

Inference in the linear model

Georgios Arampatzis

March 19, 2019

1 The model

In this tutorial we will describe how to run the Matlab code for the inference of the parameters in the linear model. We consider the model,

$$y = f(x; a, b) = ax + b. \quad (1)$$

We are given N_d data points $\mathbf{d} = \{y_i, x_i\}_{i=1}^{N_d}$. We consider the statistical model,

$$y = ax + b + \varepsilon, \quad (2)$$

where ε is a random variable that follows normal distribution with mean zeros and standard deviation σ . We assume that the data $y_i, i = 1, \dots, N_d$ are i.i.d samples from the model in eq. (2).

Set

$$\vartheta = (a, b, \sigma)$$

the vector of parameters to be inferred. The likelihood of \mathbf{d} under the assumption eq. (2) is given by,

$$p(\mathbf{d}|\vartheta) = \mathcal{N}(\mathbf{d}|\boldsymbol{\mu}, \Sigma), \quad (3)$$

where

$$\boldsymbol{\mu} = (f(x_1; a, b), \dots, f(x_{N_d}; a, b))^\top \quad \text{and} \quad \Sigma = \sigma^2 I. \quad (4)$$

The log-likelihood is given by,

$$\log p(\mathbf{d}|\vartheta) = -\frac{1}{2}N_d \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^{N_d} (y_i - f(x_i; \vartheta))^2. \quad (5)$$

The derivatives and the Fisher information matrix has been implemented in the provided Matlab functions. The user has to implement the derivatives $\frac{d}{d\vartheta_j} f(x_i, \vartheta)$.

2 Running the code

In order to validate our implementation we will first create synthetic data with known parameters and then infer these parameters. From the root directory (the directory that contains the folders **data** and **engines**),

```
1 cd data/linear
2 make_data
```

This script will create the file **data.mat**. Using these data we will infer a, b and σ . Finally, from the root directory run

```
1 addpath(engine/postprocessing/).
```

2.1 Likelihood optimization

First we can optimize the likelihood function using the CMA-ES algorithm. From the root directory

```
1 cd engine/optimize/  
2 edit run_linear_CMA.m
```

The important lines in this script are

```
1 addpath( './functions/linear/' )  
2  
3 loglike_func = 'linear_loglike_cma';  
4  
5 d = load( [ data_folder 'data.mat' ] );
```

and

```
1 opts.LBounds = [ -5 -5 0 ]';  
2 opts.UBounds = [ 5 5 10 ]';
```

The `addpath` command inserts the likelihood functions folder in the search path. In this case the log-likelihood is `linear_loglike_cma` and is located in the folder `./functions/linear/`. This function return the negative of the function `linear_loglike` that implements the log-likelihood as defined in eq. (5). The `opts.LBounds` and `opts.UBounds` variables set the search interval for the variables. Here, $a \in [-5, 5]$, $b \in [-5, 5]$ and $\sigma \in [0, 10]$. Finally, the command

```
1 start_parallel( 2, false );
```

will start a pool with two workers and run CMA in parallel. If you don't want parallel execution comment out this line. Running the script gives the results,

```
1 Maximum of log-likelihood at:  
2 a      = 1.867522  
3 b      = -1.741286  
4 sigma  = 1.770017
```

2.2 Sampling the posterior with BASIS

From the root directory

```
1 cd engine/sample/  
2 edit run_linear.m
```

uncomment the line

```
1 method = 'BASIS'; cov_check = 'NONE';
```

and run the script. After the end of execution plot the results by running

```
1 plotmatrix_hist(BASIS.theta)
```

or

```
1 load ../../data/linear/BASIS_NONE_5000.mat  
2 plotmatrix_hist(out_master.theta)
```

The results are shown in fig. 1.

2.3 Sampling the posterior with smTMCMC

First lets examine the function `linear_loglike_smMALA`. From root

```
1 cd engine/functions/linear/  
2 edit linear_loglike_smMALA.m
```

The important lines are

```

1 % the model
2 Y = theta(1)*x + theta(2);
3
4 % the derivatives (row vectors)
5 dY1 = x(:)';
6 dY2 = ones(1,Nd);

```

Here the variable `dY1` is a row vector that contains the derivative $\frac{d}{d\theta_1}f(x_i, \vartheta)$ at the i -th position and the variable `dY2` is a row vector that contains the derivative $\frac{d}{d\theta_2}f(x_i, \vartheta)$ at the i -th position.

In order to run, from the root directory

```

1 cd engine/sample/
2 edit run_linear.m

```

uncomment the line

```

1 method = 'smMALA'; cov_check = 'EIG';

```

and execute the script. After the end of execution plot the results by running

```

1 plotmatrix_hist(BASIS.theta)

```

or

```

1 load ../../data/linear/smMALA_EIG_5000.mat
2 plotmatrix_hist(out_master.theta)

```

The results are shown in fig. 2.

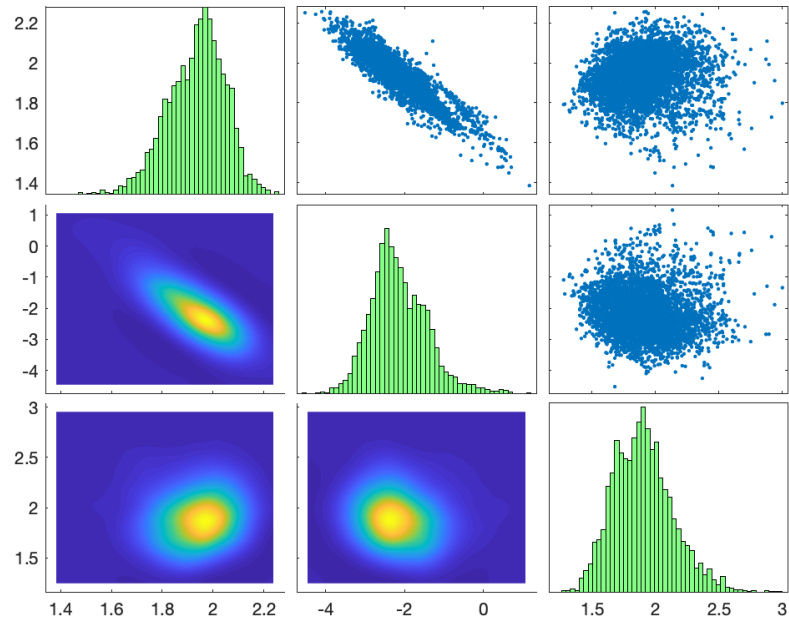


Figure 1: Results from BASIS.

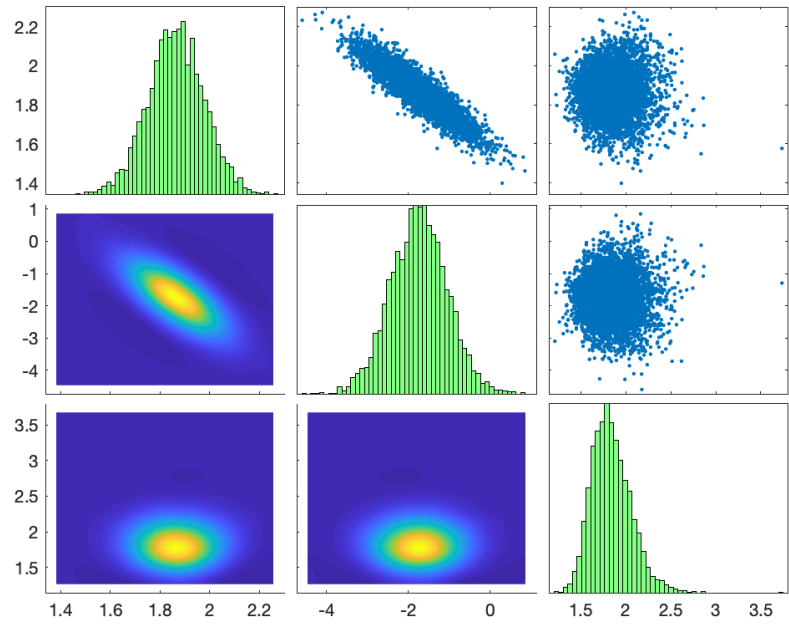


Figure 2: Results from smTMCMC.