

---

# **MeGaMix Documentation**

***Release 0.2***

**Elina Thibeaudeau-Sutre**

**Sep 18, 2017**



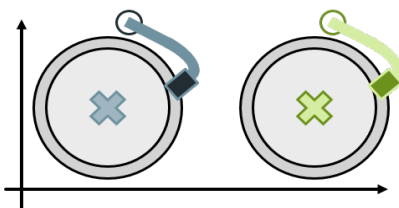
---

## Contents

---

|          |                                                              |           |
|----------|--------------------------------------------------------------|-----------|
| <b>1</b> | <b>Getting started</b>                                       | <b>3</b>  |
| 1.1      | Installation . . . . .                                       | 3         |
| 1.2      | Description . . . . .                                        | 3         |
| 1.3      | Basic usage . . . . .                                        | 5         |
| <b>2</b> | <b>API Reference</b>                                         | <b>9</b>  |
| 2.1      | Batch versions of the algorithm . . . . .                    | 9         |
| 2.2      | Online versions of the algorithm . . . . .                   | 19        |
| <b>3</b> | <b>Theory of Gaussian Mixture models</b>                     | <b>25</b> |
| 3.1      | K-means . . . . .                                            | 25        |
| 3.2      | Gaussian Mixture Model (GMM) . . . . .                       | 26        |
| 3.3      | Variational Gaussian Mixture Model (VBGMM) . . . . .         | 27        |
| 3.4      | Dirichlet Process Gaussian Mixture Model (DPGMM) . . . . .   | 28        |
| 3.5      | Pitman-Yor Process Gaussian Mixture Model (PYPGMM) . . . . . | 28        |





## Table of contents



### Installation

The package is registered on PyPI. It can be installed with the following command:

```
$ pip install megamix
```

If you want to install it manually, you can find the source code at <https://github.com/14thibea/megamix>.

MeGaMix relies on external dependencies. The setup script should install them automatically, but you may want to install them manually. The required packages are:

- NumPy 1.11.3 or newer
- scipy 0.18.1 or newer
- h5py 2.6.0 or newer
- joblib 0.11 or newer
- cython

---

**Note:** If you cannot compile the package, please download it and manually comment the line `ext_modules=cythonize(ext_modules)`, in `setup.py`. Doing so you will not compile the cython modules and only use pure python versions.

---

### Description

The MeGaMix package (Methods for Gaussian Mixtures) allows Python developers to fit different kind of models on their data. The different models are clustering methods of unsupervised machine learning. Four models have been implemented, from the most simple to the most complex:

- K-means

- GMM (Gaussian Mixture Model)
- VBGMM (Variational Bayesian Gaussian Mixture Model)
- DP-VBGMM (Dirichlet Process on Variational Bayesian Gaussian Mixture Model)
- PYP-VBGMM (Pitman-Yor Process on Variational Bayesian Gaussian Mixture Model)

## What will you be able to do ?

The main idea of clustering algorithms is to create groups by gathering points that are close to each other.

A cluster has three main parameters:

- A mean : the mean of all the points that belong to the cluster
- A weight : the number of points that belong to the cluster
- A covariance (except for K-means) : a matrix which specifies the form of the cluster

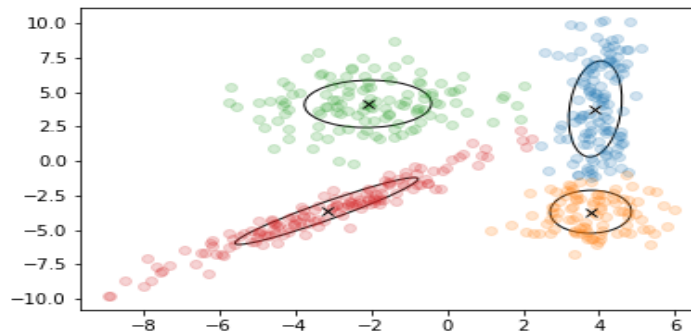


Fig. 1.1: A graphical example of a gaussian mixture model fit on a set of points

## How do the algorithms work ?

After the initialisation, the algorithms alternate between two steps, the E step (Expectation) and the M step (Maximisation).

During the **E step**, the algorithm computes the probability for each point to belong to each cluster. It produces an array of *responsibilities*. At the  $i$ th row and the  $j$ th column of this array corresponds the probability of the  $i$ th point to belong to the  $j$ th cluster.

Here is an example of responsibilities that could be obtained with 6 points and 2 clusters :

|         | Cluster 1 | Cluster 2 |
|---------|-----------|-----------|
| point 1 | 0.54      | 0.46      |
| point 2 | 0.89      | 0.11      |
| point 3 | 0.27      | 0.73      |
| point 4 | 0.01      | 0.99      |
| point 5 | 0.42      | 0.58      |
| point 6 | 0.84      | 0.16      |

*In this example, the first point has a 54% chance to belong to the first cluster and a 46% chance to belong to the second cluster.*



---

**Note:** This is not the case with K-means which is not working with probabilities but with labels. A point belongs completely to a cluster or doesn't belong to it (this is called hard assignement).

---

Then during the **M step**, the algorithm re-estimates the parameters of the model in order to maximize a convergence criterion.

Finally the algorithm will stop if the difference between the value of the convergence criterion of the current and the previous is less than a threshold fixed by the user.

This is summarized in the following pseudo-code:

```
initialize(points)
while(cc-cc_previous > tol):
    cc_previous = cc
    responsibilities = E_step(points,parameters)
    parameters = M_step(responsabilities,points)
    cc = convergence_criterion(points,responsabilities,parameters)
```

## What is it used for ?

MeGaMix has been implemented in order to process natural speech **MFCC**. Unlike the vision field where deep learning has overtaken such clustering models, they are still efficient in speech processing.

However the use of this package is more general and may serve another purpose.

## Basic usage

```
#####
# Prelude to the example
#####
"""
This example is realized with a DP-VBGMM model
The other mixtures and the K-means are working in the same way
The available classes are:
    - Kmeans (kmeans)
    - GaussianMixture (GMM)
    - VariationalGaussianMixture (VBGMM)
    - DPVariationalGaussianMixture (DP-VBGMM)
"""

from megamix import DPVariationalGaussianMixture
import numpy as np

#####
# Features used
#####

"""
Features must be numpy arrays of two dimensions:
the first dimension is the number of points
the second dimension is the dimension of the space
"""
```

```
# Here we use a random set of points for the example
n_points = 10000
dim = 39

points = np.random.randn(n_points,dim)

#####
# Fitting the model
#####

# We choose the number of clusters that we want
n_components = 100

# The model is instantiated
GM = DPVariationalGaussianMixture(n_components)

# The model is fitting
GM.fit(points)

# It is also possible to do early stopping in order to avoid overfitting
points_data = points[:n_points//2:]
points_test = points[n_points//2::]

# In this case the model will fit only on points_data but will use points_test
# to evaluate the convergence criterion.
GM.fit(points_data,points_test)

# Some clusters may disappear with the DP-VBGMM model. You may want to
# simplify the model by removing the useless information
GM_simple = GM.simplified_model(points)

#####
# Analysis of the model
#####

other_points = np.random.randn(n_points,dim)

# We can obtain the log of the responsibilities of any set of points when the
# model is fitted (or at least initialized)
log_resp = GM.predict_log_resp(other_points)
# log_resp.shape = (n_points,n_components)

# We can obtain the value of the convergence criterion for any set of points
score = GM.score(other_points)

#####
# Writing or reading a model
#####

# It is possible to write your model in a group of a h5py file
import h5py

file = h5py.File('DP_VBGMM.h5','w')
grp = file.create_group('model_fitted')

GM.write(grp)
file.close()
```

```
# You also can read data from such h5py file to initialize new models
GM_new = DPVariationalGaussianMixture()

file = h5py.File('DP_VBGMM.h5', 'r')
grp = file['model_fitted']

GM_new.read_and_init(grp, points)
file.close()

# You can also save regularly your code while fitting the model by using
# the saving parameter

GM.fit(points, saving='log', directory='mypath', legend='wonderful_model')
```



### Batch versions of the algorithm

#### Kmeans

**class** megamix.batch.**Kmeans** (*n\_components=1, init='plus', n\_jobs=1*)  
Kmeans model.

##### Parameters

- **n\_components** (*int, defaults to 1.*) – Number of clusters used.
- **init** (*str, defaults to 'kmeans'.*) – Method used in order to perform the initialization, must be in ['random', 'plus', 'AF\_KMC'].

##### name

*str* – The name of the method : 'Kmeans'

##### means

*array of floats (n\_components,dim)* – Contains the computed means of the model.

##### log\_weights

*array of floats (n\_components,)* – Contains the logarithm of the mixing coefficient of each cluster.

##### iter

*int* – The number of iterations computed with the method fit()

##### \_is\_initialized

*bool* – Ensures that the model has been initialized before using other methods such as distortion() or predict\_assignements().

**Raises** ValueError : if the parameters are inconsistent, for example if the cluster number is negative, init\_type is not in ['resp', 'mcw']...

## References

‘Fast and Provably Good Seedings for k-Means’, O. Bachem, M. Lucic, S. Hassani, A. Krause ‘Lloyd’s algorithm <[https://en.wikipedia.org/wiki/Lloyd's\\_algorithm](https://en.wikipedia.org/wiki/Lloyd's_algorithm)>’ ‘The remarkable k-means++ <<https://normaldeviate.wordpress.com/2012/09/30/the-remarkable-k-means/>>’

**fit** (*points\_data*, *points\_test=None*, *n\_iter\_max=100*, *n\_iter\_fix=None*, *tol=0*, *saving=None*, *file\_name='model'*, *saving\_iter=2*)  
The k-means algorithm

### Parameters

- **points\_data** (*array (n\_points, dim)*) – A 2D array of points on which the model will be trained
- **tol** (*float, defaults to 0*) – The EM algorithm will stop when the difference between two steps regarding the distortion is less or equal to tol.
- **n\_iter\_max** (*int, defaults to 100*) – number of iterations maximum that can be done
- **saving\_iter** (*int | defaults 2*) – An int to know how often the model is saved (see saving below).
- **file\_name** (*str | defaults model*) – The name of the file (including the path).

### Other Parameters

- **points\_test** (*array (n\_points\_bis, dim) | Optional*) – A 2D array of points on which the model will be tested.
- **n\_iter\_fix** (*int | Optional*) – If not None, the algorithm will exactly do the number of iterations of n\_iter\_fix and stop.
- **saving** (*str | Optional*) – A string in ['log', 'linear']. In the following equations x is the parameter saving\_iter (see above).
  - If 'log', the model will be saved for all iterations which verify :  $\log(\text{iter})/\log(x)$  is an int
  - If 'linear' the model will be saved for all iterations which verify :  $\text{iter}/x$  is an int

### Returns

Return type None

**predict\_assignments** (*points*)

This function return the hard assignments of points once the model is fitted.

**score** (*points, assignments=None*)

This method returns the distortion measurement at the end of the k\_means.

### Parameters

- **points** (*an array (n\_points, dim)*) –
- **assignments** (*an array (n\_components, dim)*) – an array containing the responsibilities of the clusters

### Returns distortion

Return type (float)

## Gaussian Mixture Model (GMM)

```
class megamix.batch.GaussianMixture(n_components=1, covariance_type='full', init='kmeans',
                                    reg_covar=1e-06, type_init='resp', n_jobs=1)
```

Gaussian Mixture Model

Representation of a Gaussian mixture model probability distribution. This class allows to estimate the parameters of a Gaussian mixture distribution.

### Parameters

- **n\_components** (*int*, defaults to 1.) – Number of clusters used.
- **init** (*str*, defaults to 'kmeans'.) – Method used in order to perform the initialization, must be in ['random', 'plus', 'AF\_KMC', 'kmeans'].
- **reg\_covar** (*float*, defaults to 1e-6) – In order to avoid null covariances this float is added to the diagonal of covariance matrices.
- **type\_init** (*str*, defaults to 'resp'.) – The algorithm is initialized using this data (responsibilities if 'resp' or means, covariances and weights if 'mcw').

### name

*str* – The name of the method : 'GMM'

### cov

*array of floats (n\_components,dim,dim)* – Contains the computed covariance matrices of the mixture.

### means

*array of floats (n\_components,dim)* – Contains the computed means of the mixture.

### log\_weights

*array of floats (n\_components,)* – Contains the logarithm of the mixing coefficient of each cluster.

### iter

*int* – The number of iterations computed with the method fit()

### convergence\_criterion\_data

*array of floats (iter,)* – Stores the value of the convergence criterion computed with data on which the model is fitted.

### convergence\_criterion\_test

*array of floats (iter,) | if \_early\_stopping only* – Stores the value of the convergence criterion computed with test data if it exists.

### \_is\_initialized

*bool* – Ensures that the method \_initialize() has been used before using other methods such as score() or predict\_log\_assignements().

**Raises** ValueError : if the parameters are inconsistent, for example if the cluster number is negative, init\_type is not in ['resp','mcw']...

## References

'Pattern Recognition and Machine Learning', Bishop

**fit** (*points\_data*, *points\_test=None*, *tol=0.001*, *patience=None*, *n\_iter\_max=100*, *n\_iter\_fix=None*, *saving=None*, *file\_name='model'*, *saving\_iter=2*)  
The EM algorithm

### Parameters

- **points\_data** (*array (n\_points, dim)*) – A 2D array of points on which the model will be trained
- **tol** (*float, defaults to 1e-3*) – The EM algorithm will stop when the difference between two steps regarding the convergence criterion is less than tol.
- **n\_iter\_max** (*int, defaults to 100*) – number of iterations maximum that can be done
- **saving\_iter** (*int | defaults 2*) – An int to know how often the model is saved (see saving below).
- **file\_name** (*str | defaults model*) – The name of the file (including the path).

#### Other Parameters

- **points\_test** (*array (n\_points\_bis, dim) | Optional*) – A 2D array of points on which the model will be tested.
- **patience** (*int | Optional*) – The number of iterations performed after having satisfied the convergence criterion
- **n\_iter\_fix** (*int | Optional*) – If not None, the algorithm will exactly do the number of iterations of n\_iter\_fix and stop.
- **saving** (*str | Optional*) – A string in ['log', 'linear']. In the following equations x is the parameter saving\_iter (see above).
  - If 'log', the model will be saved for all iterations which verify :  
 $\log(\text{iter})/\log(x)$  is an int
  - If 'linear' the model will be saved for all iterations which verify :  $\text{iter}/x$  is an int

#### Returns

**Return type** None

**predict\_log\_resp** (*points*)

This function returns the logarithm of each point's responsibilities

**Parameters** **points** (*array (n\_points\_bis, dim)*) – a 1D or 2D array of points with the same dimension as the problem

**Returns** **log\_resp** – the logarithm of the responsibilities

**Return type** array (n\_points\_bis, n\_components)

**read\_and\_init** (*group, points*)

A method reading a group of an hdf5 file to initialize DPGMM

**Parameters** **group** (*HDF5 group*) – A group of a hdf5 file in reading mode

**score** (*points*)

This function return the score of the function, which is the logarithm of the likelihood for GMM and the logarithm of the lower bound of the likelihood for VBGMM and DPGMM

**Parameters** **points** (*array (n\_points\_bis, dim)*) – a 1D or 2D array of points with the same dimension as the problem

**Returns** **score**

**Return type** float



**simplified\_model** (*points*)

A method creating a new model with simplified parameters: clusters unused are removed

**Parameters** *points* (*an array (n\_points, dim)*) –

**Returns** **GM**

**Return type** an instance of the same type of self: GMM,VBGMM or DPGMM

**write** (*group*)

A method creating datasets in a group of an hdf5 file in order to save the model

**Parameters** *group* (*HDF5 group*) – A group of a hdf5 file in reading mode

## Variational Gaussian Mixture Model (VBGMM)

```
class megamix.batch.VariationalGaussianMixture (n_components=1,          init='kmeans',
                                                  alpha_0=None,          beta_0=None,
                                                  nu_0=None,          means_prior=None,
                                                  cov_wishart_prior=None, reg_covar=1e-06,
                                                  type_init='resp',      n_jobs=1,
                                                  boost=None)
```

Variational Bayesian Estimation of a Gaussian Mixture

This class allows to infer an approximate posterior distribution over the parameters of a Gaussian mixture distribution.

The weights distribution is a Dirichlet distribution with parameter alpha (see Bishop's book p474-486)

### Parameters

- **n\_components** (*int, defaults to 1.*) – Number of clusters used.
- **init** (*str, defaults to 'kmeans'.*) – Method used in order to perform the initialization, must be in ['random', 'plus', 'AF\_KMC', 'kmeans', 'GMM'].
- **reg\_covar** (*float, defaults to 1e-6*) – In order to avoid null covariances this float is added to the diagonal of covariance matrices.
- **type\_init** (*str, defaults to 'resp'.*) – The algorithm is initialized using this data (responsibilities if 'resp' or means, covariances and weights if 'mcw').

### Other Parameters

- **alpha\_0** (*float, Optional | defaults to None.*) – The prior parameter on the weight distribution (Dirichlet). A high value of alpha\_0 will lead to equal weights, while a low value will allow some clusters to shrink and disappear. Must be greater than 0.

If None, the value is set to 1/n\_components

- **beta\_0** (*float, Optional | defaults to None.*) – The precision prior on the mean distribution (Gaussian). Must be greater than 0.

If None, the value is set to 1.0

- **nu\_0** (*float, Optional | defaults to None.*) – The prior of the number of degrees of freedom on the covariance distributions (Wishart). Must be greater or equal to dim.

If None, the value is set to dim

- **means\_prior** (*array (dim,), Optional | defaults to None*) – The prior value to compute the value of the means.

If None, the value is set to the mean of points\_data

- **cov\_wishart\_prior** (*type depends on covariance\_type, Optional | defaults to None*) – If covariance\_type is ‘full’ type must be array (dim,dim) If covariance\_type is ‘spherical’ type must be float The prior value to compute the value of the precisions.

If None, the value is set to the covariance of points\_data

**name**

*str* – The name of the method : ‘VBGMM’

**alpha**

*array of floats (n\_components,)* – Contains the parameters of the weight distribution (Dirichlet)

**beta**

*array of floats (n\_components,)* – Contains coefficients which are multiplied with the precision matrices to form the precision matrix on the Gaussian distribution of the means.

**nu**

*array of floats (n\_components,)* – Contains the number of degrees of freedom on the distribution of covariance matrices.

**\_inv\_prec**

*array of floats (n\_components,dim,dim)* – Contains the equivalent of the matrix W described in Bishop’s book. It is proportional to cov.

**\_log\_det\_inv\_prec**

*array of floats (n\_components,)* – Contains the logarithm of the determinant of W matrices.

**cov**

*array of floats (n\_components,dim,dim)* – Contains the computed covariance matrices of the mixture.

**means**

*array of floats (n\_components,dim)* – Contains the computed means of the mixture.

**log\_weights**

*array of floats (n\_components,)* – Contains the logarithm of weights of each cluster.

**iter**

*int* – The number of iterations computed with the method fit()

**convergence\_criterion\_data**

*array of floats (iter,)* – Stores the value of the convergence criterion computed with data on which the model is fitted.

**convergence\_criterion\_test**

*array of floats (iter,)* | *if \_early\_stopping only* – Stores the value of the convergence criterion computed with test data if it exists.

**\_is\_initialized**

*bool* – Ensures that the method \_initialize() has been used before using other methods such as score() or predict\_log\_assignements().

**Raises** ValueError : if the parameters are inconsistent, for example if the cluster number is negative, init\_type is not in [‘resp’, ‘mcw’]...

## References

‘Pattern Recognition and Machine Learning’, Bishop

**fit** (*points\_data*, *points\_test=None*, *tol=0.001*, *patience=None*, *n\_iter\_max=100*, *n\_iter\_fix=None*, *saving=None*, *file\_name='model'*, *saving\_iter=2*)  
The EM algorithm

**Parameters**

- **points\_data** (*array (n\_points, dim)*) – A 2D array of points on which the model will be trained
- **tol** (*float, defaults to 1e-3*) – The EM algorithm will stop when the difference between two steps regarding the convergence criterion is less than tol.
- **n\_iter\_max** (*int, defaults to 100*) – number of iterations maximum that can be done
- **saving\_iter** (*int | defaults 2*) – An int to know how often the model is saved (see saving below).
- **file\_name** (*str | defaults model*) – The name of the file (including the path).

**Other Parameters**

- **points\_test** (*array (n\_points\_bis, dim) | Optional*) – A 2D array of points on which the model will be tested.
- **patience** (*int | Optional*) – The number of iterations performed after having satisfied the convergence criterion
- **n\_iter\_fix** (*int | Optional*) – If not None, the algorithm will exactly do the number of iterations of n\_iter\_fix and stop.
- **saving** (*str | Optional*) – A string in ['log', 'linear']. In the following equations x is the parameter saving\_iter (see above).
  - If 'log', the model will be saved for all iterations which verify :  $\log(\text{iter})/\log(x)$  is an int
  - If 'linear' the model will be saved for all iterations which verify :  $\text{iter}/x$  is an int

**Returns**

**Return type** None

**predict\_log\_resp** (*points*)

This function returns the logarithm of each point's responsibilities

**Parameters** **points** (*array (n\_points\_bis, dim)*) – a 1D or 2D array of points with the same dimension as the problem

**Returns** **log\_resp** – the logarithm of the responsibilities

**Return type** array (n\_points\_bis, n\_components)

**read\_and\_init** (*group, points*)

A method reading a group of an hdf5 file to initialize DPGMM

**Parameters** **group** (*HDF5 group*) – A group of a hdf5 file in reading mode

**score** (*points*)

This function return the score of the function, which is the logarithm of the likelihood for GMM and the logarithm of the lower bound of the likelihood for VBGMM and DPGMM

**Parameters** **points** (*array (n\_points\_bis, dim)*) – a 1D or 2D array of points with the same dimension as the problem

**Returns** **score**

**Return type** float

**simplified\_model** (*points*)

A method creating a new model with simplified parameters: clusters unused are removed

**Parameters** *points* (*an array (n\_points, dim)*) –

**Returns** **GM**

**Return type** an instance of the same type of self: GMM,VBGMM or DPGMM

**write** (*group*)

A method creating datasets in a group of an hdf5 file in order to save the model

**Parameters** *group* (*HDF5 group*) – A group of a hdf5 file in reading mode

## Dirichlet Process Gaussian Mixture Model (DPGMM)

```
class megamix.batch.DPVariationalGaussianMixture (n_components=1,          init='kmeans',
                                                  alpha_0=None,          beta_0=None,
                                                  nu_0=None,          means_prior=None,
                                                  cov_wishart_prior=None,
                                                  reg_covar=1e-06,          type_init='resp',
                                                  n_jobs=1, pypcoeff=0, boost=None)
```

Variational Bayesian Estimation of a Gaussian Mixture with Dirichlet Process

This class allows to infer an approximate posterior distribution over the parameters of a Gaussian mixture distribution.

The weights distribution follows a Dirichlet Process with attribute alpha.

### Parameters

- **n\_components** (*int, defaults to 1.*) – Number of clusters used.
- **init** (*str, defaults to 'kmeans'.*) – Method used in order to perform the initialization, must be in ['random', 'plus', 'AF\_KMC', 'kmeans', 'GMM', 'VBGMM'].
- **reg\_covar** (*float, defaults to 1e-6*) – In order to avoid null covariances this float is added to the diagonal of covariance matrices.
- **type\_init** (*str, defaults to 'resp'.*) – The algorithm is initialized using this data (responsibilities if 'resp' or means, covariances and weights if 'mcw').

### Other Parameters

- **alpha\_0** (*float, Optional | defaults to None.*) – The prior parameter on the weight distribution (Beta). A high value of alpha\_0 will lead to equal weights, while a low value will allow some clusters to shrink and disappear. Must be greater than 0.

If None, the value is set to 1/n\_components

- **beta\_0** (*float, Optional | defaults to None.*) – The precision prior on the mean distribution (Gaussian). Must be greater than 0.

If None, the value is set to 1.0

- **nu\_0** (*float, Optional | defaults to None.*) – The prior of the number of degrees of freedom on the covariance distributions (Wishart). Must be greater or equal to dim.

If None, the value is set to dim

- **means\_prior** (*array (dim,), Optional | defaults to None*) – The prior value to compute the value of the means.

If None, the value is set to the mean of points\_data

- **cov\_wishart\_prior** (*type depends on covariance\_type, Optional | defaults to None*) – If covariance\_type is ‘full’ type must be array (dim,dim) If covariance\_type is ‘spherical’ type must be float The prior value to compute the value of the precisions.
- **pypcoeff** (*float | defaults to 0*) – If 0 the weights are generated according to a Dirichlet Process If >0 and <=1 the weights are generated according to a Pitman-Yor Process.

**name**

*str* – The name of the method : ‘VBGMM’

**alpha**

*array of floats (n\_components,2)* – Contains the parameters of the weight distribution (Beta)

**beta**

*array of floats (n\_components,)* – Contains coefficients which are multiplied with the precision matrices to form the precision matrix on the Gaussian distribution of the means.

**nu**

*array of floats (n\_components,)* – Contains the number of degrees of freedom on the distribution of covariance matrices.

**\_inv\_prec**

*array of floats (n\_components,dim,dim)* – Contains the equivalent of the matrix W described in Bishop’s book. It is proportional to cov.

**\_log\_det\_inv\_prec**

*array of floats (n\_components,)* – Contains the logarithm of the determinant of W matrices.

**cov**

*array of floats (n\_components,dim,dim)* – Contains the computed covariance matrices of the mixture.

**means**

*array of floats (n\_components,dim)* – Contains the computed means of the mixture.

**log\_weights**

*array of floats (n\_components,)* – Contains the logarithm of weights of each cluster.

**iter**

*int* – The number of iterations computed with the method fit()

**convergence\_criterion\_data**

*array of floats (iter,)* – Stores the value of the convergence criterion computed with data on which the model is fitted.

**convergence\_criterion\_test**

*array of floats (iter,) | if \_early\_stopping only* – Stores the value of the convergence criterion computed with test data if it exists.

**\_is\_initialized**

*bool* – Ensures that the method \_initialize() has been used before using other methods such as score() or predict\_log\_assignements().

**Raises** ValueError : if the parameters are inconsistent, for example if the cluster number is negative, init\_type is not in [‘resp’, ‘mcw’]...

## References

‘Variational Inference for Dirichlet Process Mixtures’, D. Blei and M. Jordan

**fit** (*points\_data*, *points\_test=None*, *tol=0.001*, *patience=None*, *n\_iter\_max=100*, *n\_iter\_fix=None*, *saving=None*, *file\_name='model'*, *saving\_iter=2*)  
The EM algorithm

#### Parameters

- **points\_data** (*array (n\_points, dim)*) – A 2D array of points on which the model will be trained
- **tol** (*float, defaults to 1e-3*) – The EM algorithm will stop when the difference between two steps regarding the convergence criterion is less than tol.
- **n\_iter\_max** (*int, defaults to 100*) – number of iterations maximum that can be done
- **saving\_iter** (*int | defaults 2*) – An int to know how often the model is saved (see saving below).
- **file\_name** (*str | defaults model*) – The name of the file (including the path).

#### Other Parameters

- **points\_test** (*array (n\_points\_bis, dim) | Optional*) – A 2D array of points on which the model will be tested.
- **patience** (*int | Optional*) – The number of iterations performed after having satisfied the convergence criterion
- **n\_iter\_fix** (*int | Optional*) – If not None, the algorithm will exactly do the number of iterations of *n\_iter\_fix* and stop.
- **saving** (*str | Optional*) – A string in ['log', 'linear']. In the following equations *x* is the parameter *saving\_iter* (see above).
  - If 'log', the model will be saved for all iterations which verify :  $\log(\text{iter})/\log(x)$  is an int
  - If 'linear' the model will be saved for all iterations which verify :  $\text{iter}/x$  is an int

#### Returns

**Return type** None

**predict\_log\_resp** (*points*)

This function returns the logarithm of each point's responsibilities

**Parameters** **points** (*array (n\_points\_bis, dim)*) – a 1D or 2D array of points with the same dimension as the problem

**Returns** **log\_resp** – the logarithm of the responsibilities

**Return type** array (n\_points\_bis, n\_components)

**read\_and\_init** (*group, points*)

A method reading a group of an hdf5 file to initialize DPGMM

**Parameters** **group** (*HDF5 group*) – A group of a hdf5 file in reading mode

**score** (*points*)

This function return the score of the function, which is the logarithm of the likelihood for GMM and the logarithm of the lower bound of the likelihood for VBGMM and DPGMM

**Parameters** **points** (*array (n\_points\_bis, dim)*) – a 1D or 2D array of points with the same dimension as the problem

**Returns score**

**Return type** float

**simplified\_model** (*points*)

A method creating a new model with simplified parameters: clusters unused are removed

**Parameters** *points* (*an array (n\_points, dim)*) –

**Returns** GM

**Return type** an instance of the same type of self: GMM, VBGMM or DPGMM

**write** (*group*)

A method creating datasets in a group of an hdf5 file in order to save the model

**Parameters** *group* (*HDF5 group*) – A group of a hdf5 file in reading mode

## Online versions of the algorithm

### Kmeans

**class** megamix.online.**Kmeans** (*n\_components=1, window=1, kappa=1.0*)

Kmeans model.

**Parameters**

- **n\_components** (*int, defaults to 1.*) – Number of clusters used.
- **window** (*int, defaults to 1*) – The number of points used at the same time in order to update the parameters.
- **kappa** (*double, defaults to 1.0*) – A coefficient in ]0.0,1.0] which give weight or not to the new points compared to the ones already used.
  - If kappa is nearly null, the new points have a big weight and the model may take a lot of time to stabilize.
  - If kappa = 1.0, the new points won't have a lot of weight and the model may not move enough from its initialization.

**name**

*str* – The name of the method : 'Kmeans'

**log\_weights**

*array of floats (n\_components)* – Contains the logarithm of the mixing coefficients of the model.

**means**

*array of floats (n\_components, dim)* – Contains the computed means of the model.

**iter**

*int* – The number of points which have been used to compute the model.

**Raises** ValueError : if the parameters are inconsistent, for example if the cluster number is negative, init\_type is not in ['resp', 'mcw']...

### References

*Online but Accurate Inference for Latent Variable Models with Local Gibbs Sampling*, C. Dupuy & F. Bach  
 'The remarkable k-means++' <<https://normaldeviate.wordpress.com/2012/09/30/the-remarkable-k-means/>>'\_

**fit** (*points\_data*, *points\_test*=None, *saving*=None, *file\_name*='model', *check\_convergence\_iter*=None, *saving\_iter*=2)  
The k-means algorithm

#### Parameters

- **points\_data** (*array (n\_points, dim)*) – A 2D array of points on which the model will be trained.
- **saving\_iter** (*int | defaults 2*) – An int to know how often the model is saved (see saving below).
- **file\_name** (*str | defaults model*) – The name of the file (including the path).

#### Other Parameters

- **points\_test** (*an array (n\_points2, dim) | Optional*) – Data used to do early stopping (avoid overfitting)
- **check\_convergence\_iter** (*int | Optional*) – If points\_test are given, convergence criterion will be computed every check\_convergence\_iter iterations. If no value is given and points\_test is not None, it will raise an Error.
- **saving** (*str | Optional*) – A string in ['log', 'linear']. In the following equations x is the parameter saving\_iter (see above).
  - If 'log', the model will be saved for all iterations which verify :  $\log(\text{iter})/\log(x)$  is an int
  - If 'linear' the model will be saved for all iterations which verify :  $\text{iter}/x$  is an int

#### Returns

**Return type** None

**get** (*name*)

A getter to allow the user to get the attributes with the cython version.

**Parameters name** (*str*) – The name of the parameter. Must be in ['\_is\_initialized', 'log\_weights', 'means', 'iter', 'window', 'kappa', 'name']

#### Returns

**Return type** The wanted parameter (may be an array, a boolean, an int or a string)

**initialize** (*points*)

This method initializes the Gaussian Mixture by setting the values of the means, covariances and weights.

**Parameters points** (*an array (n\_points, dim)*) – Data on which the model is initialized.

**predict\_assignments** (*points*)

This function return the hard assignments of points once the model is fitted.

**score** (*points, assignments*=None)

This method returns the distortion measurement at the end of the k\_means.

#### Parameters

- **points** (*an array (n\_points, dim)*) –
- **assignments** (*an array (n\_components, dim)*) – an array containing the responsibilities of the clusters



**Returns distortion**

**Return type** (float)

## Gaussian Mixture Model (GMM)

```
class megamix.online.GaussianMixture(n_components=1, kappa=1.0, reg_covar=1e-06, window=1, update=False)
```

Gaussian Mixture Model

Representation of a Gaussian mixture model probability distribution. This class allows to estimate the parameters of a Gaussian mixture distribution (with full covariance matrices only).

### Parameters

- **n\_components** (*int*, defaults to 1) – Number of clusters used.
- **kappa** (*double*, defaults to 1.0) – A coefficient in ]0.0,1.0] which give weight or not to the new points compared to the ones already used.
  - If kappa is nearly null, the new points have a big weight and the model may take a lot of time to stabilize.
  - If kappa = 1.0, the new points won't have a lot of weight and the model may not move enough from its initialization.
- **window** (*int*, defaults to 1) – The number of points used at the same time in order to update the parameters.
- **update** (*bool*, defaults to False) – If True, the matrices of Cholesky of covariance matrices are updated, else they are computed at each iteration. Set it to True if window < dimension of the problem.
- **reg\_covar** (*float*, defaults to 1e-6) – In order to avoid null covariances this float is added to the diagonal of covariance matrices.

**name**

*str* – The name of the method : 'GMM'

**cov**

*array of floats (n\_components,dim,dim)* – Contains the computed covariance matrices of the mixture.

**means**

*array of floats (n\_components,dim)* – Contains the computed means of the mixture.

**log\_weights**

*array of floats (n\_components,)* – Contains the logarithm of weights of each cluster.

**iter**

*int* – The number of iterations computed with the method fit()

**Raises** ValueError : if the parameters are inconsistent, for example if the cluster number is negative, init\_type is not in ['resp','mcw']...

## References

*Online but Accurate Inference for Latent Variable Models with Local Gibbs Sampling*, C. Dupuy & F. Bach

**fit** (*points\_data*, *points\_test*=None, *saving*=None, *file\_name*='model', *check\_convergence\_iter*=None, *saving\_iter*=2)  
The EM algorithm

#### Parameters

- **points\_data** (*array (n\_points, dim)*) – A 2D array of points on which the model will be trained.
- **saving\_iter** (*int | defaults 2*) – An int to know how often the model is saved (see saving below).
- **file\_name** (*str | defaults model*) – The name of the file (including the path).

#### Other Parameters

- **points\_test** (*an array (n\_points2, dim) | Optional*) – Data used to do early stopping (avoid overfitting)
- **check\_convergence\_iter** (*int | Optional*) – If points\_test are given, convergence criterion will be computed every check\_convergence\_iter iterations. If no value is given and points\_test is not None, it will raise an Error.
- **saving** (*str | Optional*) – A string in ['log', 'linear']. In the following equations x is the parameter saving\_iter (see above).
  - If 'log', the model will be saved for all iterations which verify :  $\log(\text{iter})/\log(x)$  is an int
  - If 'linear' the model will be saved for all iterations which verify :  $\text{iter}/x$  is an int

#### Returns

**Return type** None

**get** (*name*)

A getter to allow the user to get the attributes with the cython version.

**Parameters name** (*str*) – The name of the parameter. Must be in ['\_is\_initialized', 'log\_weights', 'means', 'cov', 'cov\_chol', 'iter', 'window', 'kappa', 'name']

#### Returns

**Return type** The wanted parameter (may be an array, a boolean, an int or a string)

**initialize** (*points*)

This method initializes the Gaussian Mixture by setting the values of the means, covariances and weights.

**Parameters points** (*an array (n\_points, dim)*) – Data on which the model is initialized using the seeds of kmeans++.

**predict\_log\_resp** (*points*)

This function returns the logarithm of each point's responsibilities

**Parameters points** (*array (n\_points\_bis, dim)*) – a 1D or 2D array of points with the same dimension as the problem

**Returns log\_resp** – the logarithm of the responsibilities

**Return type** array (n\_points\_bis, n\_components)

**read\_and\_init** (*group, points*)

A method reading a group of an hdf5 file to initialize DPGMM

**Parameters** `group` (*HDF5 group*) – A group of a hdf5 file in reading mode

**score** (*points*)

This function return the score of the function, which is the logarithm of the likelihood for GMM and the logarithm of the lower bound of the likelihood for VBGMM and DPGMM

**Parameters** `points` (*array (n\_points\_bis, dim)*) – a 1D or 2D array of points with the same dimension as the problem

**Returns** `score`

**Return type** float

**write** (*group*)

A method creating datasets in a group of an hdf5 file in order to save the model

**Parameters** `group` (*HDF5 group*) – A group of a hdf5 file in reading mode



---

## Theory of Gaussian Mixture models

---

In this part are detailed the equations used in each algorithm. We use the same notations as Bishop's *Pattern Recognition and Machine Learning*.

Features:

- $\{x_1, x_2, \dots, x_N\}$  is the set of points

Parameters:

- $\mu_k$  is the center of the  $k^{th}$  cluster
- $\pi_k$  is the weight of the  $k^{th}$  cluster
- $\Sigma_k$  is the covariance matrix of the  $k^{th}$  cluster
- $K$  is the number of clusters
- $N$  is the number of points
- $d$  is the dimension of the problem

Other notations specific to the methods will be introduced later.

### K-means

An iteration of K-means includes:

- The *E step* : a label is assigned to each point (hard assignement) according to the means.
- The *M step* : means are computed according to the parameters.
- The computation of the *convergence criterion* : the algorithm uses the distortion as described below.

## E step

The algorithm produces a matrix of responsibilities according to the following equation:

$$r_{nk} = \begin{cases} 1 & \text{if } k = \arg \min_{1 \leq j \leq K} \|x_n - \mu_j\|^2 \\ 0 & \text{otherwise} \end{cases}$$

The value of the case at the  $i^{th}$  row and  $j^{th}$  column is 1 if the  $i^{th}$  point belongs to the  $j^{th}$  cluster and 0 otherwise.

## M step

The mean of a cluster is simply the mean of all the points belonging to this latter:

$$\mu_k = \frac{\sum_{n=1}^N r_{nk} x_n}{\sum_{n=1}^N r_{nk}}$$

The weight of the cluster  $k$ , which is the number of points belonging to this latter, can be expressed as:

$$N_k = \sum_{n=1}^N r_{nk}$$

The mixing coefficients, which represent the proportion of points in a cluster, can be expressed as:

$$\pi_k = \frac{N_k}{N}$$

## Convergence criterion

The convergence criterion is the distortion defined as the sum of the norms of the difference between each point and the mean of the cluster it is belonging to:

$$D = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2$$

The distortion should only decrease during the execution of the algorithm. The model stops when the difference between the value of the convergence criterion at the previous iteration and the current iteration is less or equal to a threshold  $tol$  :

$$D_{previous} - D_{current} \leq tol$$

## Gaussian Mixture Model (GMM)

An iteration of GMM includes:

- The *E step* :  $K$  probabilities of belonging to each cluster are assigned to each point
- The *M step* : weights, means and covariances are computed according to the parameters.
- The computation of the *convergence criterion* : the algorithm uses the loglikelihood as described below.

## E step

The algorithm produces a matrix of responsibilities according to the following equation:

$$r_{nk} = \frac{\pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(x_n | \mu_j, \Sigma_j)}$$

The value of the case at the  $i^{th}$  row and  $j^{th}$  column is the probability that the point  $i$  belongs to the cluster  $j$ .

## M step

The weight of the cluster  $k$ , which is the number of points belonging to this latter, can be expressed as:

$$N_k = \sum_{n=1}^N r_{nk}$$

The mixing coefficients, which represent the proportion of points in a cluster, can be expressed as:

$$\pi_k = \frac{N_k}{N}$$

As in the Kmeans algorithm, the mean of a cluster is the mean of all the points belonging to this latter:

$$\mu_k = \frac{\sum_{n=1}^N r_{nk} x_n}{N_k}$$

The covariance of the cluster  $k$  can be expressed as:

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N r_{nk} (x_n - \mu_k)(x_n - \mu_k)^T$$

These results have been obtained by derivating the maximum loglikelihood described in the following section.

## Convergence criterion

The convergence criterion used in the Gaussian Mixture Model algorithm is the maximum log likelihood:

$$\sum_{n=1}^N \ln \sum_{k=1}^K \pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k)$$

Setting its derivatives to 0 gives the empirical terms described in the M step.

## Variational Gaussian Mixture Model (VBGMM)

In this model, we introduce three new hyperparameters and two distributions which governs the three essential parameters of the model: the mixing coefficients, the means and the covariances.

The mixing coefficients are generated with a Dirichlet Distribution:

$$q(\pi_k) = \text{Dir}(\pi | \alpha_k) = C(\alpha_k) \pi_k^{\alpha_k - 1}$$

The computation of  $\alpha_k$  is described in the M step.

Then we introduce an independant Gaussian-Wishart law governing the mean and precision of each gaussian component:

$$\begin{aligned} q(\mu_k, \Sigma_k) &= q(\mu_k | \Sigma_k) q(\Sigma_k) \\ &= \mathcal{N}(\mu_k | m_k, (\beta_k \Sigma_k)^{-1}) \mathcal{W}(\Gamma_k | W_k, \nu_k) \end{aligned}$$

The computation of the terms involved in this equation are described in the M step.

**E step**

**M step**

**Convergence criterion**

## **Dirichlet Process Gaussian Mixture Model (DPGMM)**

**E step**

**M step**

**Convergence criterion**

## **Pitman-Yor Process Gaussian Mixture Model (PYPGMM)**



## Symbols

`_inv_prec` (DPVariationalGaussianMixture attribute), 17  
`_inv_prec` (VariationalGaussianMixture attribute), 14  
`_is_initialized` (DPVariationalGaussianMixture attribute), 17  
`_is_initialized` (GaussianMixture attribute), 11  
`_is_initialized` (Kmeans attribute), 9  
`_is_initialized` (VariationalGaussianMixture attribute), 14  
`_log_det_inv_prec` (DPVariationalGaussianMixture attribute), 17  
`_log_det_inv_prec` (VariationalGaussianMixture attribute), 14

## A

`alpha` (DPVariationalGaussianMixture attribute), 17  
`alpha` (VariationalGaussianMixture attribute), 14

## B

`beta` (DPVariationalGaussianMixture attribute), 17  
`beta` (VariationalGaussianMixture attribute), 14

## C

`convergence_criterion_data` (DPVariationalGaussianMixture attribute), 17  
`convergence_criterion_data` (GaussianMixture attribute), 11  
`convergence_criterion_data` (VariationalGaussianMixture attribute), 14  
`convergence_criterion_test` (DPVariationalGaussianMixture attribute), 17  
`convergence_criterion_test` (GaussianMixture attribute), 11  
`convergence_criterion_test` (VariationalGaussianMixture attribute), 14  
`cov` (DPVariationalGaussianMixture attribute), 17  
`cov` (GaussianMixture attribute), 11, 21  
`cov` (VariationalGaussianMixture attribute), 14

## D

`DPVariationalGaussianMixture` (class in megamix.batch), 16

## F

`fit()` (megamix.batch.DPVariationalGaussianMixture method), 17  
`fit()` (megamix.batch.GaussianMixture method), 11  
`fit()` (megamix.batch.Kmeans method), 10  
`fit()` (megamix.batch.VariationalGaussianMixture method), 14  
`fit()` (megamix.online.GaussianMixture method), 21  
`fit()` (megamix.online.Kmeans method), 19

## G

`GaussianMixture` (class in megamix.batch), 11  
`GaussianMixture` (class in megamix.online), 21  
`get()` (megamix.online.GaussianMixture method), 22  
`get()` (megamix.online.Kmeans method), 20

## I

`initialize()` (megamix.online.GaussianMixture method), 22  
`initialize()` (megamix.online.Kmeans method), 20  
`iter` (DPVariationalGaussianMixture attribute), 17  
`iter` (GaussianMixture attribute), 11, 21  
`iter` (Kmeans attribute), 9, 19  
`iter` (VariationalGaussianMixture attribute), 14

## K

`Kmeans` (class in megamix.batch), 9  
`Kmeans` (class in megamix.online), 19

## L

`log_weights` (DPVariationalGaussianMixture attribute), 17  
`log_weights` (GaussianMixture attribute), 11, 21  
`log_weights` (Kmeans attribute), 9, 19  
`log_weights` (VariationalGaussianMixture attribute), 14

## M

means (DPVariationalGaussianMixture attribute), 17  
means (GaussianMixture attribute), 11, 21  
means (Kmeans attribute), 9, 19  
means (VariationalGaussianMixture attribute), 14

## N

name (DPVariationalGaussianMixture attribute), 17  
name (GaussianMixture attribute), 11, 21  
name (Kmeans attribute), 9, 19  
name (VariationalGaussianMixture attribute), 14  
nu (DPVariationalGaussianMixture attribute), 17  
nu (VariationalGaussianMixture attribute), 14

## P

predict\_assignments() (megamix.batch.Kmeans method), 10  
predict\_assignments() (megamix.online.Kmeans method), 20  
predict\_log\_resp() (megamix.batch.DPVariationalGaussianMixture method), 18  
predict\_log\_resp() (megamix.batch.GaussianMixture method), 12  
predict\_log\_resp() (megamix.batch.VariationalGaussianMixture method), 15  
predict\_log\_resp() (megamix.online.GaussianMixture method), 22

## R

read\_and\_init() (megamix.batch.DPVariationalGaussianMixture method), 18  
read\_and\_init() (megamix.batch.GaussianMixture method), 12  
read\_and\_init() (megamix.batch.VariationalGaussianMixture method), 15  
read\_and\_init() (megamix.online.GaussianMixture method), 22

## S

score() (megamix.batch.DPVariationalGaussianMixture method), 18  
score() (megamix.batch.GaussianMixture method), 12  
score() (megamix.batch.Kmeans method), 10  
score() (megamix.batch.VariationalGaussianMixture method), 15  
score() (megamix.online.GaussianMixture method), 23  
score() (megamix.online.Kmeans method), 20  
simplified\_model() (megamix.batch.DPVariationalGaussianMixture method), 19  
simplified\_model() (megamix.batch.GaussianMixture method), 12  
simplified\_model() (megamix.batch.VariationalGaussianMixture method), 15

## V

VariationalGaussianMixture (class in megamix.batch), 13

## W

write() (megamix.batch.DPVariationalGaussianMixture method), 19  
write() (megamix.batch.GaussianMixture method), 13  
write() (megamix.batch.VariationalGaussianMixture method), 16  
write() (megamix.online.GaussianMixture method), 23