

# Matlab functions for replicate regression

## 1 Installation

### Getting started

1. Unpack the files from the github repository.
2. Include the path to the matlab directory 'replicate\_regression' into your MATLAB path.
3. Run 'demo\_replicate\_regression.m' to see a single replicate regression.
4. Run 'demo\_omics\_data.m' to see the analysis of a small example omics data set (data and options files are provided in the same directory).

### Requirements

The functions were developed and tested with matlab6.

### State of the software

The functions are under development and provided 'as is'. If you would like to contribute extensions to the toolbox, please let me know.

### Documentation

Documentation (in directory 'doc') has been built automatically with M2HTML

### License

The toolbox is free software; you can redistribute it and/or modify it under the terms of the GNU General Public License as published by the Free Software Foundation; either version 2, or (at your option) any later version. The toolbox is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the file COPYING for more details.

### Any questions?

Please send questions, comments, and bug reports to wolfram.liebermeister@gmail.com

## 2 Running a single replicate regression

The MATLAB function `replicate_regression.m` allows you to run a single replicate regression. All data are directly given in the form of MATLAB variables (vectors and structs).

### MATLAB function `replicate_regression.m`

The function is called as follows:

**[result, options] = replicate\_regression(t, y, sigma, r, flag\_fix\_parameters, varargin)**

Bayesian replicate regression for multiple time series measured in replicate.

Data must be provided as vectors and are transformed to logarithmic scale if desired.

#### Function arguments

**t, y, sigma, r** input data (times, values, standard errors, replicate labels)  
given as row vectors (see `replicate_regression_core.m`)

**flag\_fix\_parameters** (Boolean, optional) If set to 1, the options given in the following argument(s) will be accepted  
without changes (otherwise they will be checked and updated)

**varargin** (optional) Either a list of property/value pairs for algorithm options (list see appendix).  
or a structure containing the property/value pairs (this is mandatory if `flag_fix_parameters` is set to 1)

#### Function output

**result** matlab struct with results from replicate regression

**options** matlab struct with options values that were used in the calculation

The options list is supposed to be ordered by priority; earlier options override later options. The function is a wrapper for the function '`replicate_regression_core`'. In converting the data to logarithms, `y` and `sigma` are either taken to be medians and geometric standard deviations, or means and standard deviations of the data values. The choice is defined by the argument '`options.transformation`'

The basis functions are adjusted to the final time interval  $[t_a, t_b]$  (from `tt`)

'cos'	cosine function, zero slope at $t=t_a$ and $t=t_b$
'sin'	sine function, zero value at $t=t_a$ and $t=t_b$
'sin_half'	sine function, zero value at $t=t_a$
'sin_horizontal'	sine function, zero value at $t=t_a$ , zero slope at $t=t_b$
'cos+sin'	cosine and sine functions, no restriction
'polynomial'	polynomial function, zero value at $t=t_a$
'exp'	exponentially relaxing functions ( $t < 0 \Rightarrow f=0$ ; $t > 0 \Rightarrow f = 1 - \exp(t/\tau)$ )

The entire curves are shifted by a constant basis function This can be suppressed by setting `options.use_offset = 0`

### 3 How to run a replicate regression for omics data

The MATLAB function `replicate_regression_omics.m` allows you to run replicate regression for an entire omics data set, including iterative updating of the curve priors. Data are given in the form of table files (both for omics data and for function options).

#### Matlab function `replicate_regression_omics.m`

The function is called as follows:

**`replicate_regression_omics(data_file, user_options_file, base_directory)`**

Bayesian replicate regression for omics data

#### Function arguments

**`data_file`** Omics data file (full directory path)

**`user_options_file`** Table file containing the options (full directory path). For list of options, see appendix.

**`base_directory`** Directory name for results (full directory path)

**Function output** Output data and graphics are written to files

#### How to prepare data and run `replicate_regression_omics.m`

1. Create a directory for the analysis
2. Create in this directory subdirectories “data”, “options”, “results”, and “graphics”
3. Create a data file (tab-separated text table in the format described below) and save it to the “data” subdirectory
4. Create an options file (tab-separated text table in the format described below) and save it to the “options” subdirectory
5. Start matlab and run replicate regression  
(see matlab script `replicate_regression/demo/omics_data_example/omics_data_example.m`)

```
% Directory name for the omics set
base_DIR = [replicate_regression_DIR '/demo/omics_data_example/'];

% Name of options file
foptions_file = [base_DIR '/options/options_omics_data_example.csv'];

% Run script for replicate regression of omics data
replicate_regression_omics_analysis;
```

Examples of data and options files can be found in the subdirectory `replicate_regression/demo/omics_data_example`

### Format of data file (tab-separated text file)

The following example table shows the format for data tables:

!UniprotID	S1R1	S2R1	S3R1	S4R1	S1R2	S2R2	S3R2	S4R2	S1R3	S2R3	S3R3	S4R3
!Time	0	10	20	30	0	10	20	30	0	10	20	30
!Replicate	R1	R1	R1	R1	R2	R2	R2	R2	R3	R3	R3	R3
Q04747	1.05	1.19	0.64	1.09	1.39	0.69	1.33	1.19	1.3	1.07	0.93	0.99
P27206	0.98	1.16	0.65	1.29	1.27	0.48	1.12	1.12	1.3	1.14	1.04	1.02
...	...	...	...	...	...	...	...	...	...	...	...	...

The first line contains the headers of the protein names column (e.g., !BSUnumber, !BGnumber, !Gene-Name, !UniprotID), followed by sample names (as headers of data columns). The second line contains the time points as numbers. The third line contains replicate names, and the fourth line (which is optional) contains the type of value given ('Value', 'Mean', or 'Std'). All further lines contain numerical data. The entry “!UniprotID” can be replaced by the name of gene identifiers used. The entries “!Time” and “!Replicate” are fixed. The sample names (in first line) and replicate names (third line) can be freely chosen; however, they must not start with a number and may not contain any special characters (e.g., “:” or “.”) or whitespace characters.

### Format of options file (tab-separated text file)

Options can be given in a tab-separated table file like in the following example:

% THIS IS A COMMENT	
data_dir	[base_DIR '/data/']
result_dir	[base_DIR '/results/']
graphics_dir	[base_DIR '/graphics/']
data_file_csv	my_data_set_data.csv
data_file_matlab	my_data_set_data.mat
options_file	[replicate_regression_DIR '/resources/mcr_options_nonlogarithmic_data.csv']
options_out_csv	my_data_set_result_regression_options.csv
result_file_matlab	my_data_set_result.mat
result_file_csv	my_data_set_result.csv
graphics_file	my_data_set_graphics
log_file	my_data_set_log.txt
data_time_unit	min
data_scale	absolute
normalise_by_median	1
convert_to_logarithm	1
log_transformation	arithmetic
data_std_relative	0.25
data_min_data_points	3
fixed_prior	0
prior_updating	3
updating_factor	1.2
t_smooth	30
run_crossvalidation	1

Each line contains one attribute. The first entry contains the attribute name, the second column the attribute value (string or number). All further entries are ignored. Lines starting with the '%' character are ignored (can be used for comments). The attribute 'options\_file' allows to declare another options file containing default options. The attribute 'data\_file\_csv' contains the name of the data file

## A Function options

### A.1 Options for matlab function `replicate_regression.m`

Options for the matlab function `replicate_regression` are given in a struct called `options`. Possible fields with default values are listed below.

OPTION	IN CORE	TYPE	DEFAULT	MEANING
<code>options.verbose</code>		Boolean	1	Output information during regression
<code>options.is_logarithmic</code>		Boolean	0	Declare that data are logarithmic
<code>options.convert_to_logarithm</code>		Boolean	1	Convert data to logarithms for regression
<code>options.log_transformation</code>		string	'arithmetic'	'arithmetic', 'geometric'
<code>options.run_crossvalidation</code>		Boolean	0	Run crossvalidation
<code>options.set_std</code>		float	nan	Value to replace all data standard deviations
<code>options.insert_std</code>		float	1	Value to replace missing data standard deviations
<code>options.start_at_t</code>		float	0	Start regression curves at starting time 'start_at_t' (instead of t=0)
<code>options.start_value</code>		float	nan	Fixed start value for regression curves
<code>options.shift_data</code>		string	'mean'	Policy for shifting data before regression 'none', 'fixed_start_value', 'mean', 'initial', 'fixed_1'
<code>options.shift_value</code>		float	nan	Shift used when shifting the data
<code>options.basis</code>	X	string	'cos+sin'	Type of basis functions (see table below)
<code>options.n_comp</code>	X	int	nan	Fixed number of basis functions
<code>options.n_comp_min</code>		int	1	Minimal number of basis functions
<code>options.n_comp_max</code>		int	20	Maximal number of basis functions
<code>options.use_offset</code>	X	Boolean	1	Use constant function as one of the basis functions
<code>options.constant_before_start</code>	X	Boolean	0	Set all basis functions constant for $t \leq 0$
<code>options.deviation_same_start</code>		Boolean	0	Enforce identical start values for all replicates
<code>options.remove_offset</code>	X	Boolean	0	Omit offset when creating the regression curves
<code>options.t_smooth</code>		float	nan	Time constant for setting decreasing prior widths
<code>options.t_jump</code>	X	float	nan	Time constant for initial jump basis function
<code>options.t_interp</code>		float	t	Time points for interpolated regression curves
<code>options.average_std</code>	X	string	'std_dev_mean'	Type of uncertainty to be reported for average curve
<code>options.central_offset_mean</code>	X	float	0	Prior mean $\sigma_{\alpha_0}$ (for $\alpha_0$ )
<code>options.central_offset_width</code>	X	float	1	Prior width $\sigma_{\alpha_0}$ (for $\alpha_0$ )
<code>options.central_first_mode_mean</code>	X	float	0	Prior mean $\sigma_{\alpha_1}$ (for $\alpha_1$ )
<code>options.central_first_mode_width</code>	X	float	1	Prior width $\sigma_{\alpha_1}$ (for $\alpha_1$ )
<code>options.central_mode_mean</code>	X	vector	[]	Prior means $\sigma_{\alpha_m}$ (for $\alpha_m$ )
<code>options.central_mode_width</code>	X	vector	[]	Prior widths $\sigma_{\alpha_m}$ (for $\alpha_m$ )
<code>options.central_jump_mean</code>	X	float	nan	Prior means $\sigma_{\alpha_{\text{jump}}}$ (for $\alpha_{\text{jump}}$ )
<code>options.central_jump_width</code>	X	float	nan	Prior widths $\sigma_{\alpha_{\text{jump}}}$ (for $\alpha_{\text{jump}}$ )
<code>options.deviation_offset_mean</code>	X	float	0	Prior mean $\sigma_{\beta_0}$ (for $\beta_0$ )
<code>options.deviation_offset_width</code>	X	float	1	Prior width $\sigma_{\beta_0}$ (for $\beta_0$ )
<code>options.deviation_first_mode_mean</code>	X	float	0	Prior mean $\sigma_{\beta_1}$ (for $\beta_1$ )
<code>options.deviation_first_mode_width</code>	X	float	1	Prior width $\sigma_{\beta_1}$ (for $\beta_1$ )
<code>options.deviation_mode_mean</code>	X	float	[]	Prior means $\sigma_{\beta_m}$ (for $\beta_m$ )
<code>options.deviation_mode_width</code>	X	float	[]	Prior widths $\sigma_{\beta_m}$ (for $\beta_m$ )
<code>options.deviation_jump_mean</code>	X	float	0	Prior means $\sigma_{\beta_{\text{jump}}}$ (for $\beta_{\text{jump}}$ )
<code>options.deviation_jump_width</code>	X	float	1	Prior widths $\sigma_{\beta_{\text{jump}}}$ (for $\beta_{\text{jump}}$ )
<code>options.flag_draw_sample</code>	X	Boolean	1	Draw sample curve parameters and curve from the posterior
<code>options.flag_time_derivative</code>	X	Boolean	0	Compute time derivative curves

The options marked in column "IN CORE" are used by the underlying function `replicate_regression_core.m`

## A.2 Options for matlab function replicate\_regression\_omics.m

Options for the matlab function `replicate_regression_omics_analysis` are given in a struct called `foptions`. Possible fields with default values (see function `replicate_regression_omics_default_options`) are listed below. The same options can also be set in an options file.

OPTION	MEANING
<code>data_dir</code>	directory name for data files
<code>result_dir</code>	directory name for result files
<code>graphics_dir</code>	directory name for graphics
<code>data_file_csv</code>	filename for data file (tsv format, see examples)
<code>data_file_matlab</code>	filename for matlab data file (written during the analysis)
<code>options_file</code>	filename for default options file (tsv format)
<code>options_out_csv</code>	filename for completed options file (tsv format, written during analysis)
<code>translation_table_file</code>	filename for ID mapping table (see example)
<code>result_file_matlab</code>	filename for
<code>result_file_csv</code>	filename for <code>hahne_salt_stress_cytosol_result.tsv</code>
<code>result_file_zip</code>	filename for <code>hahne_salt_stress_result.zip</code>
<code>graphics_file</code>	file basename for graphics
<code>data_time_unit</code>	time unit ('min')
<code>data_scale</code>	'absolute' or 'log2' (also 'ln','log','log10','log2 ratio'; these are all treated like 'log2');
<code>data_min_num_replicates</code>	minimal number of valid replicates (genes with less valid replicates are discarded; default 1)
<b>For non-logarithmic data:</b>	
<code>abs_data_adjust_std_upper</code>	upper threshold; points above are outliers (increase std dev by factor of 3)
<code>abs_data_adjust_std_lower</code>	lower threshold; points below are outliers (increase std dev by factor of 3)
<code>data_std_relative</code>	default for relative standard deviation
<code>data_std_minimal</code>	minimal standard deviation
<b>For logarithmic data</b>	i.e., ( 'log2', 'ln','log','log10','log2 ratio')
<code>data_std_log</code>	default for standard deviation (on log scale)
<code>log_data_adjust_std_threshold</code>	threshold for data values (on chosen log scale) for which std dev is modified criterion: $ [\text{data value}] - [\text{median for this gene and replicate}]  < \text{threshold}$ for the inserted std dev, see next entry
<code>log_data_adjust_std_factor</code>	new std width = factor * absolute deviation from median
<code>data_min_data_points</code>	minimal number of data points required in the analysis (default 3) at least one replicate has to reach this number, points are times $t_0$ do not count replicates with less data points are ignored
<code>convert_to_logarithm</code>	convert (nonlogarithmic) data to logarithms for replicate regression (Boolean)
<code>log_transformation</code>	type of transformation 'arithmetic': data=mean values and plotting on absolute scale 'geometric' : data = median values and plotting on log scale (but data on absolute scale)
<code>ignore_std_deviations</code>	Boolean, ignore standard deviations given in data
<code>fixed_prior</code>	keeping the prior fixed? (Boolean, default 0)
<code>prior_updating</code>	number of prior updating iterations (default 10)
<code>updating_factor</code>	updating factor, default 1.1
<code>update_prior_means</code>	change parameter means from 0 to posterior means while updating? default 0
<code>t_smooth</code>	time constant defining how prior widths depend on the frequency
<code>options_start_value</code>	fixed starting value; to be inserted into options as <code>options.start_value</code>
<code>options_start_at_t</code>	Starting time point for changes (after constant behaviour) to be inserted into options as <code>options.start_at_t</code>
<code>options_constant_before_start</code>	Boolean (keep curves constant before starting time) to be inserted into options as <code>options.constant_before_start</code>
..	..

..	..
regression_t_interp	time points for regression (optional)
regression_tmin	start time for regression (optional)
regression_tmax	end time for regression (optional)
crossvalidation	run crossvalidation? (Boolean, default 0)
postprocess_normalise	Boolean, default 1
graphics_individual	file basename (used in script replicate_regression_omics_selected)
graphics_scale	default 'log2', 'linear'
graphics_format	'eps', 'png' (for technical reasons, 'eps' needs to be written in single quotes)
convenience_name	type of protein names to be used in graphics (default 'SubtiWiki_20090701')
normalise_by_median	(Boolean, default TRUE)
mark_outliers_percentage	percentage of data points to be marked as outliers based on crossvalidation error

### Additional attributes in options file for individual graphics (function 'replicate\_regression\_omics\_selected')

OPTION	MEANING
graphics_scale	'log2', 'linear'
postprocess_normalise	1
element_id	id (or list, selected by —)
element_name	name (or list, selected by —)
delimiter_symbol	symbol for delimiting list of elements (in element_id, element_name)
title_string	title for graphics
x_label	x label for graphics
y_label	y label for graphics
plot_data	produce plot for data (single element)
plot_replicates	produce plot with replicates (single element)
plot_regression	produce plot for regression curves (single element)
plot_all	produce joint plots for all elements