

---

# **Accelerated exploration of multinary systems**

*Release 1.1*

**Elise Garel, Jean-Luc Parouty**

**Feb 21, 2022**



# CONTENTS

<b>1</b>	<b>Experiments Plannification</b>	<b>3</b>
1.1	Principle . . . . .	3
1.2	Modules . . . . .	4
<b>2</b>	<b>pyterk package</b>	<b>11</b>
2.1	Module contents . . . . .	11
2.2	Submodules . . . . .	12
2.3	pyterk.config module . . . . .	12
2.4	pyterk.models module . . . . .	13
2.5	pyterk.reporter module . . . . .	13
2.6	pyterk.task_manager module . . . . .	16
2.7	pyterk.worker module . . . . .	18
<b>3</b>	<b>MultipleRegression module</b>	<b>19</b>
3.1	Description . . . . .	19
3.2	Functions . . . . .	19
<b>4</b>	<b>Indices and tables</b>	<b>21</b>
	<b>Python Module Index</b>	<b>23</b>
	<b>MATLAB Module Index</b>	<b>25</b>
	<b>Index</b>	<b>27</b>



This documentation aims at helping use the codes developped for “Accelerated exploration of multinary systems using high-throughput experiments and machine learning” project.

This project combines material Science and AI, and consists in producing experimental data using combinatorial high-throughput methods and mixture design, and then to analyse them using Machine Learning tools to extract predictive models linking compositions, structures and properties.

### **Experiments Planification**

This Matlab GUI interface allow to prepare the experimental work. It gives a set of 1D linear gradients or 2D planar gradients for 3 to 7 element systems. It can optimize the number of samples or the experiments price.

*Requirements:* Matlab2019 or more

### **PyTerK**

These Python modules and notebooks perform iterative k-fold crossvalidation for Scikit-learn or Keras models, on any datasets, with parallelization of works. Settings are written in yml files, trainings are performed in “run” notebooks and results are visualized in “report” notebooks.

*Requirements:*

- **Install following libraries, via pip or conda**
  - tensorflow - keras
  - pandas
  - scikit-learn
- **Create environment variables:**
  - path to folder that contains the datasets DATASETS\_DIR : \$ path/to/datasets/
  - path to folder that will contain the training results : RUN\_DIR : \$ path/to/run/

### **Multiple Regression**

These Python modules and notebooks perform iterative k-fold crossvalidation for statsmodel multilinear regression, on any datasets. Training and report are performed in one notebook.

*Requirements:*

- **Install following libraries, via pip or conda**
  - statsmodels.formula.api
  - pandas
  - scikit-learn
- **Create environment variables in bash\_profile:**
  - path to folder that contains the datasets: export DATASETS\_DIR = \$path/to/datasets/
  - path to folder that will contain the training results : export RUN\_DIR= \$path/to/run/



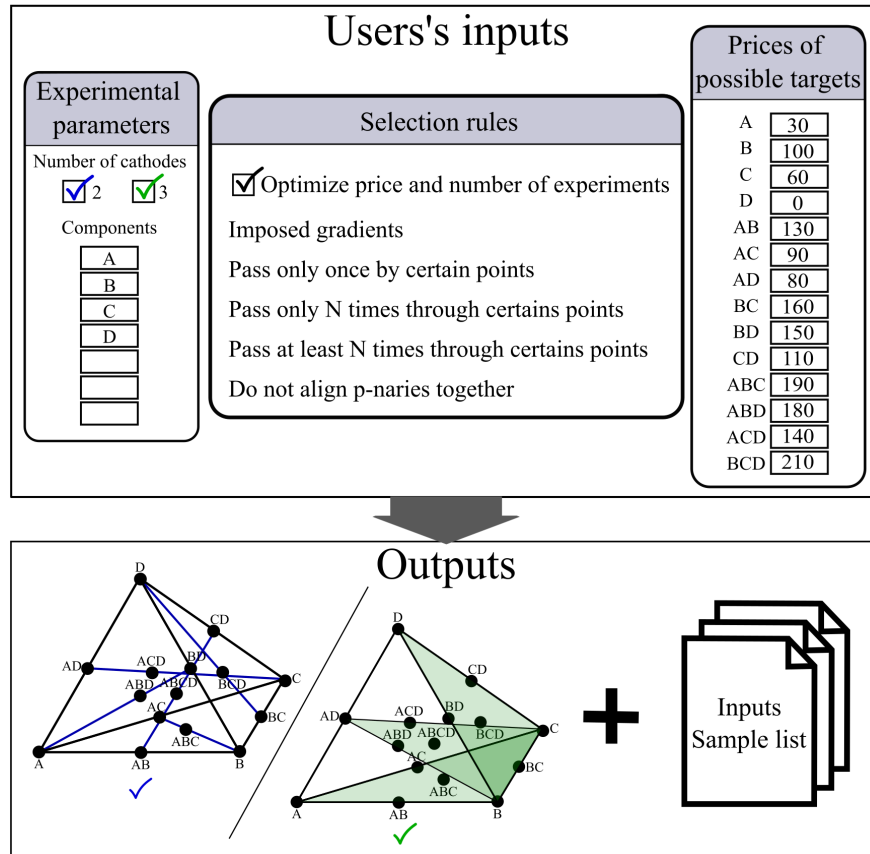
## EXPERIMENTS PLANNIFICATION

### 1.1 Principle

A Matlab GUI interface was developed in order to automatically generate a set of experiments to screen a N-element composition space using a combinatorial approach and a 2- to 3-cathodes magnetron sputtering. The starting point of this method is based on simplex centroid mixture design, in order to screen the space as uniformly as possible. From the composition points given by the mixture design, all linear/planar gradients passing by 3/7 of them are computed. Then a set of gradients/planes is chosen in order to pass at least once by each point and to respect the user inputs. The choice is done with a random exploration of all possible gradients, that starts all over again if the set does not meet the requirements.

**Features :**

- **Adaptability to the user's needs**
  - The user enters the elements of the composition space they want to explore (from 3 to 7 elements).
  - Chooses if they are using two or three cathodes
  - Indicates if they want to preferentially explore some point of the mixture design.
- Representation of the composition space and of the gradients/planes that are explored.
- Give the list of targets that allow to perform the experiments



## 1.2 Modules

`modules.check_do_not_align(name_alignments, index, do_not_align)`

Check user condition to not align certain mixtures in the same gradient

**Parameters**

- **name\_alignments** (array(str)) – name of points through which the gradients passes
- **index** (int) – index of alignments in the list of all alignments
- **do\_not\_align** (cell(list(str))) – list of mixtures that must not be aligned

**Returns** indicator: “ok” if alignment respects the user condition; else return “not ok”.

**Return type** str

`modules.check_not_repeat(name_alignments, index, name_alignement_opt, not_repeat)`

Check user condition to not repeat certain mixtures in gradients set

**Parameters**

- **name\_alignments** (array(str)) – name of points through which the gradients passes
- **index** (int) – index of alignments in the list of all alignments
- **name\_alignement\_opt** (array(str)) – gradients set that are already selected.
- **not\_repeat** (array(str)) – mixtures to no repeat in the gradients set

**Returns** indicator: “ok” if alignment respects the user condition; else return “not ok”.



**Return type** str

`modules.check_repeat_only(name_alignments, index, name_alignement_opt, repeat_only)`

Check user condition to not repeat certain mixtures in gradients set

**Parameters**

- **name\_alignments** (array(str)) – name of points through which the gradients passes
- **index** (int) – index of alignments in the list of all alignments
- **name\_alignement\_opt** (array(str)) – gradients set that are allready selected.
- **repeat\_only** (array(str, int)) – name of mixtures that must be repeated a limited number of time and this limited number of time

**Returns** indicator: “ok” if alignement respects the user condition; else return “not ok”.

**Return type** str

`modules.compute_alignments(mixture, name_mixture, nb_type_mixture)`

For a reference mixture, the function calculates the vector coefficient between this reference mixture and all the other mixtures with same or higher order. Then it looks for equals vector coefficients for segments with a common point to determine which two other mixture points are aligned with the reference mixture

**Parameters**

- **mixture** (cell{array}) – coordinates of mixtures, cell index being the mixture order  
*eg: mixtures{2} contains the binaries coordinates*
- **name\_mixture** (cell{str}) – name of mixture, cell index being the mixture order
- **nb\_type\_mixture** (int) – number of type/order of mixtures to explore

**Returns**

- **alignments** : coordinates of the mixtures through which the gradient pass (3x3 columns)
- **name\_alignments**: mixture names through which the gradient pass

**Return type** array(float), array(str)

`modules.compute_planes(name_alignment, alignments, nb_type_mixture)`

From the gradients, planes are defines in the composition space made by 3 gradients with common points, encompassing 7 points of the mixtrure design. This means that the plane is centered on one of the point of the mixture design

*eg: Nb-NbTi-Ti, Ti-TiZr-Zr and Nb-NbZr-Zr are forming a plane in a compositional space center on the ternary NbTiZr wich is a point of the mixture design: the plane is valid*

**Parameters**

- **name\_alignment** (array(str)) – points through which the gradient go
- **alignments** (array(float)) – coordinates of the points through which the gradient go (3x3 columns)
- **nb\_type\_mixture** (int) – number of type/order of mixtures to explore

**Returns** plane\_points: mixture names encompassed by the planes

**Returns** plane\_coord: coordinates of the mixtures encompassed by the planes (7x3 columns)

`modules.coordinates_name_centroid_points(nb_elements, name_elements)`

From the number and the name of the system elements, the function calculates the coordinates of the pure elements (standard uniform distribution in space) and of the equiolar mixtures of the Simplex Centroid mixture Design (all binaries, ternaries...).

**Parameters** `nb_elements` (int) – number of components

**Name\_element list(str)** namer of components

**Returns** mixture: cell coordinates of all equimolar mixture

**Return type** cell

**Returns** name\_mixture: containing the names of the equimolar mixtures.

**Return type** cell

`modules.count_occur(element, list)`

Count the number of occurrence of an element in a list

**Parameters**

- **element** – counted number or string
- **list** (list) – list in which the element is counted

**Returns** count: number of repetition of the element in list

**Return type** int

`modules.fix_nb_repetition(repeat_list, fig, position)`

**This function is a callbacks of push buttons associated to listboxes** When the buttons are pushed, the function identifies which mixture should be repeated Then it display in the interface the names of the mixtures that should be repeated and an edit box in which the user can enter the number of repetitions.

**Parameters**

- **repeat\_list** (UIcontrol) – contains the mixture that should be repeated
- **fig** (figure) – working interface / window
- **position** (list(float)) – position features of repeated list

**Returns** nb\_repet: edit boxes in which the user will enter the number of repetition of each mixture

`modules.get_elements(elements, fig1)`

Acquire the components name entered by the user

**Parameters**

- **elements** (UIcontrol) – edit boxes in which the user has entered the elements names
- **fig1** (figure) – interface window

**Returns** name\_elements: name of elements

**Return type** list(str)

`modules.gradients_set(name_mixture, mixture, alignments, name_alignment)`

Selection of a gradients set that pass at least once through each point of the mixture design and that respect user condition inputs

**Parameters**

- **name\_mixture** (cell(str)) – name of mixture, cell index being the mixture order

- **mixture** (cell(float)) – coordinates of mixtures, cell index being the mixture order  
*eg: mixtures{2} contains the binaries coordinates*
- **alignments** (cell(float)) – coordinates of the points through which the gradient pass  
(3x3 columns)
- **name\_alignement** (cell(str)) – points through which the gradient pass

**Returns**

- name\_alignement\_opt: name of mixture through which the set of gradients pass
- alignment\_opt: coordinates of mixture through which the set of gradients pass

**Return type** array(str),array(float)

`modules.index_alignments(cell_coeff_dir)`

**Called in *compute\_alignments*** [we get cell structure with vector coefficient between one reference mixture and all the mixtures with the same or higher orders .] This function compares all the coefficients one by one to find equal ones

**Parameters** `cell_coeff_dir` (cell) – contains director coefficient of vectors between one reference mixtures and all the others with same or higher order.

**Returns** cell indice\_cell, indice\_list: indices of the cell and list where two coefficients are equals.  
Allow to identify pair of equal coefficient to identify aligned points.

`modules.kill_program()`

Kill the program is the user pushed STOP button

**Returns** display the message “kill” to indicate state

`modules.lineIntersect3D(PA, PB)`

Find intersection point of lines in 3D space, in the least squares sense.

**Parameters**

- **PA** – Nx3-matrix containing starting point of N lines
- **PB** – Nx3-matrix containing end point of N lines

**Returns**

- P\_Intersect: Best intersection point of the N lines, in least squares sense.
- distances: Distances from intersection point to the input lines

Anders Eikenes (2022). Intersection point of lines in 3D space (<https://www.mathworks.com/matlabcentral/fileexchange/37192-intersection-point-of-lines-in-3d-space>), MATLAB Central File Exchange. Retrieved February 10, 2022.

`modules.listing_targets(name_alignement_opt)`

Lists the targets to use from the selected optimized set of gradients

**Parameters** `name_alignement_opt` (array(str)) – name of mixtures through which pass the gradients

**Returns** list(str) list\_target: list the target compositions to use to deposit these gradients

`modules.listing_targets_3cath(name_planes_opt)`

**Lists the targets to use from the selected optimized set of planar** gradients

**Parameters** `name_planes_opt` (array(str)) – name of mixtures encompassed by planar gradients

**Returns** list(str) list\_target: list the target compositions to use to deposit these gradients

`modules.plot_compo_space_gradients(nb_elements, mixture, name_mixture, name_elements, gradients, gradients_color)`

Plot the composition space with all the simplexe centroid points and linear gradients

**Parameters**

- **nb\_elements** (int) – number of components
- **mixture** (cell(float)) – mixture points coordinates
- **name\_mixture** (cell(str)) – mixture names
- **name\_elements** (list(str)) – name of the components
- **gradients** (array(str)) – coordinates of the gradients points
- **gradients\_color** – color of the gradients for plot

**Returns** fig: plot the compositions space dans gradients

`modules.plot_compo_space_planes(nb_elements, mixture, name_mixture, name_elements, plane_coord, plane_color)`

Plot the composition space with all the simplexe centroid points and planar gradients

**Parameters**

- **nb\_elements** (int) – number of components
- **mixture** (cell(float)) – mixture points coordinates
- **name\_mixture** (cell(str)) – mixture names
- **name\_elements** (list(str)) – name of the components
- **plane\_coord** (array(str)) – coordinates of the planes points
- **plane\_color** – color of the plane for plot

**Returns** fig: plot the compositions space dans gradients

`modules.parameters_file()`

Write the users inputs and chosen parameters for one run of the interface in text file.

`modules.planes_set(name_mixture, mixture, planes, name_planes)`

Selection of a planes set that encompass at least once each point of the mixture design and that respect user condition inputs

**Parameters**

- **name\_mixture** (cell(str)) – name of mixture, cell index being the mixture order
- **mixture** (cell(float)) – coordinates of mixtures, cell index being the mixture order  
*eg: mixtures{2} contains the binaries coordinates*
- **plane** (cell(float)) – coordinates of the points through which the planes pass (3x3 columns)
- **name\_planes** (cell(str)) – points through which the planes pass

**Returns** array(str) name\_planes\_opt: name of mixture through which the set of planes pass

**Returns** array(float) planes\_opt: coordinates of mixture through which the set of planes pass

`modules.price_calculation(prices_list, target_list)`

Calculate the price of a set of experiment

**Parameters**

- **prices\_list** (`list(str, float)`) – list of possible targets and associated price
- **targets** (`list(str)`) – list of targets associated to one set of linear gradients or planar gradients

**Returns** price: total price of the targets required for a set of linear gradients or planar gradients

**Return type** float

`modules.vector_coeff(A, B)`

Compute normed vector coefficients between two points.

**Parameters** **A,B** (`list(float)`) – coordinates of two points

**Returns** coordinates of the normed vector corresponding to (AB) line



## PYTERK PACKAGE

### 2.1 Module contents

**PyTerK** - A Python Iterated K-fold cross validation with shuffling

By E Garel / JL Parouty - SIMaP 2021

This package allows you to perform a **statistical evaluation** of different learning strategies (Keras/sklearn) by varying different (hyper)parameters.

#### 2.1.1 Description :

It is possible to combine the following (hyper)parameters :

- datasets
- models (with their characteristics...)
- batch size
- epochs
- iterations
- k fold
- seed (to control pseudo random generator)

It is possible, for example, to combine 3 datasets, with 3 models and to perform for each combination, 5 iterations of a cross validation of KFold type, with k=10. In this case, the total number of models to test would be  $3 \times 3 \times 5 \times 10 = 450$  training sessions... So, be careful, the number of model.fit can quickly be very important !

The tasks will be run in **parallel** on the different CPUs/cores available.

#### 2.1.2 Documentation and examples :

Here is a basic example, detailed in a notebook :

```
import pyterk.config      as config
import pyterk.reporter    as reporter
import pyterk.task_manager as task_manager

settings = config.load('settings_example.yml')
```

(continues on next page)

(continued from previous page)

```
task_manager.add_combinational_iterative_manyfold(settings, run_key= 'Example-03.1')
task_manager.run()

reporter.show_run_reports(settings)
```

This will retrieve all settings from *settings\_example.yml*, prepare the different tasks and execute them. The last call, intended to be used from a Jupyter lab notebook, displays a complete execution report.

You can find **3 full example notebooks**, with a setting file :

- settings\_example.yml
- 01-Example-01.ipynb
- 02-Example-02.ipynb
- 03-Example-03.ipynb

pyterk.**VERSION** = 2.14  
pyterk version

## 2.2 Submodules

## 2.3 pyterk.config module

Configuration management.

The settings files allow to specify datasets and models.

### 2.3.1 Utilisation:

Loading a settings file:

```
settings = config.load('settings_example.yml')
```

or:

```
settings = config.load('settings_example.yml',
                      datasets_dir_env='DATASETS_DIR')
```

where DATASETS\_DIR is an environment variable that will override *datasets\_dir* directive in settings file.

pyterk.config.**datasets** = None  
datasets profiles

pyterk.config.**load**(filename, datasets\_dir\_env='DATASETS\_DIR', run\_dir\_env='RUN\_DIR', verbose=0)  
Load a setting file and dfined datasets. If given, environment variable can be use to override *datasets\_dir* directive from setting file. Usefull for portability between several sites.

#### Parameters

- **filename** (*string*) – Filename of the yaml setting file
- **datasets\_dir\_env** (*string*) – Name of the overriding environment variable
- **verbose** (*int*) – verbose mode for loaded datasets (0).



**Returns** A dict from setting file, completed by datasets and more.

`pyterk.config.models = None`  
models profiles

`pyterk.config.run_dir = None`  
run\_dir, the place to put all output directories

`pyterk.config.runs = None`  
dict of runs section

`pyterk.config.settings = None`  
Dict of settings

## 2.4 pyterk.models module

This module is for internal use only - You do not have to interact with.

`pyterk.models.get_model(profile)`

Get a model from a model profile. The profile contains the module and function name of the model, and the arguments. The model will be retrieved by calling the function with the arguments.

**Parameters** `profile` (*dict*) – a model profile

**Returns** keras or scikitlearn model as defined in the profile.

**Return type** model (keras/scikitlearn model)

## 2.5 pyterk.reporter module

Module to generate execution reports.

During the run of the tasks, the bestmodel and results are saved in h5 and json files:

- *about.json* : information and description of the task
- *history.json* : history from model.fit()
- *evaluation.json* : evaluation from model.evaluate()
- *bestmodel.h5* : best model

### 2.5.1 Example :

```
reporter.show_run_reports(settings,
                           args      = ['dataset_id', 'model_id', 'batch_size'],
                           evaluation = [2])
```

This module will retrieve information from json files and generate a report.

`pyterk.reporter.plot_confusion(run_dir, predict_type='softmax', normalize='pred', figsize=(5, 5), savefig=True, mplstyle='pyterk')`

Plot a confusion matrix

**Parameters**

- **iterations\_dir** – a directory with iterations subdirs (iter-000, iter-001, ...)

- **predict\_type** – sigmoid, softmax or classes
- **normalize** – true, pred, all or None (pred)
- **figsize** – figure size
- **savefig** – save fig (True) or not (False)

**Returns** Just plot the matrix and print report and hamming loss

```
pyterk.reporter.plot_distribution(run_dir, metric_id=0, bins=10, min=None, max=None, figsize=(10, 8),  
                                savefig=False, mplstyle='pyterk')
```

Plot distribution of a given metric from an evaluation.json saved file. For a kfold or an iterative kfold, all evaluation data are concatenated in an evaluation.json file in main run\_dir.

#### Parameters

- **run\_dir** (*string*) – directory path of json evaluation file
- **metricid** (*int*) – number of metric to plot. Example : 2
- **min** (*int*) – min value
- **max** (*int*) – max value
- **bins** (*int*) – number of bins
- **figsize** (*tuple*) – figure size, default is (10,8)
- **savefig** (*boolean*) – if True, figure will be save in run\_dir.
- **mplstyle** (*string*) – name of matplotlib style. default is 'pyterk', but all matplotlib are ok (default, bmh, ...)

**Returns** Nothing, but display a beautifull distribution plot !

```
pyterk.reporter.plot_history(run_dir, metric='val_mae', min=None, max=None, figsize=(10, 8),  
                             savefig=False, mplstyle='pyterk')
```

Plot history evolution from history.json saved file. For a kfold or an iterative kfold, all history data are concatenated in history.json file in main run\_dir. This will plot a curve for each one in a common plot.

#### Parameters

- **run\_dir** (*string*) – directory path of json history file
- **metric** (*string*) – metric name to plot. Example : 'val\_mae'
- **figsize** (*tuple*) – figure size, default is (10,8)
- **savefig** (*boolean*) – if True, figure will be save in run\_dir.
- **mplstyle** (*string*) – name of matplotlib style. default is 'pyterk', but all matplotlib are ok (default, bmh, ...)

**Returns** Nothing, but display a beautifull plot !

```
pyterk.reporter.plot_kfold_correlation(run_dir, channel=0, figsize=(8, 6), axes_min='auto',  
                                       axes_max='auto', yy_deltamax=None, marker='o', markersize=8,  
                                       alpha=0.7, color='auto', savefig=True, mplstyle='pyterk')
```

Plot a correlation for a (y\_test, y\_pred) saved json file.

#### Parameters

- **run\_file** – a manyfold directory where kfold subdirectories are
- **channel** – composant of y to plot
- **figsize** (*tuple*) – figure size, default is (10,8)

- **axes\_min** – min value for x and y axe. ‘auto’ or float
- **axes\_max** – max value for x and y axe. ‘auto’ or float
- **mplstyle** (*string*) – name of matplotlib style. default is ‘pyterk’, but all matplotlib are ok (default, bmh, ...)
- **marker** – marker, default is ‘.’
- **markersize** – marker size
- **alpha** – marker alpha
- **color** – plot color or ‘auto’
- **savefig** – if True, save fig in run\_dir

**Returns** Nothing, but display a beautiful correlation plot

```
pyterk.reporter.show_report(run_dir, padding="", sections=['title', 'context', 'args', 'settings', 'evaluation',
    'monitoring', 'history', 'distribution', 'correlation'], context=['function', 'version',
    'date', 'description', 'seed'], args=['run_dir', 'dataset_id', 'model_id', 'n_iter',
    'k_fold', 'epochs', 'batch_size'], settings=['file', 'version', 'description',
    'datasets_dir', 'run_dir'], evaluation=['all'], monitoring=['duration',
    'used_data'], history=[{'metric': 'val_mae', 'min': None, 'max': None, 'figsize':
    (8, 6), 'savefig': True, 'mplstyle': 'pyterk'}], distribution=[{'metric_id': 2, 'bins':
    4, 'min': None, 'max': None, 'figsize': (8, 6), 'savefig': True, 'mplstyle': 'pyterk'}],
    correlation=[{'axes_min': 'auto', 'axes_max': 'auto', 'figsize': (8, 6), 'marker':
    '.', 'markersize': 8, 'alpha': 0.7, 'color': 'auto', 'savefig': True, 'mplstyle':
    'pyterk'}], confusion=[{'normalize': 'pred', 'predict_type': 'softmax', 'figsize': (5,
    5), 'savefig': True, 'mplstyle': 'pyterk'}])
```

Builds and displays a report from the json data of a given run\_dir.

#### Parameters

- **run\_dir** (*string*) – directory path of json report file
- **sections** (*list*) – list of sections to include in the report
- **context** (*list*) – informations to include in context section
- **args** (*list*) – informations to include in args section
- **settings** (*list*) – informations to include in settings section
- **evaluation** (*list*) – #metrics to include in evaluation section. ‘all’ mean all. Example : [0,1,2]
- **history** (*dict*) – parameters for history plot - see *plot\_history*
- **distribution** (*dict*) – parameters for metrics distribution plot
- **correlation** (*dict*) – parameters for correlation plot
- **confusion** (*dict*) – parameters for confusion matrix (need yytest files)

```
pyterk.reporter.show_run_reports(run_config, run_filter='*', sections=['title', 'context', 'args', 'settings',
                              'evaluation', 'monitoring', 'history', 'distribution', 'correlation',
                              'confusion'], context=['function', 'version', 'date', 'description', 'seed'],
                              args=['run_dir', 'dataset_id', 'model_id', 'n_iter', 'k_fold', 'epochs',
                              'batch_size'], settings=['file', 'version', 'description', 'datasets_dir',
                              'run_dir'], evaluation=['all'], monitoring=['duration', 'used_data'],
                              history=[{'metric': 'val_mae', 'min': None, 'max': None, 'figsize': (8, 6),
                              'savefig': True, 'mplstyle': 'pyterk'}], distribution=[{'metric_id': 2, 'bins':
                              4, 'min': None, 'max': None, 'figsize': (8, 6), 'savefig': True, 'mplstyle':
                              'pyterk'}], correlation=[{'axes_min': 'auto', 'axes_max': 'auto', 'figsize':
                              (8, 6), 'marker': '.', 'markersize': 8, 'alpha': 0.7, 'color': 'auto', 'savefig':
                              True, 'mplstyle': 'pyterk'}], confusion=[{'normalize': 'pred',
                              'predict_type': 'softmax', 'figsize': (5, 5), 'savefig': True, 'mplstyle':
                              'pyterk'}])
```

Displays a full report in two parts, short and long, for all runs defined in the settings. Very simple to use...

#### Parameters

- **run\_config** (*dict*) – settings, issued from config.load()
- **run\_filter** (*regex*) – regex to filter run entries from yaml settings file (\*.\*)
- **sections** (*list*) – list of sections to include in the report
- **context** (*list*) – informations to include in context section
- **args** (*list*) – informations to include in args section
- **settings** (*list*) – informations to include in settings section
- **evaluation** (*list*) – #metrics to include in evaluation section. 'all' mean all. Example : [0,1,2]
- **history** (*dict*) – parameters for history plot - see *plot\_history*
- **distribution** (*dict*) – parameters for metrics distribution plot
- **correlation** (*dict*) – parameters for correlation plot
- **confusion** (*dict*) – parameters for confusion matrix (need yytest files)

**Returns** Nothing, but display a short and long report, with index.

## 2.6 pyterk.task\_manager module

Allows to generate tasks and to execute them in a distributed way.

See example notebook : *03-Example-03.ipynb*

Example :

```
task_manager.add_combinational_iterative_manyfold(settings = settings,
                                                  run_key = 'Example-03.3')
```

```
pyterk.task_manager.add_combinational_iterative_manyfold(settings=None, run_key=None,
                                                         verbose=1)
```

Add tasks for a combinational iterative manyfold - see *03-Example-03.ipynb* Generates all the tasks of the combinatorial described in the run section of the settings file.

#### Parameters

- **settings** (*dict*) – settings
- **run\_key** (*string*) – name of the config run section
- **verbose** (*int*) – verbosity of generated tasks

**Returns** Nothings. Task are added to the pending tasks queue.

```
pyterk.task_manager.add_iterative_manyfold(settings=None, run_dir=None, dataset_id=None,
                                           model_id=None, n_iter=2, k_fold=10, epochs=10,
                                           batch_size=10, description=None, save_xctest=False,
                                           save_ytest=False, verbose=1)
```

Add tasks for an iterative manyfold - see *02-Example-02.ipynb* Generate  $n\_iter * k\_fold$  tasks, each iteration will generate a subdirectory in `run_dir`.

#### Parameters

- **settings** (*dict*) – settings
- **run\_dir** (*string*) – run directory to output k results (json files and best model)
- **dataset\_id** (*string*) – datasets id in settings file
- **model\_id** (*string*) – model id in settings file
- **n\_iter** (*int*) – number of iteration
- **k\_fold** (*int*) – number of fold
- **epochs** (*int*) – number of epochs
- **batch\_size** (*int*) – size of batch
- **description** (*string*) – description of the action
- **save\_xctest** (*Boolean*) – save `x_test` as json file, or not
- **save\_ytest** (*Boolean*) – save `y_test` and `y_pred` as json file, or not
- **verbose** (*int*) – verbosity of generated tasks

**Returns** Nothings. Task are added to the pending tasks queue.

```
pyterk.task_manager.add_manyfold(settings=None, run_dir=None, dataset_id=None, model_id=None,
                                 k_fold=10, epochs=10, batch_size=10, description=None,
                                 save_xctest=False, save_ytest=False, verbose=1)
```

Add tasks for a manyfold - see *01-Example-01.ipynb* Generate  $k\_fold$  tasks, each task will generate one subdirectory in `run_dir`.

#### Parameters

- **settings** (*dict*) – settings
- **run\_dir** (*string*) – run directory to output k results (json files and best model)
- **dataset\_id** (*string*) – datasets id in settings file
- **model\_id** (*string*) – model id in settings file
- **k\_fold** (*int*) – number of fold
- **epochs** (*int*) – number of epochs
- **batch\_size** (*int*) – size of batch
- **description** (*string*) – description of the action
- **save\_xctest** (*Boolean*) – save `x_test` as json file, or not

- **save\_yytest** (*Boolean*) – save `y_test` and `y_pred` as json file, or not
- **verbose** (*int*) – verbosity of generated tasks

**Returns** Nothings. Task are added to the pending tasks queue.

`pyterk.task_manager.reset()`

Reset pending tasks. Suppress all of them !

`pyterk.task_manager.run(processes=None, maxtasksperchild=10, verbose=1)`

`pyterk.task_manager.seed(seed=None)`

Init random generators with given seed

`pyterk.task_manager.show_tasks_size()`

Print pending tasks size

## 2.7 pyterk.worker module

This module is for internal use only - You do not have to interact with ;-).

`pyterk.worker.get_model_family(model)`

Should return the model family : 'tensorflow', 'keras' or 'sklearn'

`pyterk.worker.init(s, l, v)`

`pyterk.worker.model_fit(run_dir=None, dataset_id=None, train_index=None, test_index=None,  
model_id=None, epochs=None, batch_size=None, seed=None, description=None,  
save_xxtest=False, save_yytest=False)`

`pyterk.worker.model_fit_sklearn(model, run_dir=None, x_train=None, y_train=None, x_test=None,  
y_test=None, save_xxtest=False, save_yytest=False)`

`pyterk.worker.model_fit_tensorflow(model, run_dir=None, x_train=None, y_train=None, x_test=None,  
y_test=None, epochs=None, batch_size=None, save_xxtest=False,  
save_yytest=False)`

## MULTIPLEREGRESSION MODULE

### 3.1 Description

Module to train Multiple Linear Regression with Scheffe interaction terms, with iterative k-fold crossvalidation. The method that lies behind the model training is identical to PyTerK one, with iterative k-fold cross-validation.

**Contains functions to :**

- generate interactions
- train regression models
- plot iterative k-fold crossvalidation results

### 3.2 Functions

`MultipleRegression.Scheffe_interactions_terms(data, in_percent=True, compo_columns=['Zr', 'Nb', 'Mo', 'Ti', 'Cr'])`

Shaping composition in percentage rate into percentage Compute interaction terms for all Scheffe interactions for quartic multiple regression and add it to dataframe data

**Parameters** `panda.DataFrame` – dataset that contains compositions in Zr, Nb, Mo, Ti, Cr in columns of the same name

**Returns** extended input dataset with interactions

**Return type** `DataFrame`

`MultipleRegression.fit_outputs(model_expression, k, nb_it, output, X, y)`

Takes an OLS model expression, and use it to perform regression between X and y . Model regression is performed using iterative k-fold crossvalidation Evaluation is performed through R2 and MAE computation

**Parameters**

- **OLS-formula** – contain OLS formula for regression
- **k** (`int`) – number of folds for iterative k-fold crossvalidation
- **nb\_it** (`int`) – number of iterations for iterative k-fold crossvalidation
- **output** (`str`) – name of the Y output to fit
- **X** (`panda.DataFrame`) – contains composition and interaction terms for regression input
- **y** (`panda.DataFrame`) – contains single column dataframe with regression output

**Return model** model coefficients and p-values

**Return type** statsmodels.regression.linear\_model.RegressionResultsWrapper

**Return MAE\_list** list of MAE for every run of iterative k-fold crossvalidation, between expected vs predicted value on test set

**Return type** list

**Return R2\_adj\_list** list of R2 adjusted for every run of iterative k-fold crossvalidation, between expected vs predicted value on test set

**Return type** list

**Returns** Y\_pred : list of predicted values on test set

**Return type** list

**Returns** Y\_test : list of expected values on test set

**Return type** list

MultipleRegression.**plot\_result**(*metric, output, val\_metric, Y\_pred, Y\_test, min\_hist, max\_hist, iter, kfold, save\_distri, save\_regression*)

Plot metric histogram and regression between predictions and test values and save graphs

**Parameters**

- **metric** (*str*) – name of the metric distribution to plot
- **output** (*str*) – name of the Y output to fit
- **val\_metric** (*list*) – list of MAE for every run of iterative k-fold crossvalidation, between expected vs predicted value on test set
- **Y\_pred** (*list*) – list of predicted values on test set
- **Y\_test** (*list*) – list of expected values on test set
- **min\_hist** (*int*) – minimum of abscissa for metric distribution histogram
- **max\_hist** (*int*) – maximum of abscissa for metric distribution histogram
- **iter** (*int*) – plot regression over a certain number of iterations
- **save\_distri** (*str*) – path to save metric distribution
- **save\_regression** (*std*) – path to save regression

**Parm int kfold** plot regression over a certain number of k-fold for each iteration



## INDICES AND TABLES

- `genindex`
- `modindex`
- `search`



## PYTHON MODULE INDEX

### m

`MultipleRegression`, 19

### p

`pyterk`, 11

`pyterk.config`, 12

`pyterk.models`, 13

`pyterk.reporter`, 13

`pyterk.task_manager`, 16

`pyterk.worker`, 18



## MATLAB MODULE INDEX

m

modules, [9](#)



## A

add\_combinational\_iterative\_manyfold() (in module *pyterk.task\_manager*), 16  
 add\_iterative\_manyfold() (in module *pyterk.task\_manager*), 17  
 add\_manyfold() (in module *pyterk.task\_manager*), 17

## C

check\_do\_not\_align() (in module *modules*), 4  
 check\_not\_repeat() (in module *modules*), 4  
 check\_repeat\_only() (in module *modules*), 5  
 compute\_alignments() (in module *modules*), 5  
 compute\_planes() (in module *modules*), 5  
 coordinates\_name\_centroid\_points() (in module *modules*), 5  
 count\_occur() (in module *modules*), 6

## D

datasets (in module *pyterk.config*), 12

## F

fit\_outputs() (in module *MultipleRegression*), 19  
 fix\_nb\_repetition() (in module *modules*), 6

## G

get\_elements() (in module *modules*), 6  
 get\_model() (in module *pyterk.models*), 13  
 get\_model\_family() (in module *pyterk.worker*), 18  
 gradients\_set() (in module *modules*), 6

## I

index\_alignments() (in module *modules*), 7  
 init() (in module *pyterk.worker*), 18

## K

kill\_program() (in module *modules*), 7

## L

lineIntersect3D() (in module *modules*), 7  
 listing\_targets() (in module *modules*), 7  
 listing\_targets\_3cath() (in module *modules*), 7

load() (in module *pyterk.config*), 12

## M

model\_fit() (in module *pyterk.worker*), 18  
 model\_fit\_sklearn() (in module *pyterk.worker*), 18  
 model\_fit\_tensorflow() (in module *pyterk.worker*), 18  
 models (in module *pyterk.config*), 13  
 module  
   MultipleRegression, 19  
   pyterk, 11  
   pyterk.config, 12  
   pyterk.models, 13  
   pyterk.reporter, 13  
   pyterk.task\_manager, 16  
   pyterk.worker, 18  
 modules (module), 4–9  
 MultipleRegression  
   module, 19

## P

parameters\_file() (in module *modules*), 8  
 planes\_set() (in module *modules*), 8  
 plot\_compo\_space\_gradients() (in module *modules*), 8  
 plot\_compo\_space\_planes() (in module *modules*), 8  
 plot\_confusion() (in module *pyterk.reporter*), 13  
 plot\_distribution() (in module *pyterk.reporter*), 14  
 plot\_history() (in module *pyterk.reporter*), 14  
 plot\_kfold\_correlation() (in module *pyterk.reporter*), 14  
 plot\_result() (in module *MultipleRegression*), 20  
 price\_calculation() (in module *modules*), 8  
 pyterk  
   module, 11  
   pyterk.config  
     module, 12  
   pyterk.models  
     module, 13  
   pyterk.reporter  
     module, 13  
   pyterk.task\_manager

module, 16  
pyterk.worker  
    module, 18

## R

reset() (in module pyterk.task\_manager), 18  
run() (in module pyterk.task\_manager), 18  
run\_dir (in module pyterk.config), 13  
runs (in module pyterk.config), 13

## S

Scheffe\_interactions\_terms() (in module MultipleRegression), 19  
seed() (in module pyterk.task\_manager), 18  
settings (in module pyterk.config), 13  
show\_report() (in module pyterk.reporter), 15  
show\_run\_reports() (in module pyterk.reporter), 15  
show\_tasks\_size() (in module pyterk.task\_manager),  
    18

## V

vector\_coeff() (in module modules), 9  
VERSION (in module pyterk), 12