X_selected will only contain the features that were relevant during the training.

If you are interested in more information on the features, you can look into the member `relevant_features` after the fit.

**fit**(*X*, *y*)

> Extract the information, which of the features are relevent using the given target.
>
> For more information, please see the check_fs_sig_bh() function. All columns in the input data sample are treated as feature. The index of all rows in X must be present in y.
>
> > **Parameters**
> >
> > - **X** (*pandas.DataFrame or numpy.array*) – data sample with the features, which will be classified as relevant or not
> >
> > - **y** (*pandas.Series or numpy.array*) – target vector to be used, to classify the features
> >
> > **Returns**  the fitted estimator with the information, which features are relevant
> >
> > **Return type**  *FeatureSelector*

**transform**(*X*)

> Delete all features, which were not relevant in the fit phase.
>
> > **Parameters X** (*pandas.DataSeries or numpy.array*) – data sample with all features, which will be reduced to only those that are relevant
> >
> > **Returns**  same data sample as X, but with only the relevant features
> >
> > **Return type**  pandas.DataFrame or numpy.array

## tsfresh.transformers.per_column_imputer module

**class** tsfresh.transformers.per_column_imputer.**PerColumnImputer**(*col_to_NINF_repl_preset=None*, *col_to_PINF_repl_preset=None*, *col_to_NAN_repl_preset=None*)

> Bases: sklearn.base.BaseEstimator, sklearn.base.TransformerMixin
>
> Sklearn-compatible estimator, for column-wise imputing DataFrames by replacing all `NaNs` and `infs` with with average/extreme values from the same columns. It is basically a wrapper around *impute()*.
>
> Each occurring `inf` or `NaN` in the DataFrame is replaced by
>
> - `-inf` -> `min`
>
> - `+inf` -> `max`
>
> - `NaN` -> `median`

This estimator - as most of the sklearn estimators - works in a two step procedure. First, the `.fit` function is called where for each column the min, max and median are computed. Secondly, the `.transform` function is called which replaces the occurances of `NaNs` and `infs` using the column-wise computed min, max and median values.

**fit**(*X*, *y=None*)

Compute the min, max and median for all columns in the DataFrame. For more information, please see the *get_range_values_per_column()* function.

> **Parameters**
>
> - **X** (*pandas.DataFrame*) – DataFrame to calculate min, max and median values on
>
> - **y** (*Any*) – Unneeded.
>
> **Returns** the estimator with the computed min, max and median values
>
> **Return type** Imputer

**transform**(*X*)

Column-wise replace all `NaNs`, `-inf` and `+inf` in the DataFrame *X* with average/extreme values from the provided dictionaries.

> **Parameters** **X** (*pandas.DataFrame*) – DataFrame to impute
>
> **Returns** imputed DataFrame
>
> **Return type** pandas.DataFrame
>
> **Raises** **RuntimeError** – if the replacement dictionaries are still of None type. This can happen if the transformer was not fitted.

### tsfresh.transformers.relevant_feature_augmenter module

**class** tsfresh.transformers.relevant_feature_augmenter.**RelevantFeatureAugmenter**(*filter_only_tsfres, de-fault_fc_paramet, kind_to_fc_parar, col-umn_id=None, col-umn_sort=None, col-umn_kind=None, col-umn_value=None, time-series_container= chunk-size=None, n_jobs=2, show_warnings= dis-able_progressbar pro-file=False, pro-fil-ing_filename='pr pro-fil-ing_sorting='cur test_for_binary_i test_for_binary_i test_for_real_tar, test_for_real_tar, fdr_level=0.05, hy-pothe-ses_independent= ml_task='auto'*)

Bases: `sklearn.base.BaseEstimator`, `sklearn.base.TransformerMixin`

Sklearn-compatible estimator to calculate relevant features out of a time series and add them to a data sample.

As many other sklearn estimators, this estimator works in two steps:

In the fit phase, all possible time series features are calculated using the time series, that is set by the set_timeseries_container function (if the features are not manually changed by handing in a fea-ture_extraction_settings object). Then, their significance and relevance to the target is computed using statistical methods and only the relevant ones are selected using the Benjamini Hochberg procedure. These features are stored internally.

In the transform step, the information on which features are relevant from the fit step is used and those features are extracted from the time series. These extracted features are then added to the input data sample.

This estimator is a wrapper around most of the functionality in the tsfresh package. For more information on the subtasks, please refer to the single modules and functions, which are:

- Settings for the feature extraction: *ComprehensiveFCParameters*
- Feature extraction method: *extract_features()*
- Extracted features: *feature_calculators*
- Feature selection: check_fs_sig_bh()

This estimator works analogue to the *FeatureAugmenter* with the difference that this estimator does only output and calculate the relevant features, whereas the other outputs all features.

Also for this estimator, two datasets play a crucial role:

1. the time series container with the timeseries data. This container (for the format see *extraction*) contains the data which is used for calculating the features. It must be groupable by ids which are used to identify which feature should be attached to which row in the second dataframe:

2. the input data, where the features will be added to.

Imagine the following situation: You want to classify 10 different financial shares and you have their development in the last year as a time series. You would then start by creating features from the metainformation of the shares, e.g. how long they were on the market etc. and filling up a table - the features of one stock in one row.

```
>>> # Fill in the information of the stocks and the target
>>> X_train, X_test, y_train = pd.DataFrame(), pd.DataFrame(), pd.Series()
```

You can then extract all the relevant features from the time development of the shares, by using this estimator:

```
>>> train_time_series, test_time_series = read_in_timeseries() # get the
↪development of the shares
>>> from tsfresh.transformers import RelevantFeatureAugmenter
>>> augmenter = RelevantFeatureAugmenter()
>>> augmenter.set_timeseries_container(train_time_series)
>>> augmenter.fit(X_train, y_train)
>>> augmenter.set_timeseries_container(test_time_series)
>>> X_test_with_features = augmenter.transform(X_test)
```

X_test_with_features will then contain the same information as X_test (with all the meta information you have probably added) plus some relevant time series features calculated on the time series you handed in.

Please keep in mind that the time series you hand in before fit or transform must contain data for the rows that are present in X.

If your set filter_only_tsfresh_features to True, your manually-created features that were present in X_train (or X_test) before using this estimator are not touched. Otherwise, also those features are evaluated and may be rejected from the data sample, because they are irrelevant.

For a description what the parameters column_id, column_sort, column_kind and column_value mean, please see *extraction*.

You can control the feature extraction in the fit step (the feature extraction in the transform step is done automatically) as well as the feature selection in the fit step by handing in settings. However, the default settings which are used if you pass no flags are often quite sensible.

**fit**(*X*, *y*)

Use the given timeseries from *set_timeseries_container()* and calculate features from it and add them to the data sample X (which can contain other manually-designed features).

Then determine which of the features of X are relevant for the given target y. Store those relevant features internally to only extract them in the transform step.

If filter_only_tsfresh_features is True, only reject newly, automatically added features. If it is False, also look at the features that are already present in the DataFrame.

**Parameters**

- **X** (*pandas.DataFrame or numpy.array*) – The data frame without the time series features. The index rows should be present in the timeseries and in the target vector.

- **y** (*pandas.Series or numpy.array*) – The target vector to define, which features are relevant.

**Returns** the fitted estimator with the information, which features are relevant.

**Return type** *RelevantFeatureAugmenter*

**set_timeseries_container**(*timeseries_container*)

Set the timeseries, with which the features will be calculated. For a format of the time series container, please refer to *extraction*. The timeseries must contain the same indices as the later DataFrame, to which the features will be added (the one you will pass to *transform()* or *fit()*). You can call this function as often as you like, to change the timeseries later (e.g. if you want to extract for different ids).

**Parameters timeseries_container** (*pandas.DataFrame or dict*) – The timeseries as a pandas.DataFrame or a dict. See *extraction* for the format.

**Returns** None

**Return type** None

**transform**(*X*)

After the fit step, it is known which features are relevant, Only extract those from the time series handed in with the function *set_timeseries_container()*.

If filter_only_tsfresh_features is False, also delete the irrelevant, already present features in the data frame.

**Parameters X** (*pandas.DataFrame or numpy.array*) – the data sample to add the relevant (and delete the irrelevant) features to.

**Returns** a data sample with the same information as X, but with added relevant time series features and deleted irrelevant information (only if filter_only_tsfresh_features is False).

**Return type** pandas.DataFrame

## Module contents

The module *transformers* contains several transformers which can be used inside a sklearn pipeline.

## tsfresh.utilities package

## Submodules

## tsfresh.utilities.dataframe_functions module

Utility functions for handling the DataFrame conversions to the internal normalized format (see `normalize_input_to_internal_representation`) or on how to handle `NaN` and `inf` in the DataFrames.

tsfresh.utilities.dataframe_functions.**check_for_nans_in_columns**(*df*, *columns=None*)

Helper function to check for `NaN` in the data frame and raise a `ValueError` if there is one.

**Parameters**

- **df** (*pandas.DataFrame*) – the pandas DataFrame to test for NaNs

- **columns** (*list*) – a list of columns to test for NaNs. If left empty, all columns of the DataFrame will be tested.

**Returns** None

**Return type** [None](#)

**Raise** `ValueError` of NaNs are found in the DataFrame.

`tsfresh.utilities.dataframe_functions.`**`get_range_values_per_column`**(*df*)

Retrieves the finite max, min and mean values per column in the DataFrame *df* and stores them in three dictionaries. Those dictionaries *col_to_max*, *col_to_min*, *col_to_median* map the columnname to the maximal, minimal or median value of that column.

If a column does not contain any finite values at all, a 0 is stored instead.

**Parameters df** ([*pandas.DataFrame*](#)) – the Dataframe to get columnswise max, min and median from

**Returns** Dictionaries mapping column names to max, min, mean values

**Return type** ([dict](#), [dict](#), [dict](#))

`tsfresh.utilities.dataframe_functions.`**`impute`**(*df_impute*)

Columnwise replaces all `NaNs` and `infs` from the DataFrame *df_impute* with average/extreme values from the same columns. This is done as follows: Each occurring `inf` or `NaN` in *df_impute* is replaced by

- `-inf` -> `min`

- `+inf` -> `max`

- `NaN` -> `median`

If the column does not contain finite values at all, it is filled with zeros.

This function modifies *df_impute* in place. After that, df_impute is guaranteed to not contain any non-finite values. Also, all columns will be guaranteed to be of type `np.float64`.

**Parameters df_impute** ([*pandas.DataFrame*](#)) – DataFrame to impute

**Return df_impute** imputed DataFrame

**Rtype df_impute** pandas.DataFrame

`tsfresh.utilities.dataframe_functions.`**`impute_dataframe_range`**(*df_impute*,
*col_to_max*,
*col_to_min*,
*col_to_median*)

Columnwise replaces all `NaNs`, `-inf` and `+inf` from the DataFrame *df_impute* with average/extreme values from the provided dictionaries.

This is done as follows: Each occurring `inf` or `NaN` in *df_impute* is replaced by

- `-inf` -> by value in col_to_min

- `+inf` -> by value in col_to_max

- `NaN` -> by value in col_to_median

If a column of df_impute is not found in the one of the dictionaries, this method will raise a ValueError. Also, if one of the values to replace is not finite a ValueError is returned

This function modifies *df_impute* in place. Afterwards df_impute is guaranteed to not contain any non-finite values. Also, all columns will be guaranteed to be of type `np.float64`.

**Parameters**

- **df_impute** (*pandas.DataFrame*) – DataFrame to impute
- **col_to_max** (*dict*) – Dictionary mapping column names to max values
- **col_to_min** – Dictionary mapping column names to min values
- **col_to_median** – Dictionary mapping column names to median values

**Return df_impute** imputed DataFrame

**Rtype df_impute** pandas.DataFrame

**Raises ValueError** – if a column of df_impute is missing in col_to_max, col_to_min or col_to_median or a value to replace is non finite

tsfresh.utilities.dataframe_functions.**impute_dataframe_zero**(*df_impute*)

Replaces all NaNs, -infs and +infs from the DataFrame *df_impute* with 0s. The *df_impute* will be modified in place. All its columns will be into converted into dtype np.float64.

**Parameters df_impute** (*pandas.DataFrame*) – DataFrame to impute

**Return df_impute** imputed DataFrame

**Rtype df_impute** pandas.DataFrame

tsfresh.utilities.dataframe_functions.**make_forecasting_frame**(*x*, *kind*, *max_timeshift*, *rolling_direction*)

Takes a singular time series x and constructs a DataFrame df and target vector y that can be used for a time series forecasting task.

The returned df will contain, for every time stamp in x, the last max_timeshift data points as a new time series, such can be used to fit a time series forecasting model.

See *Time series forecasting* for a detailed description of the rolling process and how the feature matrix and target vector are derived.

The returned time series container df, will contain the rolled time series as a flat data frame, the first format from *Data Formats*.

When x is a pandas.Series, the index will be used as id.

**Parameters**

- **x** (*np.array or pd.Series*) – the singular time series
- **kind** (*str*) – the kind of the time series
- **rolling_direction** (*int*) – The sign decides, if to roll backwards (if sign is positive) or forwards in "time"
- **max_timeshift** (*int*) – If not None, shift only up to max_timeshift. If None, shift as often as possible.

**Returns** time series container df, target vector y

**Return type** (pd.DataFrame, pd.Series)

tsfresh.utilities.dataframe_functions.**restrict_input_to_index**(*df_or_dict*, *column_id*, *index*)

Restrict df_or_dict to those ids contained in index.

**Parameters**

- **df_or_dict** (*pandas.DataFrame or dict*) – a pandas DataFrame or a dictionary.

- **column_id** (*basestring*) – it must be present in the pandas DataFrame or in all DataFrames in the dictionary. It is not allowed to have NaN values in this column.

- **index** (*Iterable or `pandas.Series`*) – Index containing the ids

**Return df_or_dict_restricted** the restricted df_or_dict

**Rtype df_or_dict_restricted** dict or pandas.DataFrame

**Raise** `TypeError` if df_or_dict is not of type dict or pandas.DataFrame

`tsfresh.utilities.dataframe_functions.`**`roll_time_series`**(*df_or_dict*, *column_id*, *column_sort*, *column_kind*, *rolling_direction*, *max_timeshift=None*)

This method creates sub windows of the time series. It rolls the (sorted) data frames for each kind and each id separately in the "time" domain (which is represented by the sort order of the sort column given by *column_sort*).

For each rolling step, a new id is created by the scheme "id={id}, shift={shift}", here id is the former id of the column and shift is the amount of "time" shifts.

A few remarks:

- This method will create new IDs!

- The sign of rolling defines the direction of time rolling, a positive value means we are going back in time

- It is possible to shift time series of different lengths but

- We assume that the time series are uniformly sampled

- For more information, please see *Time series forecasting*.

**Parameters**

- **df_or_dict** (*pandas.DataFrame or dict*) – a pandas DataFrame or a dictionary. The required shape/form of the object depends on the rest of the passed arguments.

- **column_id** (*basestring or None*) – it must be present in the pandas DataFrame or in all DataFrames in the dictionary. It is not allowed to have NaN values in this column.

- **column_sort** (*basestring or None*) – if not None, sort the rows by this column. It is not allowed to have NaN values in this column.

- **column_kind** (*basestring or None*) – It can only be used when passing a pandas DataFrame (the dictionary is already assumed to be grouped by the kind). Is must be present in the DataFrame and no NaN values are allowed. If the kind column is not passed, it is assumed that each column in the pandas DataFrame (except the id or sort column) is a possible kind.

- **rolling_direction** (*int*) – The sign decides, if to roll backwards or forwards in "time"

- **max_timeshift** (*int*) – If not None, shift only up to max_timeshift. If None, shift as often as possible.

**Returns** The rolled data frame or dictionary of data frames

**Return type** the one from df_or_dict

### tsfresh.utilities.distribution module

This module contains the Distributor class, such objects are used to distribute the calculation of features. Essentially, a Distributor organizes the application of feature calculators to data chunks.

Design of this module by Nils Braun

**class** tsfresh.utilities.distribution.**ClusterDaskDistributor**(*address*)

Bases: *tsfresh.utilities.distribution.DistributorBaseClass*

Distributor using a dask cluster, meaning that the calculation is spread over a cluster

**calculate_best_chunk_size**(*data_length*)

Uses the number of dask workers in the cluster (during execution time, meaning when you start the extraction) to find the optimal chunk_size.

**Parameters data_length** (*int*) – A length which defines how many calculations there need to be.

**close**()

Closes the connection to the Dask Scheduler

**distribute**(*func*, *partitioned_chunks*, *kwargs*)

Calculates the features in a parallel fashion by distributing the map command to the dask workers on a cluster

**Parameters**

- **func** (*callable*) – the function to send to each worker.

- **partitioned_chunks** (*iterable*) – The list of data chunks - each element is again a list of chunks - and should be processed by one worker.

- **kwargs** (*dict of string to parameter*) – parameters for the map function

**Returns** The result of the calculation as a list - each item should be the result of the application of func to a single element.

**class** tsfresh.utilities.distribution.**DistributorBaseClass**

The distributor abstract base class.

The main purpose of the instances of the DistributorBaseClass subclasses is to evaluate a function (called map_function) on a list of data items (called data).

This is done on chunks of the data, meaning, that the DistributorBaseClass classes will chunk the data into chunks, distribute the data and apply the feature calculator functions from *tsfresh.feature_extraction.feature_calculators* on the time series.

Dependent on the implementation of the distribute function, this is done in parallel or using a cluster of nodes.

**calculate_best_chunk_size**(*data_length*)

Calculates the best chunk size for a list of length data_length. The current implemented formula is more or less an empirical result for multiprocessing case on one machine.

**Parameters data_length** (*int*) – A length which defines how many calculations there need to be.

**Returns** the calculated chunk size

**Return type** int

TODO: Investigate which is the best chunk size for different settings.

**close**()
> Abstract base function to clean the DistributorBaseClass after use, e.g. close the connection to a DaskScheduler

**distribute**(*func*, *partitioned_chunks*, *kwargs*)
> This abstract base function distributes the work among workers, which can be threads or nodes in a cluster. Must be implemented in the derived classes.

> **Parameters**
>> - **func** (`callable`) – the function to send to each worker.
>> - **partitioned_chunks** (`iterable`) – The list of data chunks - each element is again a list of chunks - and should be processed by one worker.
>> - **kwargs** (`dict of string to parameter`) – parameters for the map function

> **Returns** The result of the calculation as a list - each item should be the result of the application of func to a single element.

**map_reduce**(*map_function*, *data*, *function_kwargs=None*, *chunk_size=None*, *data_length=None*)
> This method contains the core functionality of the DistributorBaseClass class.

> It maps the map_function to each element of the data and reduces the results to return a flattened list.

> How the jobs are calculated, is determined by the classes *tsfresh.utilities.distribution. DistributorBaseClass.distribute()* method, which can distribute the jobs in multiple threads, across multiple processing units etc.

> To not transport each element of the data individually, the data is split into chunks, according to the chunk size (or an empirical guess if none is given). By this, worker processes not tiny but adequate sized parts of the data.

> **Parameters**
>> - **map_function** (`callable`) – a function to apply to each data item.
>> - **data** (`iterable`) – the data to use in the calculation
>> - **function_kwargs** (`dict of string to parameter`) – parameters for the map function
>> - **chunk_size** (`int`) – If given, chunk the data according to this size. If not given, use an empirical value.
>> - **data_length** (`int`) – If the data is a generator, you have to set the length here. If it is none, the length is deduced from the len of the data.

> **Returns** the calculated results

> **Return type** list

**static partition**(*data*, *chunk_size*)
> This generator chunks a list of data into slices of length chunk_size. If the chunk_size is not a divider of the data length, the last slice will be shorter than chunk_size.

> **Parameters**
>> - **data** (`list`) – The data to chunk.
>> - **chunk_size** (`int`) – Each chunks size. The last chunk may be smaller.

> **Returns** A generator producing the chunks of data.

> **Return type** generator

**class** tsfresh.utilities.distribution.**LocalDaskDistributor**(*n_workers*)

    Bases: *tsfresh.utilities.distribution.DistributorBaseClass*

    Distributor using a local dask cluster and inproc communication.

    **close**()

        Closes the connection to the local Dask Scheduler

    **distribute**(*func*, *partitioned_chunks*, *kwargs*)

        Calculates the features in a parallel fashion by distributing the map command to the dask workers on a local machine

        **Parameters**

            • **func** (*callable*) – the function to send to each worker.

            • **partitioned_chunks** (*iterable*) – The list of data chunks - each element is again a list of chunks - and should be processed by one worker.

            • **kwargs** (*dict of string to parameter*) – parameters for the map function

        **Returns** The result of the calculation as a list - each item should be the result of the application of func to a single element.

**class** tsfresh.utilities.distribution.**MapDistributor**(*disable_progressbar=False*, *progressbar_title='Feature Extraction'*)

    Bases: *tsfresh.utilities.distribution.DistributorBaseClass*

    Distributor using the python build-in map, which calculates each job sequentially one after the other.

    **calculate_best_chunk_size**(*data_length*)

        For the map command, which calculates the features sequentially, a the chunk_size of 1 will be used.

        **Parameters data_length** (*int*) – A length which defines how many calculations there need to be.

    **distribute**(*func*, *partitioned_chunks*, *kwargs*)

        Calculates the features in a sequential fashion by pythons map command

        **Parameters**

            • **func** (*callable*) – the function to send to each worker.

            • **partitioned_chunks** (*iterable*) – The list of data chunks - each element is again a list of chunks - and should be processed by one worker.

            • **kwargs** (*dict of string to parameter*) – parameters for the map function

        **Returns** The result of the calculation as a list - each item should be the result of the application of func to a single element.

**class** tsfresh.utilities.distribution.**MultiprocessingDistributor**(*n_workers*, *disable_progressbar=False*, *progressbar_title='Feature Extraction'*)

    Bases: *tsfresh.utilities.distribution.DistributorBaseClass*

    Distributor using a multiprocessing Pool to calculate the jobs in parallel on the local machine.

    **close**()

        Collects the result from the workers and closes the thread pool.

**distribute**(*func*, *partitioned_chunks*, *kwargs*)

> Calculates the features in a parallel fashion by distributing the map command to a thread pool

> > **Parameters**

> > > - **func** (`callable`) – the function to send to each worker.
> > > - **partitioned_chunks** (`iterable`) – The list of data chunks - each element is again a list of chunks - and should be processed by one worker.
> > > - **kwargs** (`dict of string to parameter`) – parameters for the map function

> > **Returns** The result of the calculation as a list - each item should be the result of the application of func to a single element.

## tsfresh.utilities.profiling module

Contains methods to start and stop the profiler that checks the runtime of the different feature calculators

tsfresh.utilities.profiling.**end_profiling**(*profiler*, *filename*, *sorting=None*)

> Helper function to stop the profiling process and write out the profiled data into the given filename. Before this, sort the stats by the passed sorting.

> > **Parameters**

> > > - **profiler** (`cProfile.Profile`) – An already started profiler (probably by start_profiling).
> > > - **filename** (`basestring`) – The name of the output file to save the profile.
> > > - **sorting** (`basestring`) – The sorting of the statistics passed to the sort_stats function.

> > **Returns** None

> > **Return type** None

> Start and stop the profiler with:

```
>>> profiler = start_profiling()
>>> # Do something you want to profile
>>> end_profiling(profiler, "out.txt", "cumulative")
```

tsfresh.utilities.profiling.**start_profiling**()

> Helper function to start the profiling process and return the profiler (to close it later).

> > **Returns** a started profiler.

> > **Return type** cProfile.Profile

> Start and stop the profiler with:

```
>>> profiler = start_profiling()
>>> # Do something you want to profile
>>> end_profiling(profiler, "cumulative", "out.txt")
```

## tsfresh.utilities.string_manipulation module

tsfresh.utilities.string_manipulation.**convert_to_output_format**(*param*)

> Helper function to convert parameters to a valid string, that can be used in a column name. Does the opposite which is used in the from_columns function.

The parameters are sorted by their name and written out in the form

&lt;param name&gt;_&lt;param value&gt;__&lt;param name&gt;_&lt;param value&gt;__ . . .

If a &lt;param_value&gt; is a string, this method will wrap it with parenthesis ", so "&lt;param_value&gt;"

> **Parameters param** (*dict*) – The dictionary of parameters to write out
>
> **Returns** The string of parsed parameters
>
> **Return type** str

tsfresh.utilities.string_manipulation.**get_config_from_string**(*parts*)

Helper function to extract the configuration of a certain function from the column name. The column name parts (split by "__") should be passed to this function. It will skip the kind name and the function name and only use the parameter parts. These parts will be split up on "_" into the parameter name and the parameter value. This value is transformed into a python object (for example is "(1, 2, 3)" transformed into a tuple consisting of the ints 1, 2 and 3).

Returns None of no parameters are in the column name.

> **Parameters parts** (*list*) – The column name split up on "__"
>
> **Returns** a dictionary with all parameters, which are encoded in the column name.
>
> **Return type** dict

## Module contents

This *utilities* submodule contains several utility functions. Those should only be used internally inside tsfresh.

## Submodules

## tsfresh.defaults module

## Module contents

At the top level we export the three most important submodules of tsfresh, which are:

- extract_features
- select_features
- extract_relevant_features

# 1.4 Data Formats

tsfresh offers three different options to specify the time series data to be used in the tsfresh. extract_features() function (and all utility functions that expect a time series, e.g. the *tsfresh. utilities.dataframe_functions.roll_time_series()* function).

Irrespective of the input format, tsfresh will always return the calculated features in the same output format described below.

All three input format options consist of pandas.DataFrame objects. There are four important column types that make up those DataFrames. Each will be described with an example from the robot failures dataset (see *Quick Start*).

Mandatory:

*column_id* This column indicates which entities the time series belong to. Features will be extracted individually for each entity. The resulting feature matrix will contain one row per entity. Each robot is a different entity, so each of it has a different id.

*column_value* This column contains the actual values of the time series. This corresponds to the measured values for different the sensors on the robots.

Optional (but strongly recommended to specify if you have this column):

*column_sort* This column contains values which allow to sort the time series (e.g. time stamps). It is not required to have equidistant time steps or the same time scale for the different ids and/or kinds. If you omit this column, the DataFrame is assumed to be already sorted in increasing order. The robot sensor measurements each have a time stamp which is used in this column.

Please note that none of the algorithms of tsfresh uses the actual values in this time column - but only their sorting order.

Optional:

*column_kind* This column indicates the names of the different time series types (E.g. different sensors in an industrial application as in the robot dataset). For each kind of time series the features are calculated individually.

Important: None of these columns is allowed to contain any `NaN`, `Inf` or `-Inf` values.

**In the following we describe the different input formats, that are build on those columns:**

- A flat DataFrame
- A stacked DataFrame
- A dictionary of flat DataFrames

The difference between a flat and a stacked DataFrame is indicated by specifying or not specifying the parameters *column_value* and *column_kind* in the `tsfresh.extract_features()` function.

If you do not know which one to choose, you probably want to try out the flat or stacked DataFrame.

## 1.4.1 Input Option 1. Flat DataFrame

If both *column_value* and *column_kind* are set to `None`, the time series data is assumed to be in a flat DataFrame. This means that each different time series must be saved as its own column.

Example: Imagine you record the values of time series x and y for different objects A and B for three different times t1, t2 and t3. Now you want to calculate some feature with tsfresh. Your resulting DataFrame may look like this:

| id | time | x | y |
|----|------|--------|--------|
| A | t1 | x(A, t1) | y(A, t1) |
| A | t2 | x(A, t2) | y(A, t2) |
| A | t3 | x(A, t3) | y(A, t3) |
| B | t1 | x(B, t1) | y(B, t1) |
| B | t2 | x(B, t2) | y(B, t2) |
| B | t3 | x(B, t3) | y(B, t3) |

and you would pass

```
column_id="id", column_sort="time", column_kind=None, column_value=None
```

to the extraction functions, to extract features separately for all ids and separately for the x and y values.

## 1.4.2 Input Option 2. Stacked DataFrame

If both *column_value* and *column_kind* are set, the time series data is assumed to be a stacked DataFrame. This means that there are no different columns for the different types of time series. This representation has several advantages over the flat Data Frame. For example, the time stamps of the different time series do not have to align.

It does not contain different columns for the different types of time series but only one value column and a kind column. The example from above would look like this:

| id | time | kind | value |
|----|------|------|--------|
| A | t1 | x | x(A, t1) |
| A | t2 | x | x(A, t2) |
| A | t3 | x | x(A, t3) |
| A | t1 | y | y(A, t1) |
| A | t2 | y | y(A, t2) |
| A | t3 | y | y(A, t3) |
| B | t1 | x | x(B, t1) |
| B | t2 | x | x(B, t2) |
| B | t3 | x | x(B, t3) |
| B | t1 | y | y(B, t1) |
| B | t2 | y | y(B, t2) |
| B | t3 | y | y(B, t3) |

Then you would set

```
column_id="id", column_sort="time", column_kind="kind", column_value="value"
```

to end up with the same extracted features as above.

## 1.4.3 Input Option 3. Dictionary of flat DataFrames

Instead of passing a DataFrame which must be split up by its different kinds by tsfresh, you can also give a dictionary mapping from the kind as string to a DataFrame containing only the time series data of that kind. So essentially you are using a singular DataFrame for each kind of time series.

The data from the example can be split into two DataFrames resulting in the following dictionary

{ "x":

| id | time | value |
|----|------|--------|
| A | t1 | x(A, t1) |
| A | t2 | x(A, t2) |
| A | t3 | x(A, t3) |
| B | t1 | x(B, t1) |
| B | t2 | x(B, t2) |
| B | t3 | x(B, t3) |

, "y":

| id | time | value |
|----|------|-------|
| A  | t1   | y(A, t1) |
| A  | t2   | y(A, t2) |
| A  | t3   | y(A, t3) |
| B  | t1   | y(B, t1) |
| B  | t2   | y(B, t2) |
| B  | t3   | y(B, t3) |

}

You would pass this dictionary to tsfresh together with the following arguments:

```
column_id="id", column_sort="time", column_kind=None, column_value="value":
```

In this case we do not need to specify the kind column as the kind is the respective dictionary key.

### 1.4.4 Output Format

The resulting feature matrix for all three input options will be the same. It will always be a `pandas.DataFrame` with the following layout

| id | x_feature_1 | … | x_feature_N | y_feature_1 | … | y_feature_N |
|----|-------------|---|-------------|-------------|---|-------------|
| A  | …           | … | …           | …           | … | …           |
| B  | …           | … | …           | …           | … | …           |

where the x features are calculated using all x values (independently for A and B), y features using all y values and so on.

This form of DataFrame is also the expected input format to the feature selection algorithms (e.g. the `tsfresh.select_features()` function).

## 1.5 scikit-learn Transformers

tsfresh includes three scikit-learn compatible transformers. You can easily add them to your existing data science pipeline. If you are not familiar with scikit-learn's pipeline we recommend you take a look at the official documentation[1].

The purpose of such a pipeline is to assemble several preprocessing steps that can be cross-validated together while setting different parameters. Our tsfresh transformer allows you to extract and filter the time series features during such a preprocessing sequence.

The first two estimators contained in tsfresh are the *FeatureAugmenter*, which extracts the features, and the *FeatureSelector*, which only performs the feature selection algorithm. It is preferable to combine extracting and filtering of the features in a single step to avoid unnecessary feature calculations. Hence, we have the `RelevantFeatureAugmenter`, which combines both the extraction and filtering of the features in a single step.

---

[1] http://scikit-learn.org/stable/modules/generated/sklearn.pipeline.Pipeline.html

### 1.5.1 Example

In the following example you see how we combine tsfresh's *RelevantFeatureAugmenter* and a RandomForestClassifier into a single pipeline. This pipeline can then fit both our transformer and the classifier in one step.

```python
from sklearn.pipeline import Pipeline
from sklearn.ensemble import RandomForestClassifier
from tsfresh.examples import load_robot_execution_failures
from tsfresh.transformers import RelevantFeatureAugmenter

pipeline = Pipeline([('augmenter', RelevantFeatureAugmenter(column_id='id', column_
→sort='time')),
            ('classifier', RandomForestClassifier())])

df_ts, y = load_robot_execution_failures()
X = pd.DataFrame(index=y.index)

pipeline.set_params(augmenter__timeseries_container=df_ts)
pipeline.fit(X, y)
```

The parameters of the augment transformer correspond to the parameters of the top-level convenience function *extract_relevant_features()*. In the example, we only set the names of two columns `column_id='id'`, `column_sort='time'` (see *Data Formats* for an explanation of those parameters).

Because we cannot pass the time series container directly as a parameter to the augmenter step when calling fit or transform on a `sklearn.pipeline.Pipeline` we have to set it manually by calling `pipeline.set_params(augmenter__timeseries_container=df_ts)`. In general, you can change the time series container from which the features are extracted by calling either the pipeline's `set_params()` method or the transformers *set_timeseries_container()* method.

For further examples, see the Jupyter Notebook pipeline_example.ipynb in the notebooks folder of the tsfresh package.

### 1.5.2 References

## 1.6 Overview on extracted features

*tsfresh* calculates a comprehensive number of features. All feature calculators are contained in the

| | |
|---|---|
| *tsfresh.feature_extraction. feature_calculators* | This module contains the feature calculators that take time series as input and calculate the values of the feature. |

| | |
|---|---|
| *abs_energy*(x) | Returns the absolute energy of the time series which is the sum over the squared values |
| *absolute_sum_of_changes*(x) | Returns the sum over the absolute value of consecutive changes in the series x |
| *agg_autocorrelation*(x, param) | Calculates the value of an aggregation function f_agg (e.g. |

Continued on next page

Table 2 – continued from previous page

| | |
|---|---|
| *agg_linear_trend*(x, param) | Calculates a linear least-squares regression for values of the time series that were aggregated over chunks versus the sequence from 0 up to the number of chunks minus one. |
| *approximate_entropy*(x, m, r) | Implements a vectorized Approximate entropy algorithm. |
| *ar_coefficient*(x, param) | This feature calculator fits the unconditional maximum likelihood of an autoregressive AR(k) process. |
| *augmented_dickey_fuller*(x, param) | The Augmented Dickey-Fuller test is a hypothesis test which checks whether a unit root is present in a time series sample. |
| *autocorrelation*(x, lag) | Calculates the autocorrelation of the specified lag, according to the formula [1] |
| *binned_entropy*(x, max_bins) | First bins the values of x into max_bins equidistant bins. |
| *c3*(x, lag) | This function calculates the value of |
| *change_quantiles*(x, ql, qh, isabs, f_agg) | First fixes a corridor given by the quantiles ql and qh of the distribution of x. |
| *cid_ce*(x, normalize) | This function calculator is an estimate for a time series complexity [1] (A more complex time series has more peaks, valleys etc.). |
| *count_above_mean*(x) | Returns the number of values in x that are higher than the mean of x |
| *count_below_mean*(x) | Returns the number of values in x that are lower than the mean of x |
| *cwt_coefficients*(x, param) | Calculates a Continuous wavelet transform for the Ricker wavelet, also known as the "Mexican hat wavelet" which is defined by |
| *energy_ratio_by_chunks*(x, param) | Calculates the sum of squares of chunk i out of N chunks expressed as a ratio with the sum of squares over the whole series |
| *fft_aggregated*(x, param) | Returns the spectral centroid (mean), variance, skew, and kurtosis of the absolute fourier transform spectrum. |
| *fft_coefficient*(x, param) | Calculates the fourier coefficients of the one-dimensional discrete Fourier Transform for real input by fast fourier transformation algorithm |
| *first_location_of_maximum*(x) | Returns the first location of the maximum value of x. |
| *first_location_of_minimum*(x) | Returns the first location of the minimal value of x. |
| *friedrich_coefficients*(x, param) | Coefficients of polynomial $h(x)$, which has been fitted to the deterministic dynamics of Langevin model |
| *has_duplicate*(x) | Checks if any value in x occurs more than once |
| *has_duplicate_max*(x) | Checks if the maximum value of x is observed more than once |
| *has_duplicate_min*(x) | Checks if the minimal value of x is observed more than once |
| *index_mass_quantile*(x, param) | Those apply features calculate the relative index i where q% of the mass of the time series x lie left of i. |
| *kurtosis*(x) | Returns the kurtosis of x (calculated with the adjusted Fisher-Pearson standardized moment coefficient G2). |
| *large_standard_deviation*(x, r) | Boolean variable denoting if the standard dev of x is higher than 'r' times the range = difference between max and min of x. |

Continued on next page

Table 2 – continued from previous page

| | |
|---|---|
| *last_location_of_maximum*(x) | Returns the relative last location of the maximum value of x. |
| *last_location_of_minimum*(x) | Returns the last location of the minimal value of x. |
| *length*(x) | Returns the length of x |
| *linear_trend*(x, param) | Calculate a linear least-squares regression for the values of the time series versus the sequence from 0 to length of the time series minus one. |
| *longest_strike_above_mean*(x) | Returns the length of the longest consecutive subsequence in x that is bigger than the mean of x |
| *longest_strike_below_mean*(x) | Returns the length of the longest consecutive subsequence in x that is smaller than the mean of x |
| *max_langevin_fixed_point*(x, r, m) | Largest fixed point of dynamics :math:argmax_x {h(x)=0}' estimated from polynomial $h(x)$, which has been fitted to the deterministic dynamics of Langevin model |
| *maximum*(x) | Calculates the highest value of the time series x. |
| *mean*(x) | Returns the mean of x |
| *mean_abs_change*(x) | Returns the mean over the absolute differences between subsequent time series values which is |
| *mean_change*(x) | Returns the mean over the absolute differences between subsequent time series values which is |
| *mean_second_derivative_central*(x) | Returns the mean value of a central approximation of the second derivative |
| *median*(x) | Returns the median of x |
| *minimum*(x) | Calculates the lowest value of the time series x. |
| *number_crossing_m*(x, m) | Calculates the number of crossings of x on m. |
| *number_cwt_peaks*(x, n) | This feature calculator searches for different peaks in x. |
| *number_peaks*(x, n) | Calculates the number of peaks of at least support n in the time series x. |
| *partial_autocorrelation*(x, param) | Calculates the value of the partial autocorrelation function at the given lag. |
| *percentage_of_reoccurring_datapoints_to_all_datapoints*(x) | Returns the percentage of unique values, that are present in the time series more than once. |
| *percentage_of_reoccurring_values_to_all_values*(x) | Returns the ratio of unique values, that are present in the time series more than once. |
| *quantile*(x, q) | Calculates the q quantile of x. |
| *range_count*(x, min, max) | Count observed values within the interval [min, max). |
| *ratio_beyond_r_sigma*(x, r) | Ratio of values that are more than r*std(x) (so r sigma) away from the mean of x. |
| *ratio_value_number_to_time_series_length*(x) | Returns a factor which is 1 if all values in the time series occur only once, and below one if this is not the case. |
| *sample_entropy*(x) | Calculate and return sample entropy of x. |
| *set_property*(key, value) | This method returns a decorator that sets the property key of the function to value |
| *skewness*(x) | Returns the sample skewness of x (calculated with the adjusted Fisher-Pearson standardized moment coefficient G1). |
| *spkt_welch_density*(x, param) | This feature calculator estimates the cross power spectral density of the time series x at different frequencies. |
| *standard_deviation*(x) | Returns the standard deviation of x |

Continued on next page

Table 2 – continued from previous page

| | |
|---|---|
| *sum_of_reoccurring_data_points*(x) | Returns the sum of all data points, that are present in the time series more than once. |
| *sum_of_reoccurring_values*(x) | Returns the sum of all values, that are present in the time series more than once. |
| *sum_values*(x) | Calculates the sum over the time series values |
| *symmetry_looking*(x, param) | Boolean variable denoting if the distribution of x *looks symmetric*. |
| *time_reversal_asymmetry_statistic*(x, lag) | This function calculates the value of |
| *value_count*(x, value) | Count occurrences of *value* in time series x. |
| *variance*(x) | Returns the variance of x |
| *variance_larger_than_standard_deviation*(x) | Boolean variable denoting if the variance of x is greater than its standard deviation. |

submodule.

The following, exhaustive list contains all features that are calculated in the current version of *tsfresh*:

| | |
|---|---|
| *abs_energy*(x) | Returns the absolute energy of the time series which is the sum over the squared values |
| *absolute_sum_of_changes*(x) | Returns the sum over the absolute value of consecutive changes in the series x |
| *agg_autocorrelation*(x, param) | Calculates the value of an aggregation function f_agg (e.g. |
| *agg_linear_trend*(x, param) | Calculates a linear least-squares regression for values of the time series that were aggregated over chunks versus the sequence from 0 up to the number of chunks minus one. |
| *approximate_entropy*(x, m, r) | Implements a vectorized Approximate entropy algorithm. |
| *ar_coefficient*(x, param) | This feature calculator fits the unconditional maximum likelihood of an autoregressive AR(k) process. |
| *augmented_dickey_fuller*(x, param) | The Augmented Dickey-Fuller test is a hypothesis test which checks whether a unit root is present in a time series sample. |
| *autocorrelation*(x, lag) | Calculates the autocorrelation of the specified lag, according to the formula [1] |
| *binned_entropy*(x, max_bins) | First bins the values of x into max_bins equidistant bins. |
| *c3*(x, lag) | This function calculates the value of |
| *change_quantiles*(x, ql, qh, isabs, f_agg) | First fixes a corridor given by the quantiles ql and qh of the distribution of x. |
| *cid_ce*(x, normalize) | This function calculator is an estimate for a time series complexity [1] (A more complex time series has more peaks, valleys etc.). |
| *count_above_mean*(x) | Returns the number of values in x that are higher than the mean of x |
| *count_below_mean*(x) | Returns the number of values in x that are lower than the mean of x |
| *cwt_coefficients*(x, param) | Calculates a Continuous wavelet transform for the Ricker wavelet, also known as the "Mexican hat wavelet" which is defined by |

Continued on next page

Table 3 – continued from previous page

| | |
|---|---|
| `energy_ratio_by_chunks`(x, param) | Calculates the sum of squares of chunk i out of N chunks expressed as a ratio with the sum of squares over the whole series |
| `fft_aggregated`(x, param) | Returns the spectral centroid (mean), variance, skew, and kurtosis of the absolute fourier transform spectrum. |
| `fft_coefficient`(x, param) | Calculates the fourier coefficients of the one-dimensional discrete Fourier Transform for real input by fast fourier transformation algorithm |
| `first_location_of_maximum`(x) | Returns the first location of the maximum value of x. |
| `first_location_of_minimum`(x) | Returns the first location of the minimal value of x. |
| `friedrich_coefficients`(x, param) | Coefficients of polynomial $h(x)$, which has been fitted to the deterministic dynamics of Langevin model |
| `has_duplicate`(x) | Checks if any value in x occurs more than once |
| `has_duplicate_max`(x) | Checks if the maximum value of x is observed more than once |
| `has_duplicate_min`(x) | Checks if the minimal value of x is observed more than once |
| `index_mass_quantile`(x, param) | Those apply features calculate the relative index i where q% of the mass of the time series x lie left of i. |
| `kurtosis`(x) | Returns the kurtosis of x (calculated with the adjusted Fisher-Pearson standardized moment coefficient G2). |
| `large_standard_deviation`(x, r) | Boolean variable denoting if the standard dev of x is higher than 'r' times the range = difference between max and min of x. |
| `last_location_of_maximum`(x) | Returns the relative last location of the maximum value of x. |
| `last_location_of_minimum`(x) | Returns the last location of the minimal value of x. |
| `length`(x) | Returns the length of x |
| `linear_trend`(x, param) | Calculate a linear least-squares regression for the values of the time series versus the sequence from 0 to length of the time series minus one. |
| `longest_strike_above_mean`(x) | Returns the length of the longest consecutive subsequence in x that is bigger than the mean of x |
| `longest_strike_below_mean`(x) | Returns the length of the longest consecutive subsequence in x that is smaller than the mean of x |
| `max_langevin_fixed_point`(x, r, m) | Largest fixed point of dynamics :math:argmax_x {h(x)=0}' estimated from polynomial $h(x)$, which has been fitted to the deterministic dynamics of Langevin model |
| `maximum`(x) | Calculates the highest value of the time series x. |
| `mean`(x) | Returns the mean of x |
| `mean_abs_change`(x) | Returns the mean over the absolute differences between subsequent time series values which is |
| `mean_change`(x) | Returns the mean over the absolute differences between subsequent time series values which is |
| `mean_second_derivative_central`(x) | Returns the mean value of a central approximation of the second derivative |
| `median`(x) | Returns the median of x |
| `minimum`(x) | Calculates the lowest value of the time series x. |
| `number_crossing_m`(x, m) | Calculates the number of crossings of x on m. |
| `number_cwt_peaks`(x, n) | This feature calculator searches for different peaks in x. |

Continued on next page

Table 3 – continued from previous page

| | |
|---|---|
| *number_peaks*(x, n) | Calculates the number of peaks of at least support n in the time series x. |
| *partial_autocorrelation*(x, param) | Calculates the value of the partial autocorrelation function at the given lag. |
| *percentage_of_reoccurring_datapoints_to_all_datapoints*(x) | Returns the percentage of unique values, that are present in the time series more than once. |
| *percentage_of_reoccurring_values_to_all_values*(x) | Returns the ratio of unique values, that are present in the time series more than once. |
| *quantile*(x, q) | Calculates the q quantile of x. |
| *range_count*(x, min, max) | Count observed values within the interval [min, max). |
| *ratio_beyond_r_sigma*(x, r) | Ratio of values that are more than r*std(x) (so r sigma) away from the mean of x. |
| *ratio_value_number_to_time_series_length*(x) | Returns a factor which is 1 if all values in the time series occur only once, and below one if this is not the case. |
| *sample_entropy*(x) | Calculate and return sample entropy of x. |
| *set_property*(key, value) | This method returns a decorator that sets the property key of the function to value |
| *skewness*(x) | Returns the sample skewness of x (calculated with the adjusted Fisher-Pearson standardized moment coefficient G1). |
| *spkt_welch_density*(x, param) | This feature calculator estimates the cross power spectral density of the time series x at different frequencies. |
| *standard_deviation*(x) | Returns the standard deviation of x |
| *sum_of_reoccurring_data_points*(x) | Returns the sum of all data points, that are present in the time series more than once. |
| *sum_of_reoccurring_values*(x) | Returns the sum of all values, that are present in the time series more than once. |
| *sum_values*(x) | Calculates the sum over the time series values |
| *symmetry_looking*(x, param) | Boolean variable denoting if the distribution of x *looks symmetric*. |
| *time_reversal_asymmetry_statistic*(x, lag) | This function calculates the value of |
| *value_count*(x, value) | Count occurrences of *value* in time series x. |
| *variance*(x) | Returns the variance of x |
| *variance_larger_than_standard_deviation*(x) | Boolean variable denoting if the variance of x is greater than its standard deviation. |

# 1.7 Feature Calculation

## 1.7.1 Feature naming

tsfresh enforces a strict naming of the created features, which you have to follow whenever you create new feature calculators. This is due to the *tsfresh.feature_extraction.settings.from_columns()* method which needs to deduce the following information from the feature name

- the time series that was used to calculate the feature

- the feature calculator method that was used to derive the feature

- all parameters that have been used to calculate the feature (optional)

Hence, to enable the *tsfresh.feature_extraction.settings.from_columns()* to deduce all the necessary conditions, the features will be named in the following format

> {time_series_name}__{feature_name}__{parameter name 1}_{parameter value 1}__[..]__{parameter name k}_{parameter value k}

(Here we assumed that {feature_name} has k parameters).

### 1.7.2 Examples for feature naming

So for example the following feature name

> temperature_1__quantile__q_0.6

is the value of the feature *tsfresh.feature_extraction.feature_calculators.quantile()* for the time series `temperature_1` and a parameter value of q=0.6. On the other hand, the feature named

> Pressure 5__cwt_coefficients__widths_(2, 5, 10, 20)__coeff_14__w_5

denotes the value of the feature *tsfresh.feature_extraction.feature_calculators.cwt_coefficients()* for the time series `Pressure 5` under parameter values of widths=(2, 5, 10, 20), coeff=14 and w=5.

## 1.8 Feature extraction settings

When starting a new data science project involving time series you probably want to start by extracting a comprehensive set of features. Later you can identify which features are relevant for the task at hand. In the final stages, you probably want to fine tune the parameter of the features to fine tune your models.

You can do all those things with tsfresh. So, you need to know how to control which features are calculated by tsfresh and how one can adjust the parameters. In this section, we will clarify this.

### 1.8.1 For the lazy: Just let me calculate some features

So, to just calculate a comprehensive set of features, call the tsfresh.extract_features() method without passing a *default_fc_parameters* or *kind_to_fc_parameters* object, which means you are using the default options (which will use all feature calculators in this package for what we think are sane default parameters).

### 1.8.2 For the advanced: How does I set the parameters for all kind of time series?

After digging deeper into your data, you maybe want to calculate more of a certain type of feature and less of another type. So, you need to use custom settings for the feature extractors. To do that with tsfresh you will have to use a custom settings object:

```
>>> from tsfresh.feature_extraction import ComprehensiveFCParameters
>>> settings = ComprehensiveFCParameters()
>>> # Set here the options of the settings object as shown in the paragraphs below
>>> # ...
>>> from tsfresh.feature_extraction import extract_features
>>> extract_features(df, default_fc_parameters=settings)
```

The *default_fc_parameters* is expected to be a dictionary, which maps feature calculator names (the function names you can find in the *tsfresh.feature_extraction.feature_calculators* file) to a list of dictionaries,

which are the parameters with which the function will be called (as key value pairs). Each function parameter combination, that is in this dict will be called during the extraction and will produce a feature. If the function does not take any parameters, the value should be set to *None*.

For example

```
fc_parameters = {
    "length": None,
    "large_standard_deviation": [{"r": 0.05}, {"r": 0.1}]
}
```

will produce three features: one by calling the `tsfresh.feature_extraction.feature_calculators.length()` function without any parameters and two by calling `tsfresh.feature_extraction.feature_calculators.large_standard_deviation()` with *r = 0.05* and *r = 0.1*.

So you can control, which features will be extracted, by adding/removing either keys or parameters from this dict. It is as easy as that. If you decide to not calculate the length feature here, you delete it from the dictionary:

```
del fc_parameters["length"]
```

And now, only the two other features are calculated.

For convenience, three dictionaries are predefined and can be used right away:

- `tsfresh.feature_extraction.settings.ComprehensiveFCParameters`: includes all features without parameters and all features with parameters, each with different parameter combinations. This is the default for *extract_features* if you do not hand in a *default_fc_parameters* at all.

- `tsfresh.feature_extraction.settings.MinimalFCParameters`: includes only a handful of features and can be used for quick tests. The features which have the "minimal" attribute are used here.

- `tsfresh.feature_extraction.settings.EfficientFCParameters`: Mostly the same features as in the `tsfresh.feature_extraction.settings.ComprehensiveFCParameters`, but without features which are marked with the "high_comp_cost" attribute. This can be used if runtime performance plays a major role.

Theoretically, you could calculate an unlimited number of features with tsfresh by adding entry after entry to the dictionary.

### 1.8.3 For the ambitious: How do I set the parameters for different type of time series?

It is also possible, to control the features to be extracted for the different kinds of time series individually. You can do so by passing another dictionary to the extract function as a

*kind_to_fc_parameters* = {"kind" : *fc_parameters*}

parameter. This dict must be a mapping from kind names (as string) to *fc_parameters* objects, which you would normally pass as an argument to the *default_fc_parameters* parameter.

So, for example using

```
kind_to_fc_parameters = {
    "temperature": {"mean": None},
    "pressure": {"max": None, "min": None}
}
```

will extract the *"mean"* feature of the *"temperature"* time series and the *"min"* and *"max"* of the *"pressure"* time series.

---

The *kind_to_fc_parameters* argument will partly override the *default_fc_parameters*. So, if you include a kind name in the *kind_to_fc_parameters* parameter, its value will be used for that kind. Other kinds will still use the *default_fc_parameters*.

### 1.8.4 A handy trick: Do I really have to create the dictionary by hand?

Not necessarily. let's assume you have a DataFrame of tsfresh features. By using feature selection algorithms you find out that only a subgroup of features is relevant.

Then, we provide the `tsfresh.feature_extraction.settings.from_columns()` method that constructs the *kind_to_fc_parameters* dictionary from the column names of this filtered feature matrix to make sure that only relevant features are extracted.

This can save a huge amount of time because you prevent the calculation of uncessary features. Let's illustrate that with an example:

```
# X_tsfresh containes the extracted tsfresh features
X_tsfresh = extract_features(...)

# which are now filtered to only contain relevant features
X_tsfresh_filtered = some_feature_selection(X_tsfresh, y, ....)

# we can easily construct the corresponding settings object
kind_to_fc_parameters = tsfresh.feature_extraction.settings.from_columns(X_tsfresh_
↪filtered)
```

this will construct you the *kind_to_fc_parameters* dictionary that corresponds to the features and parameters (!) from the tsfresh features that were filtered by the *some_feature_selection* feature selection algorithm.

## 1.9 Feature filtering

The all-relevant problem of feature selection is the identification of all strongly and weakly relevant attributes. This problem is especially hard to solve for time series classification and regression in industrial applications such as predictive maintenance or production line optimization, for which each label or regression target is associated with several time series and meta-information simultaneously.

To limit the number of irrelevant features, tsfresh deploys the fresh algorithm (fresh stands for *FeatuRe Extraction based on Scalable Hypothesis tests*)[1].

The algorithm is called by `tsfresh.feature_selection.relevance.calculate_relevance_table()`. It is an efficient, scalable feature extraction algorithm, which filters the available features in an early stage of the machine learning pipeline with respect to their significance for the classification or regression task, while controlling the expected percentage of selected but irrelevant features.

The filtering process consists of three phases which are sketched in the following figure:

---

[1] Christ, M., Kempa-Liehr, A.W. and Feindt, M. (2016). Distributed and parallel time series feature extraction for industrial big data applications. ArXiv e-prints: 1610.07717 URL: http://adsabs.harvard.edu/abs/2016arXiv161007717C

### 1.9.1 Phase 1 - Feature extraction

Firstly, the algorithm characterizes time series with comprehensive and well-established feature mappings and considers additional features describing meta-information. The feature calculators used to derive the features are contained in `tsfresh.feature_extraction.feature_calculators`.

In the figure from above, this corresponds to the change from raw time series to aggregated features.

### 1.9.2 Phase 2 - Feature significance testing

In a second step, each feature vector is individually and independently evaluated with respect to its significance for predicting the target under investigation. Those tests are contained in the submodule `tsfresh.feature_selection.significance_tests`. The result of these tests is a vector of p-values, quantifying the significance of each feature for predicting the label/target.

In the figure from above, this corresponds to the change from aggregated features to p-values.

### 1.9.3 Phase 3 - Multiple test procedure

The vector of p-values is evaluated on basis of the Benjamini-Yekutieli procedure[2] in order to decide which features to keep. This multiple testing procedure is contained in the submodule `tsfresh.feature_selection.benjamini_hochberg_test`.

In the figure from above, this corresponds to the change from p-values to selected features.

### 1.9.4 References

---

[2] Benjamini, Y. and Yekutieli, D. (2001). The control of the false discovery rate in multiple testing under dependency. Annals of statistics, 1165–1188

# 1.10 How to add a custom feature

It may be beneficial to add a custom feature to those that are calculated by tsfresh. To do so, one has to follow four simple steps:

## 1.10.1 Step 1. Decide which type of feature you want to implement

In tsfresh we differentiate between two types of feature calculation methods

> *1.* simple

> *2.* combiner

The difference lays in the number of calculated features for a singular time series. The feature_calculator returns either one (*1.*) or multiple features (*2.*). So if you want to add a singular feature stick with *1.*, the simple feature calculator class. If it is beneficial to calculate multiples features at the same time (to e.g. perform auxiliary calculations only once for all features), stick with type *2.*.

## 1.10.2 Step 2. Write the feature calculator

Depending on which type of feature you are implementing, you can use the following feature calculator skeletons:

*1.* simple features

You can write such a simple feature calculator, that returns exactly one feature, without parameter

```python
@set_property("fctype", "simple")
def your_feature_calculator(x):
    """
    The description of your feature

    :param x: the time series to calculate the feature of
    :type x: pandas.Series
    :return: the value of this feature
    :return type: bool, int or float
    """
    # Calculation of feature as float, int or bool
    f = f(x)
    return f
```

or with parameter

```python
@set_property("fctype", "simple"")
def your_feature_calculator(x, p1, p2, ...):
    """
    Description of your feature

    :param x: the time series to calculate the feature of
    :type x: pandas.Series
    :param p1: description of your parameter p1
    :type p1: type of your parameter p1
    :param p2: description of your parameter p2
    :type p2: type of your parameter p2
    ...
    :return: the value of this feature
    :return type: bool, int or float
```

(continues on next page)

```
    """
    # Calculation of feature as float, int or bool
    f = f(x)
    return f
```

*2.* combiner features

```python
@set_property("fctype", "combiner")
def your_feature_calculator(x, param):
    """
    Description of your feature

    :param x: the time series to calculate the feature of
    :type x: pandas.Series
    :param c: the time series name
    :type c: str
    :param param: contains dictionaries {"p1": x, "p2": y, ...} with p1 float, p2 int
↪...
    :type param: list
    :return: list of tuples (s, f) where s are the parameters, serialized as a string,
↪ and f the respective feature
        value as bool, int or float
    :return type: pandas.Series
    """
    # s is a function that serializes the config
    # f is a function that calculates the feature value for the config
    return [(s(config), f(x, config)) for config in param]
```

After implementing the feature calculator, please add it to the *tsfresh.feature_extraction.feature_calculators* submodule. tsfresh will only find feature calculators that are in this submodule.

## 1.10.3 Step 3. Add custom settings for your feature

Finally, you have to add custom settings if your feature is a simple or combiner feature with parameters. To do so, just append your feature with sane default parameters to the `name_to_param` dictionary inside the *tsfresh.feature_extraction.settings.ComprehensiveFCParameters* constructor:

```python
name_to_param.update({
    # here are the existing settings
    ...
    # Now the settings of your feature calculator
    "your_feature_calculator" = [{"p1": x, "p2": y, ...} for x,y in ...],
})
```

That is it, tsfresh will calculate your feature the next time you run it.

Please make sure, that the different feature extraction settings (e.g. *tsfresh.feature_extraction.settings.EfficientFCParameters*, *tsfresh.feature_extraction.settings.MinimalFCParameters* or *tsfresh.feature_extraction.settings.ComprehensiveFCParameters*) do include different sets of feature calculators to use. You can control, which feature extraction settings object will include your new feature calculator by giving your function attributes like "minimal" or "high_comp_cost". Please see the classes in *tsfresh.feature_extraction.settings* for more information.

### 1.10.4 Step 4. Add a pull request

We would very happy if you contribute your implemented features to tsfresh. So make sure to create a pull request at our github page. We happily accept partly implemented feature calculators, which we can finalize collaboratively.

## 1.11 Parallelization

The feature extraction as well as the feature selection offer the possibility of parallelization. Out of the box both tasks are parallelized by tsfresh. However, the overhead introduced with the parallelization should not be underestimated. Here we discuss the different settings to control the parallelization. To achieve best results for your use-case you should experiment with the parameters.

Please let us know about your results tuning the below mentioned parameters! It will help improve this document as well as the default settings.

### 1.11.1 Parallelization of Feature Selection

We use a `multiprocessing.Pool` to parallelize the calculation of the p-values for each feature. On instantiation we set the Pool's number of worker processes to *n_jobs*. This field defaults to the number of processors on the current system. We recommend setting it to the maximum number of available (and otherwise idle) processors.

The chunksize of the Pool's map function is another important parameter to consider. It can be set via the *chunksize* field. By default it is up to `multiprocessing.Pool` is parallelisation parameter. One data chunk is defined as a singular time series for one id and one kind. The chunksize is the number of chunks that are submitted as one task to one worker process. If you set the chunksize to 10, then it means that one worker task corresponds to calculate all features for 10 id/kind time series combinations. If it is set it to None, depending on distributor, heuristics are used to find the optimal chunksize. The chunksize can have an crucial influence on the optimal cluster performance and should be optimised in benchmarks for the problem at hand.

### 1.11.2 Parallelization of Feature Extraction

For the feature extraction tsfresh exposes the parameters *n_jobs* and *chunksize*. Both behave analogue to the parameters for the feature selection.

To do performance studies and profiling, it sometimes quite useful to turn off parallelization at all. This can be setting the parameter *n_jobs* to 0.

## 1.12 How to deploy tsfresh at scale

The high volume of time series data can demand an analysis at scale. So, time series need to be processed on a group of computational units instead of a singular machine.

Accordingly, it may be necessary to distribute the extraction of time series features to a cluster. Indeed, it is possible to extract features with *tsfresh* in a distributed fashion. This page will explain how to setup a distributed *tsfresh*.

### 1.12.1 The distributor class

To distribute the calculation of features, we use a certain object, the Distributor class (contained in the `tsfresh.utilities.distribution` module).

Essentially, a Distributor organizes the application of feature calculators to data chunks. It maps the feature calculators to the data chunks and then reduces them, meaning that it combines the results of the individual mapping into one object, the feature matrix.

So, Distributor will, in the following order,

1. calculates an optimal `chunk_size`, based on the characteristics of the time series data at hand (by `calculate_best_chunk_size()`)

2. split the time series data into chunks (by `partition()`)

3. distribute the applying of the feature calculators to the data chunks (by `distribute()`)

4. combine the results into the feature matrix (by `map_reduce()`)

5. close all connections, shutdown all resources and clean everything (by `close()`)

So, how can you use such a Distributor to extract features with *tsfresh*? You will have to pass it into as the `distributor` argument to the `extract_features()` method.

The following example shows how to define the MultiprocessingDistributor, which will distribute the calculations to a local pool of threads:

```python
from tsfresh.examples.robot_execution_failures import \
    download_robot_execution_failures, \
    load_robot_execution_failures
from tsfresh.feature_extraction import extract_features
from tsfresh.utilities.distribution import MultiprocessingDistributor

# download and load some time series data
download_robot_execution_failures()
df, y = load_robot_execution_failures()

# We construct a Distributor that will spawn the calculations
# over four threads on the local machine
Distributor = MultiprocessingDistributor(n_workers=4,
                                         disable_progressbar=False,
                                         progressbar_title="Feature Extraction")

# just to pass the Distributor object to
# the feature extraction, along the other parameters
X = extract_features(timeseries_container=df,
                     column_id='id', column_sort='time',
                     distributor=Distributor)
```

This example actually corresponds to the existing multiprocessing *tsfresh* API, where you just specify the number of jobs, without the need to construct the Distributor:

```python
from tsfresh.examples.robot_execution_failures import \
    download_robot_execution_failures, \
    load_robot_execution_failures
from tsfresh.feature_extraction import extract_features

download_robot_execution_failures()
df, y = load_robot_execution_failures()

X = extract_features(timeseries_container=df,
                     column_id='id', column_sort='time',
                     n_jobs=4)
```

### 1.12.2 Using dask to distribute the calculations

We provide distributor for the dask framework, where *"Dask is a flexible parallel computing library for analytic computing."*

Dask is a great framework to distribute analytic calculations to a cluster. It scales up and down, meaning that you can even use it on a singular machine. The only thing that you will need to run *tsfresh* on a Dask cluster is the ip address and port number of the dask-scheduler.

Lets say that your dask scheduler is running at `192.168.0.1:8786`, then we can easily construct a `ClusterDaskDistributor` that connects to the sceduler and distributes the time series data and the calculation to a cluster:

```python
from tsfresh.examples.robot_execution_failures import \
    download_robot_execution_failures, \
    load_robot_execution_failures
from tsfresh.feature_extraction import extract_features
from tsfresh.utilities.distribution import ClusterDaskDistributor

download_robot_execution_failures()
df, y = load_robot_execution_failures()

Distributor = ClusterDaskDistributor(address="192.168.0.1:8786")

X = extract_features(timeseries_container=df,
                     column_id='id', column_sort='time',
                     distributor=Distributor)
```

Compared to the `MultiprocessingDistributor` example from above, we only had to change one line to switch from one machine to a whole cluster. It is as easy as that. By changing the Distributor you can easily deploy your application to run to a cluster instead of your workstation.

You can also use a local DaskCluster on your local machine to emulate a Dask network. The following example shows how to setup a `LocalDaskDistributor` on a local cluster of 3 workers:

```python
from tsfresh.examples.robot_execution_failures import \
    download_robot_execution_failures, \
    load_robot_execution_failures
from tsfresh.feature_extraction import extract_features
from tsfresh.utilities.distribution import LocalDaskDistributor

download_robot_execution_failures()
df, y = load_robot_execution_failures()

Distributor = LocalDaskDistributor(n_workers=3)

X = extract_features(timeseries_container=df,
                     column_id='id', column_sort='time',
                     distributor=Distributor)
```

### 1.12.3 Writing your own distributor

If you want to user another framework than Dask, you will have to write your own Distributor. To construct your custom Distributor, you will have to define an object that inherits from the abstract base class `tsfresh.utilities.distribution.DistributorBaseClass`. The `tsfresh.utilities.distribution` module contains more information about what you will need to implement.

## 1.13 Time series forecasting

Features that are extracted with *tsfresh* can be used for many different tasks, such as time series classification, compression or forecasting. This section explains how one can use the features for time series forecasting tasks.

The "sort" column of a DataFrame in the supported *Data Formats* gives a sequential state to the individual measurements. In the case of time series this can be the *time* dimension while in the case of spectra the order is given by the *wavelength* or *frequency* dimensions. We can exploit this sequence to generate more input data out of a single time series, by *rolling* over the data.

Lets say you have the price of a certain stock, e.g. Apple, for 100 time steps. Now, you want to build a feature-based model to forecast future prices of the Apple stock. So you will have to extract features in every time step of the original time series while looking at a certain number of past values. A rolling mechanism will give you the sub time series of last *m* time steps to construct the features.

The following image illustrates the process:



So, we move the window that extract the features and then predict the next time step (which was not used to extract features) forward. In the above image, the window moves from left to right.

Another example can be found in streaming data, e.g. in Industry 4.0 applications. Here you typically get one new data row at a time and use this to for example predict machine failures. To train your model, you could act as if you would stream the data, by feeding your classifier the data after one time step, the data after the first two time steps etc.

Both examples imply, that you extract the features not only on the full data set, but also on all temporal coherent subsets of data, which is the process of *rolling*. In tsfresh, this is implemented in the function `tsfresh.utilities.dataframe_functions.roll_time_series()`. Further, we provide the `tsfresh.utilities.dataframe_functions.make_forecasting_frame()` method as a convenient wrapper to fast construct the container and target vector for a given sequence.

## 1.13.1 The rolling mechanism

The rolling mechanism takes a time series $x$ with its data rows $[x_1, x_2, x_3, ..., x_n]$ and creates $n$ new time series $\hat{x}^k$, each of them with a different consecutive part of $x$:

$$\hat{x}^k = [x_k, x_{k-1}, x_{k-2}, ..., x_1]$$

To see what this does in real-world applications, we look into the following example flat DataFrame in tsfresh format

| id | time | x | y |
|----|------|----|----|
| 1 | t1 | 1 | 5 |
| 1 | t2 | 2 | 6 |
| 1 | t3 | 3 | 7 |
| 1 | t4 | 4 | 8 |
| 2 | t8 | 10 | 12 |
| 2 | t9 | 11 | 13 |

where you have measured the values from two sensors x and y for two different entities (id 1 and 2) in 4 or 2 time steps (t1 to t9).

Now, we can use `tsfresh.utilities.dataframe_functions.roll_time_series()` to get consecutive sub-time series. E.g. if you set *rolling* to 0, the feature extraction works on the original time series without any rolling.

So it extracts 2 set of features,

| id | time | x | y |
|----|------|----|----|
| 1 | t1 | 1 | 5 |
| 1 | t2 | 2 | 6 |
| 1 | t3 | 3 | 7 |
| 1 | t4 | 4 | 8 |

and

| id | time | x | y |
|----|------|----|----|
| 2 | t8 | 10 | 12 |
| 2 | t9 | 11 | 13 |

If you set rolling to 1, the feature extraction works with all of the following time series:

| id | time | x | y |
|----|------|---|---|
| 1 | t1 | 1 | 5 |

| id | time | x | y |
|----|------|---|---|
| 1 | t1 | 1 | 5 |
| 1 | t2 | 2 | 6 |

| id | time | x | y |
|----|------|----|----|
| 1 | t1 | 1 | 5 |
| 1 | t2 | 2 | 6 |
| 1 | t3 | 3 | 7 |
| 2 | t8 | 10 | 12 |

| id | time | x | y |
|----|------|----|----|
| 1 | t1 | 1 | 5 |
| 1 | t2 | 2 | 6 |
| 1 | t3 | 3 | 7 |
| 1 | t4 | 4 | 8 |
| 2 | t8 | 10 | 12 |
| 2 | t9 | 11 | 13 |

If you set rolling to -1, you end up with features for the time series, rolled in the other direction

| id | time | x | y |
|----|------|----|----|
| 1 | t4 | 4 | 8 |

| id | time | x | y |
|----|------|----|----|
| 1 | t3 | 3 | 7 |
| 1 | t4 | 4 | 8 |

| id | time | x | y |
|----|------|----|----|
| 1 | t2 | 2 | 6 |
| 1 | t3 | 3 | 7 |
| 1 | t4 | 4 | 8 |
| 2 | t9 | 11 | 13 |

| id | time | x | y |
|----|------|----|----|
| 1 | t1 | 1 | 5 |
| 1 | t2 | 2 | 6 |
| 1 | t3 | 3 | 7 |
| 1 | t4 | 4 | 8 |
| 2 | t8 | 10 | 12 |
| 2 | t9 | 11 | 13 |

We only gave an example for the flat DataFrame format, but rolling actually works on all 3 *Data Formats* that are supported by tsfresh.

This process is also visualized by the following figure. It shows how the purple, rolled sub-timeseries are used as base for the construction of the feature matrix *X* (after calculation of the features by *f*). The green data points need to be predicted by the model and are used as rows in the target vector *y*.

## 1.13.2 Parameters and Implementation Notes

The above example demonstrates the overall rolling mechanism, which creates new time series. Now we discuss the naming convention for such new time series:

For identifying every subsequence, tsfresh uses the time stamp of the point that will be predicted as new "id". The above example with rolling set to 1 yields the following sub-time series:

| id | time | x | y |
|----|------|---|---|
| t1 | t1 | 1 | 5 |

| id | time | x | y |
|----|------|---|---|
| t2 | t1 | 1 | 5 |
| t2 | t2 | 2 | 6 |

| id | time | x | y |
|----|------|---|---|
| t3 | t1 | 1 | 5 |
| t3 | t2 | 2 | 6 |
| t3 | t3 | 3 | 7 |

| id | time | x | y |
|----|------|---|---|
| t4 | t1 | 1 | 5 |
| t4 | t2 | 2 | 6 |
| t4 | t3 | 3 | 7 |
| t4 | t4 | 4 | 8 |

| id | time | x | y |
|----|------|----|----|
| t8 | t8 | 10 | 12 |

| id | time | x | y |
|----|------|----|----|
| t9 | t8 | 10 | 12 |
| t9 | t9 | 11 | 13 |

The new id is the time stamp where the shift ended. So above, every table represents a sub-time series. The higher the shift value, the more steps the time series was moved into the specified direction (into the past in this example).

If you want to limit how far the time series shall be shifted into the specified direction, you can set the *max_timeshift* parameter to the maximum time steps to be shifted. In our example, setting *max_timeshift* to 1 yields the following result (setting it to 0 will create all possible shifts):

| id | time | x | y |
|----|------|----|----|
| t1 | t1 | 1 | 5 |

| id | time | x | y |
|----|------|----|----|
| t2 | t1 | 1 | 5 |
| t2 | t2 | 2 | 6 |

| id | time | x | y |
|----|------|----|----|
| t3 | t2 | 2 | 6 |
| t3 | t3 | 3 | 7 |

| id | time | x | y |
|----|------|----|----|
| t4 | t3 | 3 | 7 |
| t4 | t4 | 4 | 8 |

| id | time | x | y |
|----|------|----|----|
| t8 | t8 | 10 | 12 |

| id | time | x | y |
|----|------|----|----|
| t9 | t8 | 10 | 12 |
| t9 | t9 | 11 | 13 |

## 1.14 FAQ

1. **Does tsfresh support different time series lengths?**

   Yes, it supports different time series lengths. However, some feature calculators can demand a minimal length of the time series. If a shorter time series is passed to the calculator, a NaN is returned for those features.

2. **Is it possible to extract features from rolling/shifted time series?**

   Yes, the `tsfresh.dataframe_functions.roll_time_series()` function allows to conviniently create a rolled time series datframe from your data. You just have to transform your data into one of the supported

tsfresh *Data Formats*. Then, the `tsfresh.dataframe_functions.roll_time_series()` give you a DataFrame with the rolled time series, that you can pass to tsfresh. On the following page you can find a detailed description: *Time series forecasting*.

3. **How can I use tsfresh with windows?**

   We recommend to use Anaconda. After installing, open the Anaconda Prompt, create an environment and set up tsfresh (Please be aware that we're using multiprocessing, which can be problematic.):

   ```
   conda create -n ENV_NAME python=VERSION
   conda install -n ENV_NAME pip requests numpy pandas scipy statsmodels patsy␣
   ↪scikit-learn future six tqdm
   activate ENV_NAME
   pip install tsfresh
   ```

4. **Does tsfresh support different sampling rates in the time series?**

   Yes! The feature calculators in tsfresh do not care about the sampling frequency. You will have to use the second input format, the stacked DataFramed (see *Data Formats*)

## 1.15 Authors

### 1.15.1 Core Development Team

- Maximilian Christ (maximilianchrist.com, max.christ@me.com)
- Nils Braun (nilslennartbraun@gmail.com)
- Julius Neuffer (julius.neuffer@blue-yonder.com)

### 1.15.2 Contributions

- Andreas W. Kempa-Liehr
- Markus Frey
- Niklas Haas
- earthgecko
- Moritz Gelb
- Thibault de Boissiere
- Brian Sang
- Stephan Müller
- Vin Tang
- Chris Chow
- Ezekiel Kruglick
- Timo Klerx
- Gregor Koehler
- Matúš Tomlein
- Florian Aspart

- Sergey Shepelev
- Justin White
- 10. Kleint

## 1.16 License

```
MIT LICENCE

Copyright (c) 2016 Maximilian Christ, Blue Yonder GmbH

Permission is hereby granted, free of charge, to any person obtaining a copy of this␣
→software and associated
documentation files (the "Software"), to deal in the Software without restriction,␣
→including without limitation the
rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell␣
→copies of the Software, and to permit
persons to whom the Software is furnished to do so, subject to the following␣
→conditions:

The above copyright notice and this permission notice shall be included in all copies␣
→or substantial portions of the
Software.

THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED,␣
→INCLUDING BUT NOT LIMITED TO THE
WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT.␣
→IN NO EVENT SHALL THE AUTHORS OR
COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN␣
→ACTION OF CONTRACT, TORT OR
OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR␣
→OTHER DEALINGS IN THE SOFTWARE.
```

## 1.17 Changelog

tsfresh uses Semantic Versioning

### 1.17.1 Version 0.11.1

- general performance improvements
- removed hard pinning of dependencies
- **fixed**
    - the stock price forecasting notebook
    - the multi classification notebook

### 1.17.2 Version 0.11.0

- **new feature calculators:**

- – fft_aggregated

    – cid_ce

- renamed mean_second_derivate_central to mean_second_derivative_central

- add warning if no relevant features were found in feature selection

- add columns_to_ignore parameter to from_columns method

- add distribution module, contains support for distributed feature extraction on Dask

### 1.17.3 Version 0.10.1

- split test suite into unit and integration tests

- **fixed the following bugs**

    – use name of value column as time series kind

    – prevent the spawning of subprocesses which lead to high memory consumption

    – fix deployment from travis to pypi

### 1.17.4 Version 0.10.0

- **new feature calculators:**

    – partial autocorrelation

- added list of calculated features to documentation

- **added two ipython notebooks to**

    – illustrate PCA on features

    – illustrate the Benjamini Yekutieli procedure

- **fixed the following bugs**

    – improperly quotation of dickey fuller settings

### 1.17.5 Version 0.9.0

- **new feature calculators:**

    – ratio_beyond_r_sigma

    – energy_ratio_by_chunks

    – number_crossing_m

    – c3

    – angle & abs for fft coefficients

    – agg_autocorrelation

    – p-Value and usedLag for augmented_dickey_fuller

    – change_quantiles

- **changed the calculation of the following features:**

- fft_coefficients

- autocorrelation

- time_reversal_asymmetry_statistic

- **removed the following feature calculators:**

  - large_number_of_peak

  - mean_autocorrelation

  - mean_abs_change_quantiles

- add support for multi classification in the feature selection

- improved description of the rolling mechanism

- added function make_forecasting_frame method for forecasting tasks

- internally ditched the pandas representation of the time series, yielding drastic speed improvements

- replaced feature calculator types from aggregate/aggregate with parameter/apply to simple/combiner

- add test for the ipython notebooks

- added notebook to inspect dft features

- make sure that RelevantFeatureAugmentor always imputes

- **fixed the following bugs**

  - impute was replacing whole columns by mean

  - fft coefficient were only calculated on truncated part

  - allow to suppress warnings from impute function

  - added missing lag in time_reversal_asymmetry_statistic

### 1.17.6 Version 0.8.1

- **new features:**

  - linear trend

  - agg trend

- **new sklearn compatible transformers**

  - PerColumnImputer

- **fixed bugs**

  - make mannwhitneyu method compatible with scipy > v0.18.0

- added caching to travis

- internally, added serial calculation of features

### 1.17.7 Version 0.8.0

- **Breaking API changes:**

  - removing of feature extraction settings object, replaced by keyword arguments and a plain dictionary (fc_parameters)

- removing of feature selection settings object, replaced by keyword arguments

- added notebook with examples of new API

- added chapter in docs about the new API

- adjusted old notebooks and documentation to new API

### 1.17.8 Version 0.7.1

- added a maximum shift parameter to the rolling utility

- added a FAQ entry about how to use tsfresh on windows

- **drastically decreased the runtime of the following features**

  - cwt_coefficient

  - index_mass_quantile

  - number_peaks

  - large_standard_deviation

  - symmetry_looking

- removed baseline unit tests

- **bugfixes:**

  - per sample parallel imputing was done on chunks which gave non deterministic results

  - imputing on dtypes other that float32 did not work properly

- several improvements to documentation

### 1.17.9 Version 0.7.0

- new rolling utility to use tsfresh for time series forecasting tasks

- **bugfixes:**

  - index_mass_quantile was using global index of time series container

  - an index with same name as id_column was breaking parallelization

  - friedrich_coefficients and max_langevin_fixed_point were occasionally stalling

### 1.17.10 Version 0.6.0

- progress bar for feature selection

- new feature: estimation of largest fixed point of deterministic dynamics

- new notebook: demonstration how to use tsfresh in a pipeline with train and test datasets

- remove no logging handler warning

- fixed bug in the RelevantFeatureAugmenter regarding the evaluate_only_added_features parameters

### 1.17.11 Version 0.5.0

- new example: driftbif simulation
- further improvements of the parallelization
- language improvements in the documentation
- performance improvements for some features
- performance improvements for the impute function
- new feature and feature renaming: sum_of_recurring_values, sum_of_recurring_data_points

### 1.17.12 Version 0.4.0

- fixed several bugs: checking of UCI dataset, out of index error for mean_abs_change_quantiles
- added a progress bar denoting the progress of the extraction process
- added parallelization per sample
- added unit tests for comparing results of feature extraction to older snapshots
- added "high_comp_cost" attribute
- added ReasonableFeatureExtraction settings only calculating features without "high_comp_cost" attribute

### 1.17.13 Version 0.3.1

- fixed several bugs: closing multiprocessing pools / index out of range cwt calculator / division by 0 in index_mass_quantile
- now all warnings are disabled by default
- for a singular type time series data, the name of value column is used as feature prefix

### 1.17.14 Version 0.3.0

- fixed bug with parsing of "NUMBER_OF_CPUS" environment variable
- now features are calculated in parallel for each type

### 1.17.15 Version 0.2.0

- now p-values are calculated in parallel
- fixed bugs for constant features
- allow time series columns to be named 0
- moved uci repository datasets to github mirror
- added feature calculator sample_entropy
- added MinimalFeatureExtraction settings
- fixed bug in calculation of fourier coefficients

---

### 1.17.16 Version 0.1.2

- added support for python 3.5.2
- fixed bug with the naming of the features that made the naming of features non-deterministic

### 1.17.17 Version 0.1.1

- mainly fixes for the read-the-docs documentation, the pypi readme and so on

### 1.17.18 Version 0.1.0

- Initial version :)

## 1.18 How to contribute

We want tsfresh to become the biggest archive of feature extraction methods in python. To achieve this goal, we need your help!

All contributions, bug reports, bug fixes, documentation improvements, enhancements and ideas are welcome. If you want to add one or two interesting feature calculators, implement a new feature selection process or just fix 1-2 typos, your help is appreciated.

If you want to help, just create a pull request on our github page. To the new user, working with Git can sometimes be confusing and frustrating. If you are not familiar with Git you can also contact us by *email*.

### 1.18.1 Guidelines

There are three general coding paradigms that we believe in:

1. **Keep it simple**. We believe that *"Programs should be written for people to read, and only incidentally for machines to execute."*.

2. **Keep it documented** by at least including a docstring for each method and class. Do not describe what you are doing but why you are doing it.

3. **Keep it tested**. We aim for a high test coverage.

There are two important copyright guidelines:

4. Please do not include any data sets for which a licence is not available or commercial use is even prohibited. Those can undermine the licence of the whole projects.

5. Do not use code snippets for which a licence is not available (e.g. from stackoverflow) or commercial use is even prohibited. Those can undermine the licence of the whole projects.

Further, there are some technical decisions we made:

6. Clear the Output of iPython notebooks. This improves the readability of related Git diffs.

### 1.18.2 Test framework

After making your changes, you probably want to test your changes locally. To run our comprehensive suite of unit tests you have to install all the relevant python packages with

```
cd /path/to/tsfresh
pip install -r requirements.txt
pip install -r rdocs-requirements.txt
pip install -r test-requirements.txt
pip install -e .
```

The last command will dynamically link the tsfresh package which means that changes to the code will directly show up for example in your test run.

Then, if you have everything installed, you can run the tests with

```
python setup.py test
```

or build the documentation with

```
python setup.py docs
```

The finished documentation can be found in the docs/_build/html folder.

On Github we use a Travis CI Folder that runs our test suite every time a commit or pull request is sent. The configuration of Travi is controlled by the .travis.yml file.

We are looking forward to hear from you! =)

# Indices and tables

- genindex
- modindex
- search

# CHAPTER 3

## Acknowledgements

# Python Module Index

# Index

## M

## N

## P

## Q

## R

# V