

---

# **AutoMix**

***Release 2.0***

**Jan 29, 2020**



---

## Contents

---

<b>1</b>	<b>Tutorial</b>	<b>1</b>
1.1	Example . . . . .	1
1.2	AutoMix Sampler . . . . .	3
1.3	Run Statistics . . . . .	4
<b>2</b>	<b>Compilation</b>	<b>7</b>
2.1	Library Compilation . . . . .	7
2.2	Compiling your program against libautomix . . . . .	7
2.3	Compiling test and tutorial . . . . .	8
2.4	Compiling “user*” legacy programs . . . . .	8
2.5	Clean . . . . .	8
<b>3</b>	<b>Public API</b>	<b>9</b>
3.1	Functions: . . . . .	9
3.2	C struct’s . . . . .	10
3.3	Typedef’s and Enum’s . . . . .	12
<b>4</b>	<b>License</b>	<b>15</b>
<b>5</b>	<b>What is Reversible Jump MCMC?</b>	<b>17</b>
<b>6</b>	<b>Main advantages of AutoMix:</b>	<b>19</b>
<b>7</b>	<b>Where to go from here?</b>	<b>21</b>
	<b>Index</b>	<b>23</b>



## 1.1 Example

Suppose we have 10 samples from a random process  $X$  that we know not much about.

$$X = \{0.2, 0.13, 0.35, 0.17, 0.89, 0.33, 0.78, 0.23, 0.54, 0.16\}$$

We will assume the samples are independent, identically distributed (IID) from some unknown probability density distribution (PDF). The nature of the problem indicates that it could come from a Normal distribution, a Gamma distribution or perhaps a Beta distribution. We don't have any other information on the parameters,

So we propose three models to account for what we know:

- **Model 1:** Normal distribution with arbitrary parameters.

$$p_1(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-x_0)^2}{2\sigma^2}\right)$$

This model has undertermined parameters  $\sigma$  and  $x_0$ .

- **Model 2:** Beta distribution

$$p_2(x) = \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}$$

This model has undetermined parameters  $\alpha$  and  $\beta$  under the restriction  $\alpha > 0$  and  $\beta > 0$ .

- **Model 3:** Gamma distribution

$$p_3(x) = \frac{\beta^\alpha x^{\alpha-1} e^{-\beta x}}{\Gamma(\alpha)}$$

This model has undetermined parameters  $\alpha$  and  $\beta$  under the restriction  $\alpha > 0$  and  $\beta > 0$ .

---

**Note:** The dimensions for the problems are 2 for all 3 models.

---

We would like to find the probability of each model given the evidence (called posterior in a Bayesian analysis.)

Assuming independent, identically distributed (IID) samples, the probability of the evidence given the model (or likelihood) can be easily calculated

$$p(X|M) = \prod_i p_M(x_i)$$

The probability of the model given the evidence is then

$$p(M|X) \propto \prod_i p_M(x_i)p(M)$$

In what follows we will assume all models are equally likely and set  $p(M) = 1/3 \forall M$ .

The AutoMix sampler needs to be provided with a function that returns the target distribution (the posterior in our example), for each model. The function prototype must be in the form:

```
double targetDistribution(int model, int model_dim, double *params);
```

The log-posteriors need to be implemented up to an additive constant. A quick implementation of the log-PDFs could be:

```
#include "float.h"
#include <math.h>

int nsamples = 10;
double data_samples[] = {0.2, 0.13, 0.35, 0.17, 0.89,
                        0.33, 0.78, 0.23, 0.54, 0.16};

double logp_normal(double sigma, double x0) {
    double prod = 0;
    for (int i = 0; i < nsamples; i++) {
        double x = data_samples[i];
        prod += -(x - x0) * (x - x0);
    }
    prod = -nsamples * log(sigma) + prod / (2.0 * sigma * sigma);
    return prod;
}

double logp_beta(double alpha, double beta) {
    if (alpha <= 0.0 || beta <= 0.0) {
        return -DBL_MAX;
    }
    double prod = 0.0;
    for (int i = 0; i < nsamples; i++) {
        double x = data_samples[i];
        prod += (alpha - 1.0) * log(x) + (beta - 1.0) * log(1.0 - x);
    }
    prod +=
        nsamples * (loggamma(alpha + beta) - loggamma(alpha) - loggamma(beta));
    return prod;
}

double logp_gamma(double alpha, double beta) {
```

(continues on next page)

(continued from previous page)

```
double prod = 0.0;
for (int i = 0; i < nsamples; ++i) {
    double x = data_samples[i];
    prod += (alpha - 1.0) * log(x) - beta * x;
}
prod += nsamples * (alpha * log(beta) - loggamma(alpha));
return prod;
}
```

Now that we have defined our problem and log-posteriors, we can set up AutoMix to generate samples from the posteriors of each model.

## 1.2 AutoMix Sampler

Below is an example of a minimal call to generate 1,000 samples from a given problem:

```
#include "automix.h"

int main() {
    int nmodels = 3;
    int model_dims[] = {2, 2, 2};
    double initRWM[] = {0.5, 0.5, 2.0, 2.0, 9.0, 2.0};
    amSampler am;
    initAMSampler(&am, nmodels, model_dims, logposterior, initRWM);
    estimate_conditional_probs(&am, 100000);
    burn_samples(&am, 10000);
    int nsweeps = 100000;
    rjcmc_samples(&am, nsweeps);
    freeAMSampler(&am);
    return 0;
}
```

This is a simple set-up for an AutoMix program. Let's analyze it by parts.

The first line includes the AutoMix header file, where the *amSampler* structure and automix functions are defined:

```
#include "automix.h"
```

To initiate the *amSampler* struct, we need to set 4 things:

- The number of models we will use (3 in our example):

```
int nmodels = 3;
```

- The dimensions of each model:

```
int model_dims[] = {2, 2, 2};
```

- A function that returns the **logarithm** of the **probability density function** (PDF), (or the log-posterior distribution in the context of a Bayesian analysis) for each of the models in our problem. It must have the prototype:

```
double logposterior(int model, double *params);
```

And an implementation:

```
double logposterior(int model, double *params) {  
    if (model == 1) {  
        return logp_normal(params[0], params[1]);  
    } else if (model == 2) {  
        return logp_beta(params[0], params[1]);  
    } else if (model == 3) {  
        return logp_gamma(params[0], params[1]);  
    }  
    return 0.0;  
}
```

- The initial value of the parameters for every model. This should be given as a 1d continuous array with an appropriate initial value for each model. For example the Beta distribution is bounded to  $[0, 1]$  and any initial value has to be in that interval as well. This is to avoid starting the chain in a forbidden region of the parameter space or very far from the average values of the parameters:

```
double initRWM[] = {0.5, 0.5, 2.0, 2.0, 9.0, 2.0};
```

In our case we start our chain with values  $\sigma = 0.5$ ,  $x_0 = 0.5$  for the Gaussian model;  $\alpha = 2.0$ ,  $\beta = 2.0$  for the Beta model; and  $\alpha = 9.0$ ,  $\beta = 2.0$  for the Gamma model.

Once all of the above is defined for our problem we can init our *amSampler*:

```
amSampler am;  
initAMSampler(&am, nmodels, model_dims, logposterior, initRWM);
```

The next step is to estimate the conditional probabilities for our posterior model to create a Proposal Distribution with a multi-modal Normal mixture:

```
estimate_conditional_probs(&am, 100000);
```

We just need to pass the automix sampler struct and the number of sweeps we want for the estimation.

It is recommended to “burn” some initial samples to let the MCMC chain achieve convergence. We do this with the following line:

```
burn_samples(&am, 10000);
```

Finally we can create however many RJMCMC samples as we want:

```
int nsweeps = 100000;  
rjmc_samples(&am, nsweeps);
```

## 1.3 Run Statistics

Our previous main program has one fundamental problem: the lack of output.

AutoMix saves the run statistics in different C data structures.

The most important is the struct *runStats*. In *amSampler* is the member *st*. There, we will find several statistics, including the parameters sampled during the RJMCMC run. The two most important ones are *st.k\_summary* and *st.theta\_summary*. For a full description of all the parameters, see *Public API*.

Just as an example, let’s print out the relative probability of each model:



```
#include <stdio.h>

int main() {
    ...
    rjmcnc_samples(&am, nsweeps);

    printf("p(M=1|E) = %lf\n", (double)am.st.ksummary[0] / nsweeps);
    printf("p(M=2|E) = %lf\n", (double)am.st.ksummary[1] / nsweeps);
    printf("p(M=3|E) = %lf\n", (double)am.st.ksummary[2] / nsweeps);

    freeAMSampler(&am);
    return 0;
}
```

This should print something like the following on screen:

```
p(M=1|E) = 0.792750
p(M=2|E) = 0.023890
p(M=3|E) = 0.183360
```



### 2.1 Library Compilation

To compile the automix library, simply type `make` on the command line:

```
$ make
```

This default command compiles only the library with optimization flags `-O3` on.

If you need to link your program against a library with debug information, type:

```
$ make DEBUG=1
```

instead.

If the compilation went without problems, you should see the `libautomix.so` file created.

Typically, libraries are installed in the `/usr/local/lib` directory, and include files such as `automix.h` are installed in `/usr/local/include`. If the default install location is enough for your needs, type:

```
$ sudo make install
```

`sudo` privilege may be needed to copy files into the `/usr/local` directory.

If you prefer to install in a different location you can provide it with the `PREFIX` command line flag:

```
$ make install PREFIX=/my/preferred/path/
```

### 2.2 Compiling your program against libautomix

If your program is contained in a single `main.c` file, and assuming no extra libraries or include files are required:

```
$ cc main.c -lautomix -o myprogram
```

should be enough to compile. If `/usr/local/` is not a default search place for your libraries, add a `-L/usr/local/lib` and a `-I/usr/local/include` to the compilation line:

```
$ cc main.c -L/usr/local/lib -I/usr/local/include -lautomix -o myprogram
```

You should have a `myprogram` executable compiled.

## 2.3 Compiling test and tutorial

Along with the library you can compile and run the tests:

```
$ make test
$ ./test
```

Or the tutorial provided in this documentation:

```
$ make tutorial
$ ./tutorial
```

## 2.4 Compiling “user\*” legacy programs

**Warning:** The user examples are deprecated on version 2.x. If you want to compile AutoMix using the old interface with `user*` files, please refer to version 1.x of this package.

The original AutoMix was program-oriented and offered several example programs. To compile the example programs type:

```
$ make examples
```

After compilation you should see the programs `amtoy1`, `amtoy2`, `amcpt`, `amcpters`, `amddi`, and `amrb9`.

You can compile each one of them individually as well, for example:

```
$ make amtoy1
$ ./amtoy1
```

## 2.5 Clean

To remove the executables and object (`.o`) files, type:

```
make clean
```

### 3.1 Functions:

All public functions require as a first argument a pointer to *amSampler* structure.

```
int initAMSampler(amSampler *am, int nmodels, int *model_dims,
targetFunc logpost, double *initRWM);
```

To initialize *amSampler* you need to provide:

- the number of models in *nmodels*;
- the dimensions for each model in *model\_dims*;
- the logarithm of the target distribution (or posterior distribution in a Bayesian analysis) in *logpost* (see *targetFunc*);
- and an array with initial conditions for each model in *initRWM*.

*initRWM* is a flattened array, with the initial values in contiguous order for each model. If *initRWM* is passed as *NULL*, then the initial values are randomly selected from a uniform distribution in [0, 1).

*initAMSampler* allocates necessary memory for most array in the rest of the structures.

```
void freeAMSampler(amSampler *am);
```

Once you finish working with *amSampler*, you should free the memory with a call to *freeAMSampler*.

```
void estimate_conditional_probs(amSampler *am, int nsweep2);
```

Before running RJMCMC sweeps, you need to call *estimate\_conditional\_probs()* to estimate the conditional probabilities of your models, to adjust for an appropriate proposal distribution.

The statistics collected during the *estimate\_conditional\_probs()* calls, are stored in the *cpstats* member inside the *amSampler* struct. For a full description of *cpstats*, see *condProbStats*.

```
void burn_samples(amSampler *am, int nburn);
```

It is advised to burn some samples at the beginning of a RJMCMC run to let the Markov chain converge.

```
void rjmc_samples(amSampler *am, int nsweep);
```

This is the function that will provide the RJMCMC samples.

The statistics collected during the `rjmcmc_samples()` calls, are stored in the `st` member inside the `amSampler` struct. For a full description of `st`, see `runStats`.

## 3.2 C struct's

### **runStats**

Contains statistics of the `rjmcmc_samples` call.

#### **naccrwm**

[unsigned long] Number of accepted Block-RWM.

#### **ntryrwm**

[unsigned long] Number of tried Block-RWM.

#### **naccrws**

[unsigned long] Number of accepted Single-RWM.

#### **ntryrws**

[unsigned long] Number of tried Single-RWM.

#### **nacctd**

[unsigned long] Number of accepted Auto RJ.

#### **ntrytd**

[unsigned long] Number of tried Auto RJ.

#### **theta\_summary**

[double \*\*\*] The theta parameter for each model, for each sweep. `theta_summary[i][j][k]` is the  $\theta_k$  component for the `j`-th sweep for the `i`-th model.

#### **theta\_summary\_len**

[int \*] the number of sweeps for `theta_summary` in model `k`. `theta_summary_len[1] = 20` means that the model 1 has 20  $\theta$  samples.

#### **ksummary**

[int \*] the number of times the model was visited.

### **amSampler**

#### **ch**

A `chainState` instance containing the chain state of the RJMCMC.

#### **jd**

A `proposalDist` instance containing the proposal distribution for the chain.

#### **cpstats**

A `condProbStats` instance that holds statistics for the estimation of conditional probabilities (if ran).

#### **st**

A `runStats` instance that holds statistics for the RJMCMC runs (if ran).

#### **doAdapt**

A `bool` value indicating if Adaptation is to be performed or not during the RJMCMC runs. Default value is `true`.

#### **doPerm**

A `bool` value indicating if Permutation is to be performed or not during the RJMCMC runs. Default value is `false`.

**student\_T\_dof**

int type. The degree of freedom of the Student's T distribution to sample from. If 0, sample from a Normal distribution instead. Default value is 0.

**am\_mixfit**

A *automix\_mix\_fit* value. Default is *FIGUEREIDO\_MIX\_FIT*.

**seed**

An unsigned long value to initialize the random number generator. Defaults to clock time.

**condProbStats****rwm\_summary\_len**

TBD.

**sig\_k\_rwm\_summary**

TBD.

**nacc\_ntry\_rwm**

TBD.

**nfitmix**

TBD.

**fitmix\_annulations**

TBD.

**fitmix\_costfnnew**

TBD.

**fitmix\_lpn**

TBD.

**fitmix\_Lkk**

TBD.

**proposalDist**

The proposal distribution is of the form:

$$p(\theta) = \sum_{k=1}^{k_{max}} L_k \exp \left( -\frac{1}{2} (\theta_k - \mu_k) B \cdot B^T (\theta_k - \mu_k) \right).$$

**nmodels**

The number of models in M (referred as  $K_{max}$  in thesis, p 143)

**nMixComps**

The number of mixture components for each model (referred as  $L_k$  in thesis, p144)

**model\_dims**

The dimension of the  $\theta$  parameter space for each model.

**lambda**

The relative weights for each mixture component for each model.  $\text{lambda}[i][j]$  is the weight of the j-th mixture component for model  $M_i$ .

**mu**

The mean vector for the mixture component for each model.  $\text{mu}[i][j][k]$  is the k-th vector index for the j-th mixture component for model  $M_i$ .

**B**

$B \cdot B^T$  is the covariance matrix for each mixture component for each model.  $B[i][j][k][m]$  is the k,m index of the B matrix for mixture component j of model  $M_i$ .

**sig**

Vector of adapted RWM scale parameters for each model.

**chainState**

Struct to hold the MCMC chain state

**theta**

The current value of  $\theta_k$  in the chain.

**pk**

TBD.

**log\_posterior**

The current value of the log-posterior.

**current\_model\_k**

The current model  $k$  in the chain.

**mdim**

The dimension of the current model in the chain.

**current\_Lkk**

TBD

**nreinit**

TBD.

**reinit**

TBD.

**pkllim**

TBD.

**doBlockRWM**

A *bool* value to indicate if the chain should do a Block-RWM. This is done every 10 sweeps.

**isBurning**

A *bool* that indicates whether the chain is burning samples.

**sweep\_i**

The index of the chain.

## 3.3 Typedef's and Enum's

**targetFunc**

targetFunc is the prototype of the log-posterior function that must be passed to *amSampler* upon initialization.

```
typedef double (*targetFunc) (int model_k, double *x);
```

**bool**

bool is a typedef for int. C does not have a bool type but int is just as good.

```
typedef int bool;
```

**automix\_mix\_fit**

Enum to specify whether using Figueredo or AutoRJ in conditional probabilities estimation.



```
typedef enum { FIGUEREIDO_MIX_FIT = 0, AUTORJ_MIX_FIT } automix_mix_fit;
```



## CHAPTER 4

---

### License

---

The AutoMix sampler is free for personal and academic use, but users must reference the sampler as instructed below. For commercial use permission must be sought from the author. To seek permission for such use please send an e-mail to [d\\_hastie@hotmail.com](mailto:d_hastie@hotmail.com) outlining the desired usage.

Use of the AutoMix sampler is entirely at the user's own risk. It is the responsibility of the user to ensure that any conclusions made through the use of the AutoMix sampler are valid. The author accepts no responsibility whatsoever for any loss, financial or otherwise, that may arise in connection with the use of the sampler.

The AutoMix sampler is available from <http://www.davidhastie.me.uk/AutoMix>. The sampler may be modified and redistributed as desired but the author encourages users to register at the above site so that notice can be received of updates of the software.

Before use, please read this README file bundled with this software.

The AutoMix package is a C program for Unix-like systems, implementing the automatic Reversible Jump MCMC sampler of the same name described in Chapters 4, 5, and 6 of David Hastie's Ph.D. thesis (included in `docs/thesis`).

While the original AutoMix is highly useful, the fact that it can only be used as an executable can limit its applicability. **LibAutoMix makes the core algorithms of Source Extractor available as a library of stand-alone functions and data structures.** The code is completely derived from the original AutoMix code base and aims to produce results compatible with it whenever possible. LibAutoMix is a C library with no dependencies outside the standard library.



---

### What is Reversible Jump MCMC?

---

Reversible Jump Markov Chain Monte Carlo (RJMCMC) extends the standard MCMC sampler to include a discrete random variable  $k$  that represents the index of a *model*. So, instead of sampling from the usual parameter space of a given distribution, RJMCMC will also sample across different models (distributions).

The final samples of a RJMCMC reflect the probabilities of the parameters of each model, but also the relative probabilities of the models themselves.



---

### Main advantages of AutoMix:

---

- Reversible Jump MCMC allows sampling from several distributions (models) simultaneously.
- The different models may have different number of free parameters (dimension).
- The relative frequency of sampling for different models is proportional to the probability of the model. That means that AutoMix can be used as a model selection sampler.
- AutoMix requires minimum input from the user.
- AutoMix automatically adapts proposal distributions with a multi-modal Normal mixture.

**Caution:** Potential users should carefully understand the limitations of using the AutoMix sampler. The reliability of results from this sampler depends on many factors, including the scaling of the parameters and the degree of multimodality of the within-model conditionals of the target distribution.





## CHAPTER 7

---

### Where to go from here?

---

For a quick tour of the library, check out the [Tutorial](#).

To compile your program against the AutoMix library see [Compilation](#).

For a description of the full public API with statistics information, see [Public API](#).



## A

amSampler (*C type*), 10  
amSampler.am\_mixfit (*C member*), 11  
amSampler.ch (*C member*), 10  
amSampler.cpstats (*C member*), 10  
amSampler.doAdapt (*C member*), 10  
amSampler.doPerm (*C member*), 10  
amSampler.jd (*C member*), 10  
amSampler.seed (*C member*), 11  
amSampler.st (*C member*), 10  
amSampler.student\_T\_dof (*C member*), 10  
automix\_mix\_fit (*C type*), 12

## B

bool (*C type*), 12

## C

chainState (*C type*), 12  
chainState.current\_Lkk (*C member*), 12  
chainState.current\_model\_k (*C member*), 12  
chainState.doBlockRWM (*C member*), 12  
chainState.isBurning (*C member*), 12  
chainState.log\_posterior (*C member*), 12  
chainState.mdims (*C member*), 12  
chainState.nreinit (*C member*), 12  
chainState.pk (*C member*), 12  
chainState.pkllim (*C member*), 12  
chainState.reinit (*C member*), 12  
chainState.sweep\_i (*C member*), 12  
chainState.theta (*C member*), 12  
condProbStats (*C type*), 11  
condProbStats.fitmix\_annulations (*C member*), 11  
condProbStats.fitmix\_costfnnew (*C member*), 11  
condProbStats.fitmix\_Lkk (*C member*), 11  
condProbStats.fitmix\_lpn (*C member*), 11  
condProbStats.nacc\_ntry\_rwm (*C member*), 11  
condProbStats.nfitmix (*C member*), 11

condProbStats.rwm\_summary\_len (*C member*), 11  
condProbStats.sig\_k\_rwm\_summary (*C member*), 11

## P

proposalDist (*C type*), 11  
proposalDist.B (*C member*), 11  
proposalDist.lambda (*C member*), 11  
proposalDist.model\_dims (*C member*), 11  
proposalDist.mu (*C member*), 11  
proposalDist.nMixComps (*C member*), 11  
proposalDist.nmodels (*C member*), 11  
proposalDist.sig (*C member*), 12

## R

runStats (*C type*), 10  
runStats.ksummary (*C member*), 10  
runStats.naccrwm (*C member*), 10  
runStats.naccrwms (*C member*), 10  
runStats.nacctd (*C member*), 10  
runStats.ntryrwm (*C member*), 10  
runStats.ntryrwms (*C member*), 10  
runStats.ntrytd (*C member*), 10  
runStats.theta\_summary (*C member*), 10  
runStats.theta\_summary\_len (*C member*), 10

## T

targetFunc (*C type*), 12