

# Package ‘MagmaClustR’

August 29, 2022

**Title** Clustering and Prediction using Multi-Task Gaussian Processes with Common Mean

**Version** 1.0.1

**Description** An implementation for the multi-task Gaussian processes with common mean framework. Two main algorithms, called 'Magma' and 'MagmaClust', are available to perform predictions for supervised learning problems, in particular for time series or any functional/continuous data applications. The corresponding articles has been respectively proposed by Arthur Leroy, Pierre Latouche, Benjamin Guedj and Servane Gey (2022) <[doi:10.1007/s10994-022-06172-1](https://doi.org/10.1007/s10994-022-06172-1)>, and Arthur Leroy, Pierre Latouche, Benjamin Guedj and Servane Gey (2020) <[arXiv:2011.07866](https://arxiv.org/abs/2011.07866)>. These approaches leverage the learning of cluster-specific mean processes, which are common across similar tasks, to provide enhanced prediction performances (even far from data) at a linear computational cost (in the number of tasks). 'MagmaClust' is a generalisation of 'Magma' where the tasks are simultaneously clustered into groups, each being associated to a specific mean process. User-oriented functions in the package are decomposed into training, prediction and plotting functions. Some basic features (classic kernels, training, prediction) of standard Gaussian processes are also implemented.

**License** MIT + file LICENSE

**URL** <https://github.com/ArthurLeroy/MagmaClustR>,  
<https://arthurleroy.github.io/MagmaClustR/>

**BugReports** <https://github.com/ArthurLeroy/MagmaClustR/issues>

**Imports** broom, dplyr, ggplot2, magrittr, methods, mvtnorm, Rcpp, rlang, stats, tibble, tidyr, tidyselect

**Suggests** ganimate, gifski, knitr, plotly, png, rmarkdown, testthat (>= 3.0.0), transformr

**LinkingTo** Rcpp

**Encoding** UTF-8

**RoxygenNote** 7.2.0

**NeedsCompilation** yes

**Author** Arthur Leroy [aut, cre] (<<https://orcid.org/0000-0003-0806-8934>>),  
 Pierre Pathé [ctb],  
 Pierre Latouche [aut]

**Maintainer** Arthur Leroy <arthur.leroy.pro@gmail.com>

**Repository** CRAN

**Date/Publication** 2022-08-29 20:20:02 UTC

## R topics documented:

data_allocate_cluster . . . . .	2
hp . . . . .	3
hyperposterior . . . . .	4
hyperposterior_clust . . . . .	6
MagmaClustR . . . . .	8
plot_db . . . . .	9
plot_gif . . . . .	10
plot_gp . . . . .	11
plot_magmaclust . . . . .	13
pred_gif . . . . .	15
pred_gp . . . . .	17
pred_magma . . . . .	19
pred_magmaclust . . . . .	21
proba_max_cluster . . . . .	23
sample_gp . . . . .	24
select_nb_cluster . . . . .	25
simu_db . . . . .	26
train_gp . . . . .	28
train_gp_clust . . . . .	30
train_magma . . . . .	32
train_magmaclust . . . . .	34
<b>Index</b>	<b>38</b>

---

data\_allocate\_cluster *Allocate training data into the most probable cluster*

---

### Description

Allocate training data into the most probable cluster

### Usage

```
data_allocate_cluster(trained_model)
```

**Arguments**

`trained_model` A list, containing the information coming from a MagmaClust model, previously trained using the `train_magmaclust` function.

**Value**

The original dataset used to train the MagmaClust model, with additional 'Cluster' and associated 'Proba' columns, indicating the most probable cluster for each individual/task at the end of the training procedure.

**Examples**

```
TRUE
```

---

hp	<i>Generate random hyper-parameters</i>
----	---

---

**Description**

Generate a set of random hyper-parameters, specific to the chosen type of kernel, under the format that is used in Magma.

**Usage**

```
hp(
  kern = "SE",
  list_ID = NULL,
  list_hp = NULL,
  noise = FALSE,
  common_hp = FALSE
)
```

**Arguments**

`kern` A function, or a character string indicating the chosen type of kernel among:

- "SE": the Squared Exponential kernel,
- "LIN": the Linear kernel,
- "PERIO": the Periodic kernel,
- "RQ": the Rational Quadratic kernel. Compound kernels can be created as sums or products of the above kernels. For combining kernels, simply provide a formula as a character string where elements are separated by whitespaces (e.g. "SE + PERIO"). As the elements are treated sequentially from the left to the right, the product operator '\*' shall always be used before the '+' operators (e.g. 'SE \* LIN + RQ' is valid whereas 'RQ + SE \* LIN' is not).

	In case of a custom kernel function, the argument <code>list_hp</code> has to be provided as well, for designing a tibble with the correct names of hyper-parameters.
<code>list_ID</code>	A vector, associating an ID value with each individual for whom hyper-parameters are generated. If <code>NULL</code> (default) only one set of hyper-parameters is returned without the ID column.
<code>list_hp</code>	A vector of characters, providing the name of each hyper-parameter, in case where <code>kern</code> is a custom kernel function.
<code>noise</code>	A logical value, indicating whether a 'noise' hyper-parameter should be included.
<code>common_hp</code>	A logical value, indicating whether the set of hyper-parameters is assumed to be common to all individuals.

**Value**

A tibble, providing a set of random hyper-parameters associated with the kernel specified through the argument `kern`.

**Examples**

```
TRUE
```

---

`hyperposterior`      *Compute the hyper-posterior distribution in Magma*

---

**Description**

Compute the parameters of the hyper-posterior Gaussian distribution of the mean process in Magma (similarly to the expectation step of the EM algorithm used for learning). This hyper-posterior distribution, evaluated on a grid of inputs provided through the `grid_inputs` argument, is a key component for making prediction in Magma, and is required in the function [pred\\_magma](#).

**Usage**

```
hyperposterior(
  data,
  hp_0,
  hp_i,
  kern_0,
  kern_i,
  prior_mean = NULL,
  grid_inputs = NULL,
  pen_diag = 1e-10
)
```

**Arguments**

data	A tibble or data frame. Required columns: 'Input', 'Output'. Additional columns for covariates can be specified. The 'Input' column should define the variable that is used as reference for the observations (e.g. time for longitudinal data). The 'Output' column specifies the observed values (the response variable). The data frame can also provide as many covariates as desired, with no constraints on the column names. These covariates are additional inputs (explanatory variables) of the models that are also observed at each reference 'Input'.
hp_0	A named vector, tibble or data frame of hyper-parameters associated with kern_0.
hp_i	A tibble or data frame of hyper-parameters associated with kern_i.
kern_0	A kernel function, associated with the mean GP. Several popular kernels (see <a href="#">The Kernel Cookbook</a> ) are already implemented and can be selected within the following list: <ul style="list-style-type: none"> <li>• "SE": (default value) the Squared Exponential Kernel (also called Radial Basis Function or Gaussian kernel),</li> <li>• "LIN": the Linear kernel,</li> <li>• "PERIO": the Periodic kernel,</li> <li>• "RQ": the Rational Quadratic kernel. Compound kernels can be created as sums or products of the above kernels. For combining kernels, simply provide a formula as a character string where elements are separated by whitespaces (e.g. "SE + PERIO"). As the elements are treated sequentially from the left to the right, the product operator '*' shall always be used before the '+' operators (e.g. 'SE * LIN + RQ' is valid whereas 'RQ + SE * LIN' is not).</li> </ul>
kern_i	A kernel function, associated with the individual GPs. ("SE", "PERIO" and "RQ" are also available here)
prior_mean	Hyper-prior mean parameter of the mean GP. This argument, can be specified under various formats, such as: <ul style="list-style-type: none"> <li>• NULL (default). The hyper-prior mean would be set to 0 everywhere.</li> <li>• A number. The hyper-prior mean would be a constant function.</li> <li>• A vector of the same length as all the distinct Input values in the data argument. This vector would be considered as the evaluation of the hyper-prior mean function at the training Inputs.</li> <li>• A function. This function is defined as the hyper-prior mean.</li> <li>• A tibble or data frame. Required columns: Input, Output. The Input values should include at least the same values as in the data argument.</li> </ul>
grid_inputs	A vector, indicating the grid of additional reference inputs on which the mean process' hyper-posterior should be evaluated.
pen_diag	A number. A jitter term, added on the diagonal to prevent numerical issues when inverting nearly singular matrices.

**Value**

A list gathering the parameters of the mean processes' hyper-posterior distributions, namely:

- mean: A tibble, the hyper-posterior mean parameter evaluated at each training Input.
- cov: A matrix, the covariance parameter for the hyper-posterior distribution of the mean process.
- pred: A tibble, the predicted mean and variance at Input for the mean process' hyper-posterior distribution under a format that allows the direct visualisation as a GP prediction.

## Examples

TRUE

---

hyperposterior\_clust *Compute the hyper-posterior distribution for each cluster in MagmaClust*

---

## Description

Recompute the E-step of the VEM algorithm in MagmaClust for a new set of reference Input. Once training is completed, it can be necessary to evaluate the hyper-posterior distributions of the mean processes at specific locations, for which we want to make predictions. This process is directly implemented in the [pred\\_magmaclust](#) function but for the user might want to use hyperpost\_clust for a tailored control 'by hand' of the prediction procedure.

## Usage

```
hyperposterior_clust(
  data,
  mixture,
  hp_k,
  hp_i,
  kern_k,
  kern_i,
  prior_mean_k = NULL,
  grid_inputs = NULL,
  pen_diag = 1e-10
)
```

## Arguments

data	A tibble or data frame. Required columns: ID, Input , Output. Additional columns for covariates can be specified. The ID column contains the unique names/codes used to identify each individual/task (or batch of data). The Input column should define the variable that is used as reference for the observations (e.g. time for longitudinal data). The Output column specifies the observed values (the response variable). The data frame can also provide as many covariates as desired, with no constraints on the column names. These covariates are additional inputs (explanatory variables) of the models that are also observed at each reference Input.
------	---

mixture	A tibble or data frame, indicating the mixture probabilities of each cluster for each individual. Required column: ID.
hp_k	A tibble or data frame of hyper-parameters associated with kern_k.
hp_i	A tibble or data frame of hyper-parameters associated with kern_i.
kern_k	A kernel function, associated with the mean GPs. Several popular kernels (see <a href="#">The Kernel Cookbook</a> ) are already implemented and can be selected within the following list: <ul style="list-style-type: none"> <li>• "SE": (default value) the Squared Exponential Kernel (also called Radial Basis Function or Gaussian kernel),</li> <li>• "LIN": the Linear kernel,</li> <li>• "PERIO": the Periodic kernel,</li> <li>• "RQ": the Rational Quadratic kernel. Compound kernels can be created as sums or products of the above kernels. For combining kernels, simply provide a formula as a character string where elements are separated by whitespaces (e.g. "SE + PERIO"). As the elements are treated sequentially from the left to the right, the product operator '*' shall always be used before the '+' operators (e.g. 'SE * LIN + RQ' is valid whereas 'RQ + SE * LIN' is not).</li> </ul>
kern_i	A kernel function, associated with the individual GPs. ("SE", "LIN", "PERIO" and "RQ" are also available here)
prior_mean_k	The set of hyper-prior mean parameters (m_k) for the K mean GPs, one value for each cluster. This argument can be specified under various formats, such as: <ul style="list-style-type: none"> <li>• NULL (default). All hyper-prior means would be set to 0 everywhere.</li> <li>• A numerical vector of the same length as the number of clusters. Each number is associated with one cluster, and considered to be the hyper-prior mean parameter of the cluster (i.e. a constant function at all Input).</li> <li>• A list of functions. Each function is associated with one cluster. These functions are all evaluated at all Input values, to provide specific hyper-prior mean vectors for each cluster.</li> </ul>
grid_inputs	A vector, indicating the grid of additional reference inputs on which the mean process' hyper-posterior should be evaluated.
pen_diag	A number. A jitter term, added on the diagonal to prevent numerical issues when inverting nearly singular matrices.

## Value

A list containing the parameters of the mean processes' hyper-posterior distribution, namely:

- mean: A list of tibbles containing, for each cluster, the hyper-posterior mean parameters evaluated at each Input.
- cov: A list of matrices containing, for each cluster, the hyper-posterior covariance parameter of the mean process.
- mixture: A tibble, indicating the mixture probabilities in each cluster for each individual.

## Examples

```
TRUE
```

---

MagmaClustR	<i>MagmaClustR : Clustering and Prediction using Multi-Task Gaussian Processes</i>
-------------	--

---

## Description

The **MagmaClustR** package implements two main algorithms, called *Magma* and *MagmaClust*, using a multi-task GPs model to perform predictions for supervised learning problems. These approaches leverage the learning of cluster-specific mean processes, which are common across similar tasks, to provide enhanced prediction performances (even far from data) at a linear computational cost (in the number of tasks). *MagmaClust* is a generalisation of *Magma* where the tasks are simultaneously clustered into groups, each being associated to a specific mean process. User-oriented functions in the package are decomposed into training, prediction and plotting functions. Some basic features of standard GPs are also implemented.

## Details

For a quick introduction to **MagmaClustR**, please refer to the README at <https://github.com/ArthurLeroy/MagmaClustR>

## Author(s)

Arthur Leroy, Pierre Pathe and Pierre Latouche  
 Maintainer: Arthur Leroy - <arthur.leroy.pro@gmail.com>

## References

Arthur Leroy, Pierre Latouche, Benjamin Guedj, and Servane Gey.  
 MAGMA: Inference and Prediction with Multi-Task Gaussian Processes. *Machine Learning*, 2022,  
<https://link.springer.com/article/10.1007/s10994-022-06172-1>

Arthur Leroy, Pierre Latouche, Benjamin Guedj, and Servane Gey.  
 Cluster-Specific Predictions with Multi-Task Gaussian Processes. *PREPRINT*, Nov. 2020, <https://arxiv.org/abs/2011.07866>

## Examples

### Simulate a dataset, train and predict with Magma

```
:
set.seed(42)
data_magma <- simu_db(M = 11, N = 10, K = 1)
magma_train <- data_magma %>% subset(ID %in% 1:10)
magma_test <- data_magma %>% subset(ID == 11) %>% head(5)
```



```
magma_model <- train_magma(data = magma_train)
magma_pred <- pred_magma(data = magma_test, trained_model = magma_model, grid_inputs =
seq(0, 10, 0.01))
```

### Simulate a dataset, train and predict with MagmaClust

```
:
set.seed(42)
data_magmaclust <- simu_db(M = 4, N = 10, K = 3)
list_ID = unique(data_magmaclust$ID)
magmaclust_train <- data_magmaclust %>% subset(ID %in% list_ID[1:11])
magmaclust_test <- data_magmaclust %>% subset(ID == list_ID[12]) %>% head(5)

magmaclust_model <- train_magmaclust(data = magmaclust_train)
magmaclust_pred <- pred_magmaclust(data = magmaclust_test,
trained_model = magmaclust_model, grid_inputs = seq(0, 10, 0.01))
```

---

plot\_db

*Plot smoothed curves of raw data*

---

## Description

Display raw data under the Magma format as smoothed curves.

## Usage

```
plot_db(data, cluster = FALSE, legend = FALSE)
```

## Arguments

data	A data frame or tibble with format : ID, Input, Output.
cluster	A boolean indicating whether data should be coloured by cluster. Requires a column named 'Cluster'.
legend	A boolean indicating whether the legend should be displayed.

## Value

Graph of smoothed curves of raw data.

## Examples

```
TRUE
```

---

 plot\_gif

 Create a GIF of Magma or GP predictions
 

---

### Description

Create a GIF animation displaying how Magma or classic GP predictions evolve and improve when the number of data points increase.

### Usage

```
plot_gif(
  pred_gp,
  x_input = NULL,
  data = NULL,
  data_train = NULL,
  prior_mean = NULL,
  y_grid = NULL,
  heatmap = FALSE,
  prob_CI = 0.95,
  size_data = 3,
  size_data_train = 1,
  alpha_data_train = 0.5,
  export_gif = FALSE,
  path = "gif_gp.gif",
  ...
)
```

### Arguments

pred_gp	A tibble, typically coming from the <a href="#">pred_gif</a> function. Required columns: 'Input', 'Mean', 'Var' and 'Index'.
x_input	A vector of character strings, indicating which input should be displayed. If NULL(default) the 'Input' column is used for the x-axis. If providing a 2-dimensional vector, the corresponding columns are used for the x-axis and y-axis.
data	(Optional) A tibble or data frame. Required columns: 'Input', 'Output'. Additional columns for covariates can be specified. The 'Input' column should define the variable that is used as reference for the observations (e.g. time for longitudinal data). The 'Output' column specifies the observed values (the response variable). The data frame can also provide as many covariates as desired, with no constraints on the column names. These covariates are additional inputs (explanatory variables) of the models that are also observed at each reference 'Input'.
data_train	(Optional) A tibble or data frame, containing the training data of the Magma model. The data set should have the same format as the data argument with

	an additional column 'ID' for identifying the different individuals/tasks. If provided, those data are displayed as backward colourful points (each colour corresponding to one individual/task).
prior_mean	(Optional) A tibble or a data frame, containing the 'Input' and associated 'Output' prior mean parameter of the GP prediction.
y_grid	A vector, indicating the grid of values on the y-axis for which probabilities should be computed for heatmaps of 1-dimensional predictions. If NULL (default), a vector of length 50 is defined, ranging between the min and max 'Output' values contained in pred_gp.
heatmap	A logical value indicating whether the GP prediction should be represented as a heatmap of probabilities for 1-dimensional inputs. If FALSE (default), the mean curve and associated 95% CI are displayed.
prob_CI	A number between 0 and 1 (default is 0.95), indicating the level of the Credible Interval associated with the posterior mean curve.
size_data	A number, controlling the size of the data points.
size_data_train	A number, controlling the size of the data_train points.
alpha_data_train	A number, between 0 and 1, controlling transparency of the data_train points.
export_gif	A logical value indicating whether the animation should be exported as a .gif file.
path	A character string defining the path where the GIF file should be exported.
...	Any additional parameters that can be passed to the function <a href="#">transition_states</a> from the <code>gganimate</code> package.

### Value

Visualisation of a Magma or GP prediction (optional: display data points, training data points and the prior mean function), where data points are added sequentially for visualising changes in prediction as information increases.

### Examples

```
TRUE
```

---

plot\_gp

*Plot Magma or GP predictions*

---

### Description

Display Magma or classic GP predictions. According to the dimension of the inputs, the graph may be a mean curve + Credible Interval or a heatmap of probabilities.

**Usage**

```
plot_gp(
  pred_gp,
  x_input = NULL,
  data = NULL,
  data_train = NULL,
  prior_mean = NULL,
  y_grid = NULL,
  heatmap = FALSE,
  prob_CI = 0.95,
  size_data = 3,
  size_data_train = 1,
  alpha_data_train = 0.5
)
```

**Arguments**

pred_gp	A tibble or data frame, typically coming from <code>pred_magma</code> or <code>pred_gp</code> functions. Required columns: 'Input', 'Mean', 'Var'. Additional covariate columns may be present in case of multi-dimensional inputs.
x_input	A vector of character strings, indicating which input should be displayed. If NULL (default) the 'Input' column is used for the x-axis. If providing a 2-dimensional vector, the corresponding columns are used for the x-axis and y-axis.
data	(Optional) A tibble or data frame. Required columns: 'Input', 'Output'. Additional columns for covariates can be specified. This argument corresponds to the raw data on which the prediction has been performed.
data_train	(Optional) A tibble or data frame, containing the training data of the Magma model. The data set should have the same format as the data argument with an additional required column 'ID' for identifying the different individuals/tasks. If provided, those data are displayed as backward colourful points (each colour corresponding to one individual/task).
prior_mean	(Optional) A tibble or a data frame, containing the 'Input' and associated 'Output' prior mean parameter of the GP prediction.
y_grid	A vector, indicating the grid of values on the y-axis for which probabilities should be computed for heatmaps of 1-dimensional predictions. If NULL (default), a vector of length 50 is defined, ranging between the min and max 'Output' values contained in <code>pred_gp</code> .
heatmap	A logical value indicating whether the GP prediction should be represented as a heatmap of probabilities for 1-dimensional inputs. If FALSE (default), the mean curve and associated Credible Interval are displayed.
prob_CI	A number between 0 and 1 (default is 0.95), indicating the level of the Credible Interval associated with the posterior mean curve. If this this argument is set to 1, the Credible Interval is not displayed.
size_data	A number, controlling the size of the data points.

size\_data\_train

A number, controlling the size of the data\_train points.

alpha\_data\_train

A number, between 0 and 1, controlling transparency of the data\_train points.

### Value

Visualisation of a Magma or GP prediction (optional: display data points, training data points and the prior mean function). For 1-D inputs, the prediction is represented as a mean curve and its associated 95% Credible Interval, or as a heatmap of probabilities if heatmap = TRUE. For 2-D inputs, the prediction is represented as a heatmap, where each couple of inputs on the x-axis and y-axis are associated with a gradient of colours for the posterior mean values, whereas the uncertainty is indicated by the transparency (the narrower is the Credible Interval, the more opaque is the associated colour, and vice versa)

### Examples

TRUE

---

plot\_magmaclust      *Plot MagmaClust predictions*

---

### Description

Display MagmaClust predictions. According to the dimension of the inputs, the graph may be a mean curve (dim inputs = 1) or a heatmap (dim inputs = 2) of probabilities. Moreover, MagmaClust can provide credible intervals only by visualising cluster-specific predictions (e.g. for the most probable cluster). When visualising the full mixture-of-GPs prediction, which can be multimodal, the user should choose between the simple mean function or the full heatmap of probabilities (more informative but slower).

### Usage

```
plot_magmaclust(
  pred_clust,
  cluster = "all",
  x_input = NULL,
  data = NULL,
  data_train = NULL,
  col_clust = FALSE,
  prior_mean = NULL,
  y_grid = NULL,
  heatmap = FALSE,
  prob_CI = 0.95,
  size_data = 3,
  size_data_train = 1,
  alpha_data_train = 0.5
)
```

**Arguments**

pred_clust	A list of predictions, typically coming from <a href="#">pred_magmaclust</a> . Required elements: pred, mixture, mixture_pred.
cluster	A character string, indicating which cluster to plot from. If 'all' (default) the mixture of GPs prediction is displayed as a mean curve (1-D inputs) or a mean heatmap (2-D inputs). Alternatively, if the name of one cluster is provided, the classic mean curve + credible interval is displayed (1-D inputs), or a heatmap with colour gradient for the mean and transparency gradient for the Credible Interval (2-D inputs).
x_input	A vector of character strings, indicating which input should be displayed. If NULL (default) the 'Input' column is used for the x-axis. If providing a 2-dimensional vector, the corresponding columns are used for the x-axis and y-axis.
data	(Optional) A tibble or data frame. Required columns: Input , Output. Additional columns for covariates can be specified. This argument corresponds to the raw data on which the prediction has been performed.
data_train	(Optional) A tibble or data frame, containing the training data of the MagmaClust model. The data set should have the same format as the data argument with an additional required column ID for identifying the different individuals/tasks. If provided, those data are displayed as backward colourful points (each colour corresponding to one individual or a cluster, see col_clust below).
col_clust	A boolean indicating whether backward points are coloured according to the individuals or to their most probable cluster. If one wants to colour by clusters, a column Cluster shall be present in data_train. We advise to use <a href="#">data_allocate_cluster</a> for automatically creating a well-formatted dataset from a trained MagmaClust model.
prior_mean	(Optional) A list providing, for each cluster, a tibble containing prior mean parameters of the prediction. This argument typically comes as an outcome hyperpost\$mean, available through the <a href="#">train_magmaclust</a> , <a href="#">pred_magmaclust</a> functions.
y_grid	A vector, indicating the grid of values on the y-axis for which probabilities should be computed for heatmaps of 1-dimensional predictions. If NULL (default), a vector of length 50 is defined, ranging between the min and max 'Output' values contained in pred.
heatmap	A logical value indicating whether the GP prediction should be represented as a heatmap of probabilities for 1-dimensional inputs. If FALSE (default), the mean curve (and associated Credible Interval if available) are displayed.
prob_CI	A number between 0 and 1 (default is 0.95), indicating the level of the Credible Interval associated with the posterior mean curve. If this this argument is set to 1, the Credible Interval is not displayed.
size_data	A number, controlling the size of the data points.
size_data_train	A number, controlling the size of the data_train points.
alpha_data_train	A number, between 0 and 1, controlling transparency of the data_train points.

**Value**

Visualisation of a MagmaClust prediction (optional: display data points, training data points and the prior mean functions). For 1-D inputs, the prediction is represented as a mean curve (and its associated 95% Credible Interval for cluster-specific predictions), or as a heatmap of probabilities if heatmap = TRUE. In the case of MagmaClust, the heatmap representation should be preferred for clarity, although the default display remains mean curve for quicker execution. For 2-D inputs, the prediction is represented as a heatmap, where each couple of inputs on the x-axis and y-axis are associated with a gradient of colours for the posterior mean values, whereas the uncertainty is indicated by the transparency (the narrower is the Credible Interval, the more opaque is the associated colour, and vice versa). As for 1-D inputs, Credible Interval information is only available for cluster-specific predictions.

**Examples**

```
TRUE
```

---

```
pred_gif
```

```
Magma prediction for plotting GIFs
```

---

**Description**

Generate a Magma or classic GP prediction under a format that is compatible with a further GIF visualisation of the results. For a Magma prediction, either the trained\_model or hyperpost argument is required. Otherwise, a classic GP prediction is applied and the prior mean can be specified through the mean argument.

**Usage**

```
pred_gif(
  data,
  trained_model = NULL,
  grid_inputs = NULL,
  hyperpost = NULL,
  mean = NULL,
  hp = NULL,
  kern = "SE",
  pen_diag = 1e-10
)
```

**Arguments**

data	A tibble or data frame. Required columns: 'Input', 'Output'. Additional columns for covariates can be specified. The 'Input' column should define the variable that is used as reference for the observations (e.g. time for longitudinal data). The 'Output' column specifies the observed values (the response variable). The data frame can also provide as many covariates as desired, with no constraints on the column names. These covariates are additional inputs (explanatory variables) of the models that are also observed at each reference 'Input'.
------	--

trained_model	A list, containing the information coming from a Magma model, previously trained using the <code>train_magma</code> function.
grid_inputs	The grid of inputs (reference Input and covariates) values on which the GP should be evaluated. Ideally, this argument should be a tibble or a data frame, providing the same columns as data, except 'Output'. Nonetheless, in cases where data provides only one 'Input' column, the <code>grid_inputs</code> argument can be NULL (default) or a vector. This vector would be used as reference input for prediction and if NULL, a vector of length 500 is defined, ranging between the min and max Input values of data.
hyperpost	A list, containing the elements 'mean' and 'cov', the parameters of the hyper-posterior distribution of the mean process. Typically, this argument should come from a previous learning using <code>train_magma</code> , or a previous prediction with <code>pred_magma</code> , with the argument <code>get_hyperpost</code> set to TRUE. The 'mean' element should be a data frame with two columns 'Input' and 'Output'. The 'cov' element should be a covariance matrix with colnames and rownames corresponding to the 'Input' in 'mean'. In all cases, the column 'Input' should contain all the values appearing both in the 'Input' column of data and in <code>grid_inputs</code> .
mean	Mean parameter of the GP. This argument can be specified under various formats, such as: <ul style="list-style-type: none"> <li>• NULL (default). The mean would be set to 0 everywhere.</li> <li>• A number. The mean would be a constant function.</li> <li>• A function. This function is defined as the mean.</li> <li>• A tibble or data frame. Required columns: Input, Output. The Input values should include at least the same values as in the data argument.</li> </ul>
hp	A named vector, tibble or data frame of hyper-parameters associated with kern. The columns/elements should be named according to the hyper-parameters that are used in kern. The function <code>train_gp</code> can be used to learn maximum-likelihood estimators of the hyper-parameters,
kern	A kernel function, defining the covariance structure of the GP. Several popular kernels (see <a href="#">The Kernel Cookbook</a> ) are already implemented and can be selected within the following list: <ul style="list-style-type: none"> <li>• "SE": (default value) the Squared Exponential Kernel (also called Radial Basis Function or Gaussian kernel),</li> <li>• "LIN": the Linear kernel,</li> <li>• "PERIO": the Periodic kernel,</li> <li>• "RQ": the Rational Quadratic kernel. Compound kernels can be created as sums or products of the above kernels. For combining kernels, simply provide a formula as a character string where elements are separated by whitespaces (e.g. "SE + PERIO"). As the elements are treated sequentially from the left to the right, the product operator '*' shall always be used before the '+' operators (e.g. 'SE * LIN + RQ' is valid whereas 'RQ + SE * LIN' is not).</li> </ul>
pen_diag	A number. A jitter term, added on the diagonal to prevent numerical issues when inverting nearly singular matrices.



**Value**

A tibble, representing Magma or GP predictions as two column 'Mean' and 'Var', evaluated on the `grid_inputs`. The column 'Input' and additional covariates columns are associated to each predicted values. An additional 'Index' column is created for the sake of GIF creation using the function `plot_gif`

**Examples**

```
TRUE
```

---

pred_gp	<i>Gaussian Process prediction</i>
---------	------------------------------------

---

**Description**

Compute the posterior distribution of a simple GP, using the formalism of Magma. By providing observed data, the prior mean and covariance matrix (by defining a kernel and its associated hyper-parameters), the mean and covariance parameters of the posterior distribution are computed on the grid of inputs that has been specified. This predictive distribution can be evaluated on any arbitrary inputs since a GP is an infinite-dimensional object.

**Usage**

```
pred_gp(
  data,
  grid_inputs = NULL,
  mean = NULL,
  hp = NULL,
  kern = "SE",
  get_full_cov = FALSE,
  plot = TRUE,
  pen_diag = 1e-10
)
```

**Arguments**

<code>data</code>	A tibble or data frame. Required columns: 'Input', 'Output'. Additional columns for covariates can be specified. The 'Input' column should define the variable that is used as reference for the observations (e.g. time for longitudinal data). The 'Output' column specifies the observed values (the response variable). The data frame can also provide as many covariates as desired, with no constraints on the column names. These covariates are additional inputs (explanatory variables) of the models that are also observed at each reference 'Input'.
<code>grid_inputs</code>	The grid of inputs (reference Input and covariates) values on which the GP should be evaluated. Ideally, this argument should be a tibble or a data frame, providing the same columns as <code>data</code> , except 'Output'. Nonetheless, in cases where <code>data</code> provides only one 'Input' column, the <code>grid_inputs</code> argument can

	be NULL (default) or a vector. This vector would be used as reference input for prediction and if NULL, a vector of length 500 is defined, ranging between the min and max Input values of data.
mean	<p>Mean parameter of the GP. This argument can be specified under various formats, such as:</p> <ul style="list-style-type: none"> <li>• NULL (default). The mean would be set to 0 everywhere.</li> <li>• A number. The mean would be a constant function.</li> <li>• A function. This function is defined as the mean.</li> <li>• A tibble or data frame. Required columns: Input, Output. The Input values should include at least the same values as in the data argument.</li> </ul>
hp	A named vector, tibble or data frame of hyper-parameters associated with kern. The columns/elements should be named according to the hyper-parameters that are used in kern. If NULL (default), the function <code>train_gp</code> is called with random initial values for learning maximum-likelihood estimators of the hyper-parameters associated with kern.
kern	<p>A kernel function, defining the covariance structure of the GP. Several popular kernels (see <a href="#">The Kernel Cookbook</a>) are already implemented and can be selected within the following list:</p> <ul style="list-style-type: none"> <li>• "SE": (default value) the Squared Exponential Kernel (also called Radial Basis Function or Gaussian kernel),</li> <li>• "LIN": the Linear kernel,</li> <li>• "PERIO": the Periodic kernel,</li> <li>• "RQ": the Rational Quadratic kernel. Compound kernels can be created as sums or products of the above kernels. For combining kernels, simply provide a formula as a character string where elements are separated by whitespaces (e.g. "SE + PERIO"). As the elements are treated sequentially from the left to the right, the product operator '*' shall always be used before the '+' operators (e.g. 'SE * LIN + RQ' is valid whereas 'RQ + SE * LIN' is not).</li> </ul>
get_full_cov	A logical value, indicating whether the full posterior covariance matrix should be returned.
plot	A logical value, indicating whether a plot of the results is automatically displayed.
pen_diag	A number. A jitter term, added on the diagonal to prevent numerical issues when inverting nearly singular matrices.

### Value

A tibble, representing the GP predictions as two column 'Mean' and 'Var', evaluated on the `grid_inputs`. The column 'Input' and additional covariates columns are associated to each predicted values. If the `get_full_cov` argument is TRUE, the function returns a list, in which the tibble described above is defined as 'pred' and the full posterior covariance matrix is defined as 'cov'.

### Examples

```
TRUE
```

---

pred_magma	<i>Magma prediction</i>
------------	-------------------------

---

### Description

Compute the posterior predictive distribution in Magma. Providing data of any new individual/task, its trained hyper-parameters and a previously trained Magma model, the predictive distribution is evaluated on any arbitrary inputs that are specified through the 'grid\_inputs' argument.

### Usage

```
pred_magma(
  data,
  trained_model = NULL,
  grid_inputs = NULL,
  hp = NULL,
  kern = "SE",
  hyperpost = NULL,
  get_hyperpost = FALSE,
  get_full_cov = FALSE,
  plot = TRUE,
  pen_diag = 1e-10
)
```

### Arguments

data	A tibble or data frame. Required columns: 'Input', 'Output'. Additional columns for covariates can be specified. The 'Input' column should define the variable that is used as reference for the observations (e.g. time for longitudinal data). The 'Output' column specifies the observed values (the response variable). The data frame can also provide as many covariates as desired, with no constraints on the column names. These covariates are additional inputs (explanatory variables) of the models that are also observed at each reference 'Input'.
trained_model	A list, containing the information coming from a Magma model, previously trained using the <a href="#">train_magma</a> function.
grid_inputs	The grid of inputs (reference Input and covariates) values on which the GP should be evaluated. Ideally, this argument should be a tibble or a data frame, providing the same columns as data, except 'Output'. Nonetheless, in cases where data provides only one 'Input' column, the grid_inputs argument can be NULL (default) or a vector. This vector would be used as reference input for prediction and if NULL, a vector of length 500 is defined, ranging between the min and max Input values of data.
hp	A named vector, tibble or data frame of hyper-parameters associated with kern. The columns/elements should be named according to the hyper-parameters that are used in kern. The function <a href="#">train_gp</a> can be used to learn maximum-likelihood estimators of the hyper-parameters.

kern	<p>A kernel function, defining the covariance structure of the GP. Several popular kernels (see <a href="#">The Kernel Cookbook</a>) are already implemented and can be selected within the following list:</p> <ul style="list-style-type: none"> <li>• "SE": (default value) the Squared Exponential Kernel (also called Radial Basis Function or Gaussian kernel),</li> <li>• "LIN": the Linear kernel,</li> <li>• "PERIO": the Periodic kernel,</li> <li>• "RQ": the Rational Quadratic kernel. Compound kernels can be created as sums or products of the above kernels. For combining kernels, simply provide a formula as a character string where elements are separated by whitespaces (e.g. "SE + PERIO"). As the elements are treated sequentially from the left to the right, the product operator '*' shall always be used before the '+' operators (e.g. 'SE * LIN + RQ' is valid whereas 'RQ + SE * LIN' is not).</li> </ul>
hyperpost	<p>A list, containing the elements 'mean' and 'cov', the parameters of the hyper-posterior distribution of the mean process. Typically, this argument should come from a previous learning using <a href="#">train_magma</a>, or a previous prediction with <a href="#">pred_magma</a>, with the argument <code>get_hyperpost</code> set to TRUE. The 'mean' element should be a data frame with two columns 'Input' and 'Output'. The 'cov' element should be a covariance matrix with colnames and rownames corresponding to the 'Input' in 'mean'. In all cases, the column 'Input' should contain all the values appearing both in the 'Input' column of data and in <code>grid_inputs</code>.</p>
get_hyperpost	<p>A logical value, indicating whether the hyper-posterior distribution of the mean process should be returned. This can be useful when planning to perform several predictions on the same grid of inputs, since recomputation of the hyper-posterior can be prohibitive for high dimensional grids.</p>
get_full_cov	<p>A logical value, indicating whether the full posterior covariance matrix should be returned.</p>
plot	<p>A logical value, indicating whether a plot of the results is automatically displayed.</p>
pen_diag	<p>A number. A jitter term, added on the diagonal to prevent numerical issues when inverting nearly singular matrices.</p>

### Value

A tibble, representing Magma predictions as two column 'Mean' and 'Var', evaluated on the `grid_inputs`. The column 'Input' and additional covariates columns are associated to each predicted values. If the `get_full_cov` or `get_hyperpost` arguments are TRUE, the function returns a list, in which the tibble described above is defined as 'pred\_gp' and the full posterior covariance matrix is defined as 'cov', and the hyper-posterior distribution of the mean process is defined as 'hyperpost'.

### Examples

```
TRUE
```

---

pred\_magmaclust      *MagmaClust prediction*

---

## Description

Compute the posterior predictive distribution in MagmaClust. Providing data from any new individual/task, its trained hyper-parameters and a previously trained MagmaClust model, the multi-task posterior distribution is evaluated on any arbitrary inputs that are specified through the 'grid\_inputs' argument. Due to the nature of the model, the prediction is defined as a mixture of Gaussian distributions. Therefore the present function computes the parameters of the predictive distribution associated with each cluster, as well as the posterior mixture probabilities for this new individual/task.

## Usage

```
pred_magmaclust(
  data,
  trained_model = NULL,
  grid_inputs = NULL,
  mixture = NULL,
  hp = NULL,
  kern = "SE",
  hyperpost = NULL,
  prop_mixture = NULL,
  get_hyperpost = FALSE,
  get_full_cov = FALSE,
  plot = TRUE,
  pen_diag = 1e-10
)
```

## Arguments

data	A tibble or data frame. Required columns: Input, Output. Additional columns for covariates can be specified. The Input column should define the variable that is used as reference for the observations (e.g. time for longitudinal data). The Output column specifies the observed values (the response variable). The data frame can also provide as many covariates as desired, with no constraints on the column names. These covariates are additional inputs (explanatory variables) of the models that are also observed at each reference 'Input'.
trained_model	A list, containing the information coming from a MagmaClust model, previously trained using the <a href="#">train_magmaclust</a> function. If trained_model is set to NULL, the hyperpost and prop_mixture arguments are mandatory to perform required re-computations for the prediction to succeed.
grid_inputs	The grid of inputs (reference Input and covariates) values on which the GP should be evaluated. Ideally, this argument should be a tibble or a data frame, providing the same columns as data, except 'Output'. Nonetheless, in cases

where data provides only one 'Input' column, the `grid_inputs` argument can be NULL (default) or a vector. This vector would be used as reference input for prediction and if NULL, a vector of length 500 is defined, ranging between the min and max Input values of data.

mixture	A tibble or data frame, indicating the mixture probabilities of each cluster for the new individual/task. If NULL, the <code>train_gp_clust</code> function is used to compute these posterior probabilities according to data.
hp	A named vector, tibble or data frame of hyper-parameters associated with kern. The columns/elements should be named according to the hyper-parameters that are used in kern. The <code>train_gp_clust</code> function can be used to learn maximum-likelihood estimators of the hyper-parameters.
kern	<p>A kernel function, defining the covariance structure of the GP. Several popular kernels (see <a href="#">The Kernel Cookbook</a>) are already implemented and can be selected within the following list:</p> <ul style="list-style-type: none"> <li>• "SE": (default value) the Squared Exponential Kernel (also called Radial Basis Function or Gaussian kernel),</li> <li>• "LIN": the Linear kernel,</li> <li>• "PERIO": the Periodic kernel,</li> <li>• "RQ": the Rational Quadratic kernel. Compound kernels can be created as sums or products of the above kernels. For combining kernels, simply provide a formula as a character string where elements are separated by whitespaces (e.g. "SE + PERIO"). As the elements are treated sequentially from the left to the right, the product operator '*' shall always be used before the '+' operators (e.g. 'SE * LIN + RQ' is valid whereas 'RQ + SE * LIN' is not).</li> </ul>
hyperpost	A list, containing the elements mean, cov and mixture the parameters of the hyper-posterior distributions of the mean processes. Typically, this argument should come from a previous learning using <code>train_magmaclust</code> , or a previous prediction with <code>pred_magmaclust</code> , with the argument <code>get_hyperpost</code> set to TRUE.
prop_mixture	A tibble or a named vector of the mixture proportions. Each name of column or element should refer to a cluster. The value associated with each cluster is a number between 0 and 1. If both mixture and <code>trained_model</code> are set to NULL, this argument allows to recompute mixture probabilities, thanks to the <code>hyperpost</code> argument and the <code>train_gp_clust</code> function.
get_hyperpost	A logical value, indicating whether the hyper-posterior distributions of the mean processes should be returned. This can be useful when planning to perform several predictions on the same grid of inputs, since recomputation of the hyper-posterior can be prohibitive for high dimensional grids.
get_full_cov	A logical value, indicating whether the full posterior covariance matrices should be returned.
plot	A logical value, indicating whether a plot of the results is automatically displayed.
pen_diag	A number. A jitter term, added on the diagonal to prevent numerical issues when inverting nearly singular matrices.

**Value**

A list of GP prediction results composed of:

- pred: As sub-list containing, for each cluster:
  - pred\_gp: A tibble, representing the GP predictions as two column Mean and Var, evaluated on the grid\_inputs. The column Input and additional covariates columns are associated with each predicted values.
  - proba: A number, the posterior probability associated with this cluster.
  - cov (if get\_full\_cov = TRUE): A matrix, the full posterior covariance matrix associated with this cluster.
- mixture: A tibble, indicating the mixture probabilities of each cluster for the predicted individual/task.
- hyperpost (if get\_hyperpost = TRUE): A list, containing the hyper-posterior distributions information useful for visualisation purposes.

**Examples**

TRUE

---

<code>proba_max_cluster</code>	<i>Indicates the most probable cluster</i>
--------------------------------	--

---

**Description**

Indicates the most probable cluster

**Usage**

```
proba_max_cluster(mixture)
```

**Arguments**

`mixture`            A tibble or data frame containing mixture probabilities.

**Value**

A tibble, retaining only the most probable cluster. The column Cluster indicates the the cluster's name whereas Proba refers to its associated probability. If ID is initially a column of mixture (optional), the function returns the most probable cluster for all the different ID values.

**Examples**

TRUE

sample\_gp

*Display realisations from a posterior GP***Description**

A realisation of a posterior GP distribution is drawn and displayed. According to the dimension of the inputs, the graph may be a curve or a heatmap.

**Usage**

```
sample_gp(
  pred_gp,
  x_input = NULL,
  data = NULL,
  data_train = NULL,
  prior_mean = NULL,
  size_data = 3,
  size_data_train = 1,
  alpha_data_train = 0.5
)
```

**Arguments**

pred_gp	A tibble or data frame, typically coming from <a href="#">pred_magma</a> or <a href="#">pred_gp</a> functions. Required columns: 'Input', 'Mean', 'Var'. Additional covariate columns may be present in case of multi-dimensional inputs.
x_input	A vector of character strings, indicating which input should be displayed. If NULL(default) the 'Input' column is used for the x-axis. If providing a 2-dimensional vector, the corresponding columns are used for the x-axis and y-axis.
data	(Optional) A tibble or data frame, containing the data used in the GP prediction.
data_train	(Optional) A tibble or data frame, containing the training data of the Magma model. The data set should have the same format as the data argument with an additional column 'ID' for identifying the different individuals/tasks. If provided, those data are displayed as backward colourful points (each colour corresponding to one individual/task).
prior_mean	(Optional) A tibble or a data frame, containing the 'Input' and associated 'Output' prior mean parameter of the GP prediction.
size_data	A number, controlling the size of the data points.
size_data_train	A number, controlling the size of the data_train points.
alpha_data_train	A number, between 0 and 1, controlling transparency of the data_train points.



**Value**

Draw and visualise from a posterior distribution from Magma or GP prediction (optional: display data points, training data points and the prior mean function).

**Examples**

```
TRUE
```

---

select_nb_cluster	<i>Select the optimal number of clusters</i>
-------------------	--

---

**Description**

In MagmaClust, as for any clustering method, the number  $K$  of clusters has to be provided as an hypothesis of the model. This function implements a model selection procedure, by maximising a variational BIC criterion, computed for different values of  $K$ . A heuristic for a fast approximation of the procedure is proposed as well, although the corresponding models would not be properly trained.

**Usage**

```
select_nb_cluster(
  data,
  fast_approx = TRUE,
  grid_nb_cluster = 1:10,
  ini_hp_k = NULL,
  ini_hp_i = NULL,
  kern_k = "SE",
  kern_i = "SE",
  plot = TRUE,
  ...
)
```

**Arguments**

data	A tibble or data frame. Columns required: ID, Input , Output. Additional columns for covariates can be specified. The ID column contains the unique names/codes used to identify each individual/task (or batch of data). The Input column should define the variable that is used as reference for the observations (e.g. time for longitudinal data). The Output column specifies the observed values (the response variable). The data frame can also provide as many covariates as desired, with no constraints on the column names. These covariates are additional inputs (explanatory variables) of the models that are also observed at each reference Input.
------	---

<code>fast_approx</code>	A boolean, indicating whether a fast approximation should be used for selecting the number of clusters. If TRUE, each Magma or MagmaClust model will perform only one E-step of the training, using the same fixed values for the hyper-parameters ( <code>ini_hp_k</code> and <code>ini_hp_i</code> , or random values if not provided) in all models. The resulting models should not be considered as trained, but this approach provides a convenient heuristic to avoid a cumbersome model selection procedure.
<code>grid_nb_cluster</code>	A vector of integer, corresponding to grid of values that will be tested for the number of clusters.
<code>ini_hp_k</code>	A tibble or data frame of hyper-parameters associated with <code>kern_k</code> .
<code>ini_hp_i</code>	A tibble or data frame of hyper-parameters associated with <code>kern_i</code> .
<code>kern_k</code>	A kernel function associated to the mean processes.
<code>kern_i</code>	A kernel function associated to the individuals/tasks.
<code>plot</code>	A boolean indicating whether the plot of V-BIC values for all numbers of clusters should be displayed.
<code>...</code>	Any additional argument that could be passed to <code>train_magmaclust</code> .

### Value

A list, containing the results of model selection procedure for selecting the optimal number of clusters thanks to a V-BIC criterion maximisation. The elements of the list are:

- `best_k`: An integer, indicating the resulting optimal number of clusters
- `seq_vbic`: A vector, corresponding to the sequence of the V-BIC values associated with the models trained for each provided cluster's number in `grid_nb_cluster`.
- `trained_models`: A list, named by associated number of clusters, of Magma or MagmaClust models that have been trained (or approximated if `fast_approx = T`) during the model selection procedure.

### Examples

```
TRUE
```

---

```
simu_db
```

```
Simulate a dataset tailored for MagmaClustR
```

---

### Description

Simulate a complete training dataset, which may be representative of various applications. Several flexible arguments allow adjustment of the number of individuals, of observed inputs, and the values of many parameters controlling the data generation.

**Usage**

```

simu_db(
  M = 10,
  N = 10,
  K = 1,
  covariate = FALSE,
  grid = seq(0, 10, 0.05),
  common_input = TRUE,
  common_hp = TRUE,
  add_hp = FALSE,
  add_clust = FALSE,
  int_mu_v = c(4, 5),
  int_mu_l = c(0, 1),
  int_i_v = c(1, 2),
  int_i_l = c(0, 1),
  int_i_sigma = c(0, 0.2),
  m0_slope = c(-5, 5),
  m0_intercept = c(-50, 50),
  int_covariate = c(-5, 5)
)

```

**Arguments**

M	An integer. The number of individual per cluster.
N	An integer. The number of observations per individual.
K	An integer. The number of underlying clusters.
covariate	A logical value indicating whether the dataset should include an additional input covariate named 'Covariate'.
grid	A vector of numbers defining a grid of observations (i.e. the reference inputs).
common_input	A logical value indicating whether the reference inputs are common to all individual.
common_hp	A logical value indicating whether the hyper-parameters are common to all individual. If TRUE and $K > 1$ , the hyper-parameters remain different between the clusters.
add_hp	A logical value indicating whether the values of hyper-parameters should be added as columns in the dataset.
add_clust	A logical value indicating whether the name of the clusters should be added as a column in the dataset.
int_mu_v	A vector of 2 numbers, defining an interval of admissible values for the variance hyper-parameter of the mean process' kernel.
int_mu_l	A vector of 2 numbers, defining an interval of admissible values for the length-scale hyper-parameter of the mean process' kernel.
int_i_v	A vector of 2 numbers, defining an interval of admissible values for the variance hyper-parameter of the individual process' kernel.

int_i_l	A vector of 2 numbers, defining an interval of admissible values for the length-scale hyper-parameter of the individual process' kernel.
int_i_sigma	A vector of 2 numbers, defining an interval of admissible values for the noise hyper-parameter.
m0_slope	A vector of 2 numbers, defining an interval of admissible values for the slope of m0.
m0_intercept	A vector of 2 numbers, defining an interval of admissible values for the intercept of m0.
int_covariate	A vector of 2 numbers, defining an interval of admissible values for the covariate inputs.

**Value**

A full dataset of simulated training data.

**Examples**

```
## Generate a dataset with 3 clusters of 4 individuals, observed at 10 inputs
data = simu_db(M = 4, N = 10, K = 3)

## Generate a 2-D dataset with an additional input 'Covariate'
data = simu_db(covariate = TRUE)

## Generate a dataset where input locations are different among individuals
data = simu_db(common_input = FALSE)

## Generate a dataset with an additional column indicating the true clusters
data = simu_db(K = 3, add_clust = TRUE)
```

---

train\_gp

*Learning hyper-parameters of a Gaussian Process*


---

**Description**

Learning hyper-parameters of any new individual/task in Magma is required in the prediction procedure. This function can also be used to learn hyper-parameters of a simple GP (just let the hyperpost argument set to NULL, and use prior\_mean instead). When using within Magma, by providing data for the new individual/task, the hyper-posterior mean and covariance parameters, and initialisation values for the hyper-parameters, the function computes maximum likelihood estimates of the hyper-parameters.

**Usage**

```
train_gp(
  data,
  prior_mean = NULL,
  ini_hp = NULL,
```

```

    kern = "SE",
    hyperpost = NULL,
    pen_diag = 1e-10
  )

```

## Arguments

data	A tibble or data frame. Required columns: Input, Output. Additional columns for covariates can be specified. The Input column should define the variable that is used as reference for the observations (e.g. time for longitudinal data). The Output column specifies the observed values (the response variable). The data frame can also provide as many covariates as desired, with no constraints on the column names. These covariates are additional inputs (explanatory variables) of the models that are also observed at each reference Input.
prior_mean	Mean parameter of the GP. This argument can be specified under various formats, such as: <ul style="list-style-type: none"> <li>• NULL (default). The hyper-posterior mean would be set to 0 everywhere.</li> <li>• A number. The hyper-posterior mean would be a constant function.</li> <li>• A vector of the same length as all the distinct Input values in the data argument. This vector would be considered as the evaluation of the hyper-posterior mean function at the training Inputs.</li> <li>• A function. This function is defined as the hyper-posterior mean.</li> <li>• A tibble or data frame. Required columns: Input, Output. The Input values should include at least the same values as in the data argument.</li> </ul>
ini_hp	A named vector, tibble or data frame of hyper-parameters associated with the kern of the new individual/task. The columns should be named according to the hyper-parameters that are used in kern. In cases where the model includes a noise term, ini_hp should contain an additional 'noise' column. If NULL (default), random values are used as initialisation.
kern	A kernel function, defining the covariance structure of the GP. Several popular kernels (see <a href="#">The Kernel Cookbook</a> ) are already implemented and can be selected within the following list: <ul style="list-style-type: none"> <li>• "SE": (default value) the Squared Exponential Kernel (also called Radial Basis Function or Gaussian kernel),</li> <li>• "LIN": the Linear kernel,</li> <li>• "PERIO": the Periodic kernel,</li> <li>• "RQ": the Rational Quadratic kernel. Compound kernels can be created as sums or products of the above kernels. For combining kernels, simply provide a formula as a character string where elements are separated by whitespaces (e.g. "SE + PERIO"). As the<sup>2</sup> elements are treated sequentially from the left to the right, the product operator '*' shall always be used before the '+' operators (e.g. 'SE * LIN + RQ' is valid whereas 'RQ + SE * LIN' is not).</li> </ul>
hyperpost	A list, containing the elements 'mean' and 'cov', the parameters of the hyper-posterior distribution of the mean process. Typically, this argument should come from a previous learning using <a href="#">train_magma</a> , or from the <a href="#">hyperposterior</a>

function. If hyperpost is provided, the likelihood that is maximised is the one involved during Magma's prediction step, and the prior\_mean argument is ignored. For classic GP training, leave hyperpost to NULL.

pen\_diag A number. A jitter term, added on the diagonal to prevent numerical issues when inverting nearly singular matrices.

### Value

A tibble, containing the trained hyper-parameters for the kernel of the new individual/task.

### Examples

```
TRUE
```

---

train_gp_clust	<i>Prediction in MagmaClust: learning new HPs and mixture probabilities</i>
----------------	---

---

### Description

Learning hyper-parameters and mixture probabilities of any new individual/task is required in MagmaClust in the prediction procedure. By providing data for the new individual/task, the hyper-posterior mean and covariance parameters, the mixture proportions, and initialisation values for the hyper-parameters, train\_gp\_clust uses an EM algorithm to compute maximum likelihood estimates of the hyper-parameters and hyper-posterior mixture probabilities of the new individual/task.

### Usage

```
train_gp_clust(
  data,
  prop_mixture = NULL,
  ini_hp = NULL,
  kern = "SE",
  hyperpost = NULL,
  pen_diag = 1e-10,
  n_iter_max = 25,
  cv_threshold = 0.001
)
```

### Arguments

data A tibble or data frame. Required columns: Input, Output. Additional columns for covariates can be specified. The Input column should define the variable that is used as reference for the observations (e.g. time for longitudinal data). The Output column specifies the observed values (the response variable). The data frame can also provide as many covariates as desired, with no constraints on the column names. These covariates are additional inputs (explanatory variables) of the models that are also observed at each reference Input.

prop_mixture	A tibble or a named vector. Each name of column or element should refer to a cluster. The value associated with each cluster is a number between 0 and 1, corresponding to the mixture proportions.
ini_hp	A tibble or data frame of hyper-parameters associated with kern, the individual process kernel.
kern	A kernel function, defining the covariance structure of the GP. Several popular kernels (see <a href="#">The Kernel Cookbook</a> ) are already implemented and can be selected within the following list: <ul style="list-style-type: none"> <li>• "SE": (default value) the Squared Exponential Kernel (also called Radial Basis Function or Gaussian kernel),</li> <li>• "LIN": the Linear kernel,</li> <li>• "PERIO": the Periodic kernel,</li> <li>• "RQ": the Rational Quadratic kernel. Compound kernels can be created as sums or products of the above kernels. For combining kernels, simply provide a formula as a character string where elements are separated by whitespaces (e.g. "SE + PERIO"). As the<sup>2</sup> elements are treated sequentially from the left to the right, the product operator '*' shall always be used before the '+' operators (e.g. 'SE * LIN + RQ' is valid whereas 'RQ + SE * LIN' is not).</li> </ul>
hyperpost	A list, containing the elements mean, cov and mixture the parameters of the hyper-posterior distributions of the mean processes. Typically, this argument should come from a previous learning using <a href="#">train_magmaclust</a> , or a previous prediction with <a href="#">pred_magmaclust</a> , with the argument get_hyperpost set to TRUE.
pen_diag	A number. A jitter term, added on the diagonal to prevent numerical issues when inverting nearly singular matrices.
n_iter_max	A number, indicating the maximum number of iterations of the EM algorithm to proceed while not reaching convergence.
cv_threshold	A number, indicating the threshold of the likelihood gain under which the EM algorithm will stop.

### Value

A list, containing the results of the EM algorithm used during the prediction step of MagmaClust. The elements of the list are:

- hp: A tibble of optimal hyper-parameters for the new individual's GP.
- mixture: A tibble of mixture probabilities for the new individual.

### Examples

```
TRUE
```

train\_magma

*Training Magma with an EM algorithm***Description**

The hyper-parameters and the hyper-posterior distribution involved in Magma can be learned thanks to an EM algorithm implemented in `train_magma`. By providing a dataset, the model hypotheses (hyper-prior mean parameter and covariance kernels) and initialisation values for the hyper-parameters, the function computes maximum likelihood estimates of the HPs as well as the mean and covariance parameters of the Gaussian hyper-posterior distribution of the mean process.

**Usage**

```
train_magma(
  data,
  prior_mean = NULL,
  ini_hp_0 = NULL,
  ini_hp_i = NULL,
  kern_0 = "SE",
  kern_i = "SE",
  common_hp = TRUE,
  grid_inputs = NULL,
  pen_diag = 1e-10,
  n_iter_max = 25,
  cv_threshold = 0.001,
  fast_approx = FALSE
)
```

**Arguments**

data	A tibble or data frame. Required columns: ID, Input , Output. Additional columns for covariates can be specified. The ID column contains the unique names/codes used to identify each individual/task (or batch of data). The Input column should define the variable that is used as reference for the observations (e.g. time for longitudinal data). The Output column specifies the observed values (the response variable). The data frame can also provide as many covariates as desired, with no constraints on the column names. These covariates are additional inputs (explanatory variables) of the models that are also observed at each reference Input.
prior_mean	Hyper-prior mean parameter ( $m_0$ ) of the mean GP. This argument can be specified under various formats, such as: <ul style="list-style-type: none"> <li>• NULL (default). The hyper-prior mean would be set to 0 everywhere.</li> <li>• A number. The hyper-prior mean would be a constant function.</li> <li>• A vector of the same length as all the distinct Input values in the data argument. This vector would be considered as the evaluation of the hyper-prior mean function at the training Inputs.</li> </ul>



	<ul style="list-style-type: none"> <li>• A function. This function is defined as the hyper_prior mean.</li> <li>• A tibble or data frame. Required columns: Input, Output. The Input values should include at least the same values as in the data argument.</li> </ul>
ini_hp_0	A named vector, tibble or data frame of hyper-parameters associated with kern_0, the mean process' kernel. The columns/elements should be named according to the hyper-parameters that are used in kern_0. If NULL (default), random values are used as initialisation.
ini_hp_i	A tibble or data frame of hyper-parameters associated with kern_i, the individual processes' kernel. Required column : ID. The ID column contains the unique names/codes used to identify each individual/task. The other columns should be named according to the hyper-parameters that are used in kern_i. Compared to ini_hp_0 should contain an additional 'noise' column to initialise the noise hyper-parameter of the model. If NULL (default), random values are used as initialisation.
kern_0	A kernel function, associated with the mean GP. Several popular kernels (see <a href="#">The Kernel Cookbook</a> ) are already implemented and can be selected within the following list: <ul style="list-style-type: none"> <li>• "SE": (default value) the Squared Exponential Kernel (also called Radial Basis Function or Gaussian kernel),</li> <li>• "LIN": the Linear kernel,</li> <li>• "PERIO": the Periodic kernel,</li> <li>• "RQ": the Rational Quadratic kernel. Compound kernels can be created as sums or products of the above kernels. For combining kernels, simply provide a formula as a character string where elements are separated by whitespaces (e.g. "SE + PERIO"). As the elements are treated sequentially from the left to the right, the product operator '*' shall always be used before the '+' operators (e.g. 'SE * LIN + RQ' is valid whereas 'RQ + SE * LIN' is not).</li> </ul>
kern_i	A kernel function, associated with the individual GPs. ("SE", "PERIO" and "RQ" are also available here).
common_hp	A logical value, indicating whether the set of hyper-parameters is assumed to be common to all individuals.
grid_inputs	A vector, indicating the grid of additional reference inputs on which the mean process' hyper-posterior should be evaluated.
pen_diag	A number. A jitter term, added on the diagonal to prevent numerical issues when inverting nearly singular matrices.
n_iter_max	A number, indicating the maximum number of iterations of the EM algorithm to proceed while not reaching convergence.
cv_threshold	A number, indicating the threshold of the likelihood gain under which the EM algorithm will stop. The convergence condition is defined as the difference of likelihoods between two consecutive steps, divided by the absolute value of the last one ( $(LL_n - LL_{n-1})/ LL_n $ ).
fast_approx	A boolean, indicating whether the EM algorithm should stop after only one iteration of the E-step. This advanced feature is mainly used to provide a faster approximation of the model selection procedure, by preventing any optimisation over the hyper-parameters.

## Details

The user can specify custom kernel functions for the argument `kern_0` and `kern_i`. The hyper-parameters used in the kernel should have explicit names, and be contained within the `hp` argument. `hp` should typically be defined as a named vector or a data frame. Although it is not mandatory for the `train_magma` function to run, gradients can be provided within kernel function definition. See for example [se\\_kernel](#) to create a custom kernel function displaying an adequate format to be used in Magma.

## Value

A list, gathering the results of the EM algorithm used for training in Magma. The elements of the list are:

- `hp_0`: A tibble of the trained hyper-parameters for the mean process' kernel.
- `hp_i`: A tibble of all the trained hyper-parameters for the individual processes' kernels.
- `hyperpost`: A sub-list gathering the parameters of the mean processes' hyper-posterior distributions, namely:
  - `mean`: A tibble, the hyper-posterior mean parameter (Output) evaluated at each training reference Input.
  - `cov`: A matrix, the covariance parameter for the hyper-posterior distribution of the mean process.
  - `pred`: A tibble, the predicted mean and variance at Input for the mean process' hyper-posterior distribution under a format that allows the direct visualisation as a GP prediction.
- `ini_args`: A list containing the initial function arguments and values for the hyper-prior mean, the hyper-parameters. In particular, if those arguments were set to NULL, `ini_args` allows us to retrieve the (randomly chosen) initialisations used during training.
- `seq_loglikelihood`: A vector, containing the sequence of log-likelihood values associated with each iteration.
- `converged`: A logical value indicated whether the EM algorithm converged or not.
- `training_time`: Total running time of the complete training.

## Examples

```
TRUE
```

---

```
train_magmaclust
```

```
Training MagmaClust with a Variational EM algorithm
```

---

## Description

The hyper-parameters and the hyper-posterior distributions involved in MagmaClust can be learned thanks to a VEM algorithm implemented in `train_magmaclust`. By providing a dataset, the model hypotheses (hyper-prior mean parameters, covariance kernels and number of clusters) and initialisation values for the hyper-parameters, the function computes maximum likelihood estimates of the HPs as well as the mean and covariance parameters of the Gaussian hyper-posterior distributions of the mean processes.

**Usage**

```
train_magmaclust(
  data,
  nb_cluster = NULL,
  prior_mean_k = NULL,
  ini_hp_k = NULL,
  ini_hp_i = NULL,
  kern_k = "SE",
  kern_i = "SE",
  ini_mixture = NULL,
  common_hp_k = TRUE,
  common_hp_i = TRUE,
  grid_inputs = NULL,
  pen_diag = 1e-10,
  n_iter_max = 25,
  cv_threshold = 0.001,
  fast_approx = FALSE
)
```

**Arguments**

data	A tibble or data frame. Columns required: ID, Input , Output. Additional columns for covariates can be specified. The ID column contains the unique names/codes used to identify each individual/task (or batch of data). The Input column should define the variable that is used as reference for the observations (e.g. time for longitudinal data). The Output column specifies the observed values (the response variable). The data frame can also provide as many covariates as desired, with no constraints on the column names. These covariates are additional inputs (explanatory variables) of the models that are also observed at each reference Input.
nb_cluster	A number, indicating the number of clusters of individuals/tasks that are assumed to exist among the dataset.
prior_mean_k	The set of hyper-prior mean parameters ( $m_k$ ) for the K mean GPs, one value for each cluster. This argument can be specified under various formats, such as: <ul style="list-style-type: none"> <li>• NULL (default). All hyper-prior means would be set to 0 everywhere.</li> <li>• A numerical vector of the same length as the number of clusters. Each number is associated with one cluster, and considered to be the hyper-prior mean parameter of the cluster (i.e. a constant function at all Input).</li> <li>• A list of functions. Each function is associated with one cluster. These functions are all evaluated at all Input values, to provide specific hyper-prior mean vectors for each cluster.</li> </ul>
ini_hp_k	A tibble or data frame of hyper-parameters associated with kern_k, the mean process' kernel. Required column : ID. The ID column contains the unique names/codes used to identify each cluster. The other columns should be named according to the hyper-parameters that are used in kern_k.

ini_hp_i	A tibble or data frame of hyper-parameters associated with kern_i, the individual processes' kernel. Required column : ID. The ID column contains the unique names/codes used to identify each individual/task. The other columns should be named according to the hyper-parameters that are used in kern_i.
kern_k	A kernel function, associated with the mean GPs. Several popular kernels (see <a href="#">The Kernel Cookbook</a> ) are already implemented and can be selected within the following list: <ul style="list-style-type: none"> <li>• "SE": (default value) the Squared Exponential Kernel (also called Radial Basis Function or Gaussian kernel),</li> <li>• "LIN": the Linear kernel,</li> <li>• "PERIO": the Periodic kernel,</li> <li>• "RQ": the Rational Quadratic kernel. Compound kernels can be created as sums or products of the above kernels. For combining kernels, simply provide a formula as a character string where elements are separated by whitespaces (e.g. "SE + PERIO"). As the elements are treated sequentially from the left to the right, the product operator '*' shall always be used before the '+' operators (e.g. 'SE * LIN + RQ' is valid whereas 'RQ + SE * LIN' is not).</li> </ul>
kern_i	A kernel function, associated with the individual GPs. (See details above in kern_k).
ini_mixture	Initial values of the probability to belong to each cluster for each individual ( <a href="#">ini_mixture</a> can be used for a k-means initialisation. Used by default if NULL).
common_hp_k	A boolean indicating whether hyper-parameters are common among the mean GPs.
common_hp_i	A boolean indicating whether hyper-parameters are common among the individual GPs.
grid_inputs	A vector, indicating the grid of additional reference inputs on which the mean processes' hyper-posteriors should be evaluated.
pen_diag	A number. A jitter term, added on the diagonal to prevent numerical issues when inverting nearly singular matrices.
n_iter_max	A number, indicating the maximum number of iterations of the VEM algorithm to proceed while not reaching convergence.
cv_threshold	A number, indicating the threshold of the likelihood gain under which the VEM algorithm will stop. The convergence condition is defined as the difference of elbo between two consecutive steps, divided by the absolute value of the last one ( $(ELBO_n - ELBO_{n-1})/ ELBO_n $ ).
fast_approx	A boolean, indicating whether the VEM algorithm should stop after only one iteration of the VE-step. This advanced feature is mainly used to provide a faster approximation of the model selection procedure, by preventing any optimisation over the hyper-parameters.

## Details

The user can specify custom kernel functions for the argument kern\_k and kern\_i. The hyper-parameters used in the kernel should have explicit names, and be contained within the hp argument.

hp should typically be defined as a named vector or a data frame. Although it is not mandatory for the `train_magmaclust` function to run, gradients can be provided within kernel function definition. See for example `se_kernel` to create a custom kernel function displaying an adequate format to be used in MagmaClust.

### Value

A list, containing the results of the VEM algorithm used in the training step of MagmaClust. The elements of the list are:

- `hp_k`: A tibble containing the trained hyper-parameters for the mean process' kernel and the mixture proportions for each cluster.
- `hp_i`: A tibble containing the trained hyper-parameters for the individual processes' kernels.
- `hyperpost`: A sub-list containing the parameters of the mean processes' hyper-posterior distribution, namely:
  - `mean`: A list of tibbles containing, for each cluster, the hyper-posterior mean parameters evaluated at each Input.
  - `cov`: A list of matrices containing, for each cluster, the hyper-posterior covariance parameter of the mean process.
  - `mixture`: A tibble, indicating the mixture probabilities in each cluster for each individual.
- `ini_args`: A list containing the initial function arguments and values for the hyper-prior means, the hyper-parameters. In particular, if those arguments were set to `NULL`, `ini_args` allows us to retrieve the (randomly chosen) initialisations used during training.
- `seq_elbo`: A vector, containing the sequence of ELBO values associated with each iteration.
- `converged`: A logical value indicated whether the algorithm converged.
- `training_time`: Total running time of the complete training.

### Examples

TRUE

# Index

`data_allocate_cluster`, [2](#), [14](#)

`hp`, [3](#)  
`hyperposterior`, [4](#), [29](#)  
`hyperposterior_clust`, [6](#)

`ini_mixture`, [36](#)

`MagmaClustR`, [8](#)

`plot_db`, [9](#)  
`plot_gif`, [10](#), [17](#)  
`plot_gp`, [11](#)  
`plot_magmaclust`, [13](#)  
`pred_gif`, [10](#), [15](#)  
`pred_gp`, [12](#), [17](#), [24](#)  
`pred_magma`, [4](#), [12](#), [16](#), [19](#), [20](#), [24](#)  
`pred_magmaclust`, [6](#), [14](#), [21](#), [22](#), [31](#)  
`proba_max_cluster`, [23](#)

`sample_gp`, [24](#)  
`se_kernel`, [34](#), [37](#)  
`select_nb_cluster`, [25](#)  
`simu_db`, [26](#)

`train_gp`, [16](#), [18](#), [19](#), [28](#)  
`train_gp_clust`, [22](#), [30](#)  
`train_magma`, [16](#), [19](#), [20](#), [29](#), [32](#)  
`train_magmaclust`, [3](#), [14](#), [21](#), [22](#), [26](#), [31](#), [34](#)  
`transition_states`, [11](#)