

MLML User Guide

Jenny Qu Meng Zhou Qiang Song Elizabeth Hong Andrew Smith

April 17, 2013

mlml aims to simultaneously make consistent estimation of 5mC and 5hmC levels. It is developed to analyze DNA methylation level data. The input can be any pair of BS-seq, oxBS-seq and Tab-seq data or all of them.

1 Input file format

The input format is BED format file with 6 columns (<https://genome.ucsc.edu/FAQ/FAQformat.html#format1>). Below is an example for a line of input file:

```
chr1      3001345 3001346 CpG:9      0.777777777778      +
```

The first three columns are genomic coordinates of a region. The fourth column indicates that this site is a CpG site, and the number of reads covering this site is 9. The fifth column is the methylation level of the CpG site, ranging from 0 to 1. This level is usually calculated from the number of methylated reads covering one site dividing by the total number of reads included in the fourth column.

2 Usage

The usage of mlml is simple: specify the input files with corresponding options to their types. mlml can take various types of input combinations. Assume user has three input files ready: meth_BS-seq.bed, meth_oxBS-seq.bed and meth_Tab-seq.bed, all in the format described above. The following command will take all the inputs:

```
$ ./mlml -u meth_BS-seq.bed -m meth_oxBS-seq.bed \
      -h meth_Tab-seq.bed -o result.bed
```

Note this command should be run in the path where mlml is installed. The option -o will direct the output to file result.bed.

If only two types of input are available, e.g. meth_BS-seq.bed and meth_oxBS-seq.bed, then use the following command:

```
$ ./mlml -u meth_BS-seq.bed -m meth_oxBS-seq.bed \
      -o result.bed
```

In some cases, user might want to specify the convergence tolerance for EM algorithm. This can be done through -t option. For example:

```
$ ./mlml -u meth_BS-seq.bed -m meth_oxBS-seq.bed \  
-o result.bed -t 1e-2
```

This command will make the iteration process stop when the difference of estimation between two iterations is less than 10^{-2} . The value format can be scientific notation, e.g. $1e-5$, or float number, e.g. 0.00001.

3 Output file format

The output of mlml is also in BED format. Here is an example:

```
chr1      3001345 3001346 mC:0.125      0.652778      +
```

The fourth column is estimated 5mC level, and the fifth column is 5hmC level. Note that the output is always in this format, no matter what combination of inputs are used. Therefore user can easily extract 5mC and 5hmC levels from specific columns.

4 FAQ

- Q: I got result like this, what happened?

```
chr1      10570  10571  mC:nan  nan      +
```

A: Because some sites have 0 coverage, there is not enough information to calculate mC and hmC levels. Therefore the output will be nan instead.