

# Persistence of bias in individually normalized chromatin interaction matrices, and its effect on the detection of differential chromatin interactions

John Stansfield, Mikhail Dozmorov

## Contents

<b>Introduction</b>	<b>1</b>
The effect of normalization methods on removing global differences . . . . .	1
<b>The effect of loess normalization on removing chromosome-specific biases</b>	<b>4</b>
Common cutting enzyme . . . . .	4
Different cutting enzymes . . . . .	5
<b>The effect of normalization methods on detecting differential chromatin interactions</b>	<b>6</b>
loess . . . . .	6
ChromoR . . . . .	7
ICE . . . . .	8
KR . . . . .	9
SCN . . . . .	10
<b>Comparison of HiCcompare vs. ChromoR in detecting differential chromatin interactions</b>	<b>11</b>

## Introduction

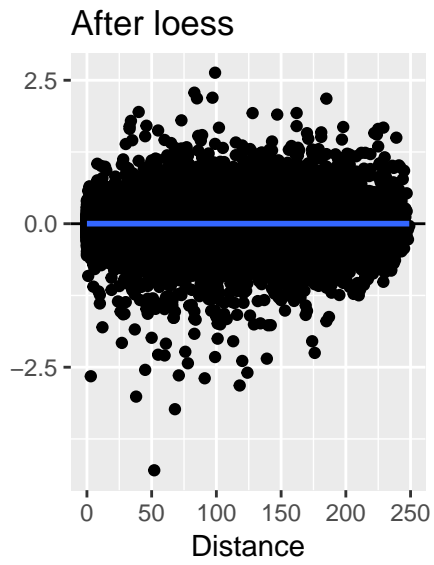
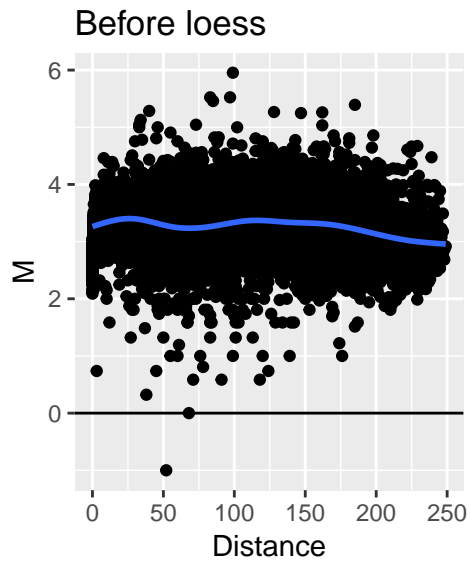
To compare the ability of methods for normalizing individual datasets to remove biases *between* chromatin interaction matrices, we compare individually normalized matrices with the jointly normalized ones. Several parameters were assessed:

- The effect of global differences. Most of the time different chromatin interaction matrices will contain different total number of reads, resulting in the overall differences. We assessed whether methods for normalizing individual datasets were able to account for the differences in the total number of reads.
- The ability of the joint normalization to account for biases under different conditions, such as when comparing matrices obtained with different cutting enzymes, or matrices from different chromosomes.
- The effect of individual and joint normalization methods on detecting chromatin interaction differences.

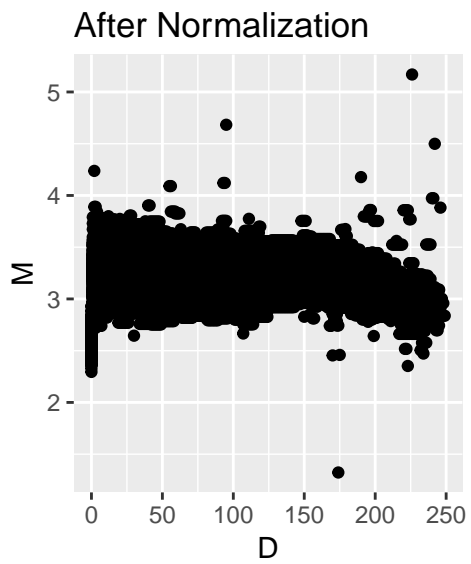
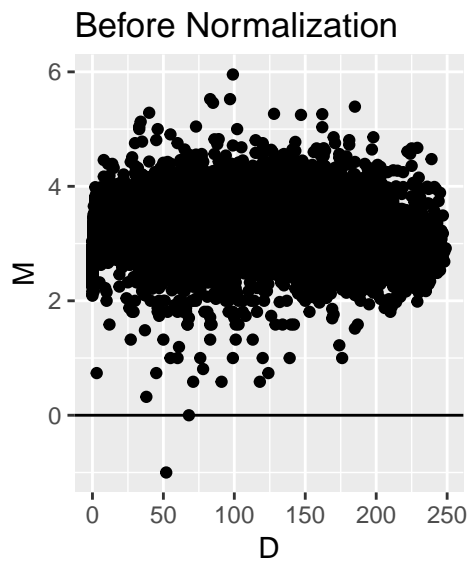
### The effect of normalization methods on removing global differences

Hi-C matrices may have a different total number of reads. This disbalance will lead to the overall difference between the two matrices, reflected by the global shift of the cloud of  $M$  differences from zero. The unscaled matrices, globally shifted from  $M = 0$ , can be successfully normalized by `loess`. However, individually normalized matrices will still contain the global shift, as shown below. By default, the `create.hic.table` function rescales the matrices to have the same total number of reads. Rescaling is accomplished by first calculating the scale factor  $\psi = \frac{\sum IF_i}{\sum IF_j}$  where  $i$  is the set of all the IFs for the upper triangle of the first Hi-C matrix and  $j$  is the set of all IFs for the upper triangle of the second Hi-C matrix. Next,  $IF_j$  is scaled by setting  $IF_{j_{new}} = \frac{IF_j}{\psi}$ .

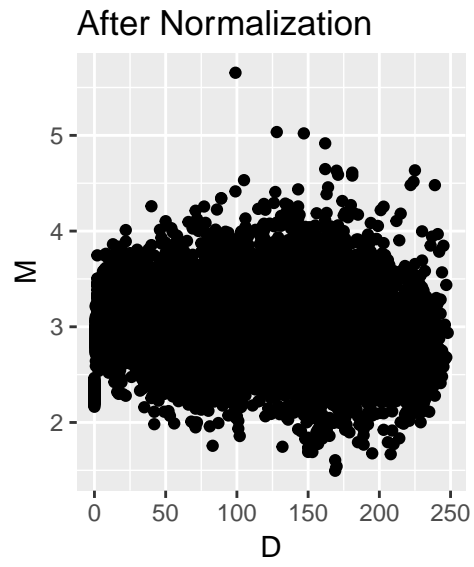
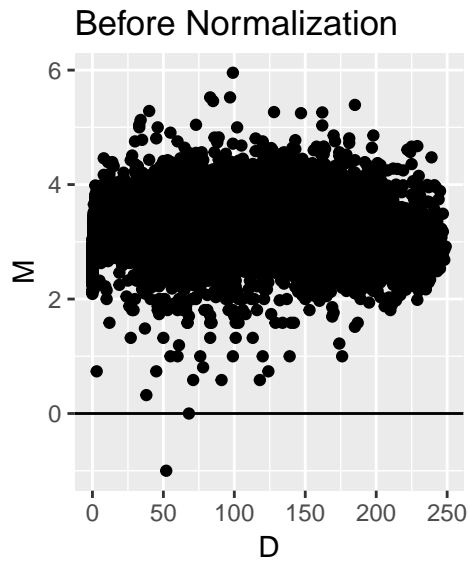
loess



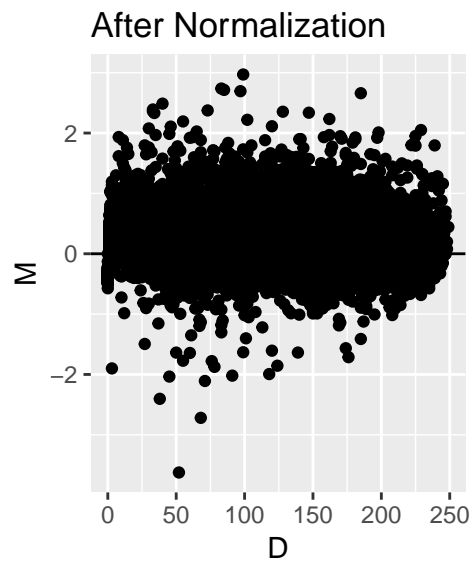
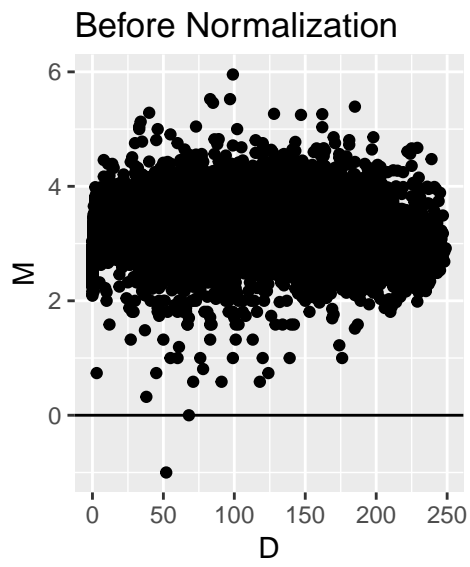
ChromoR



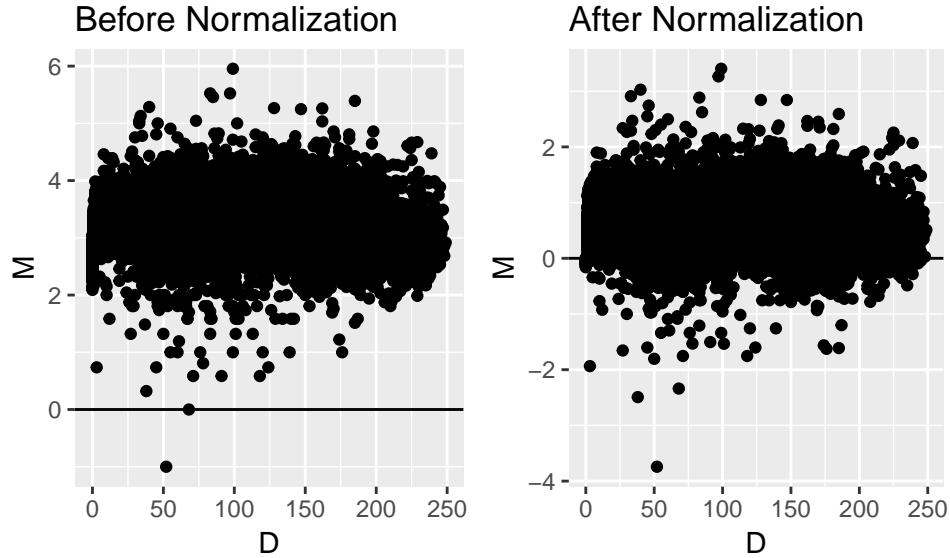
ICE



KR



SCN



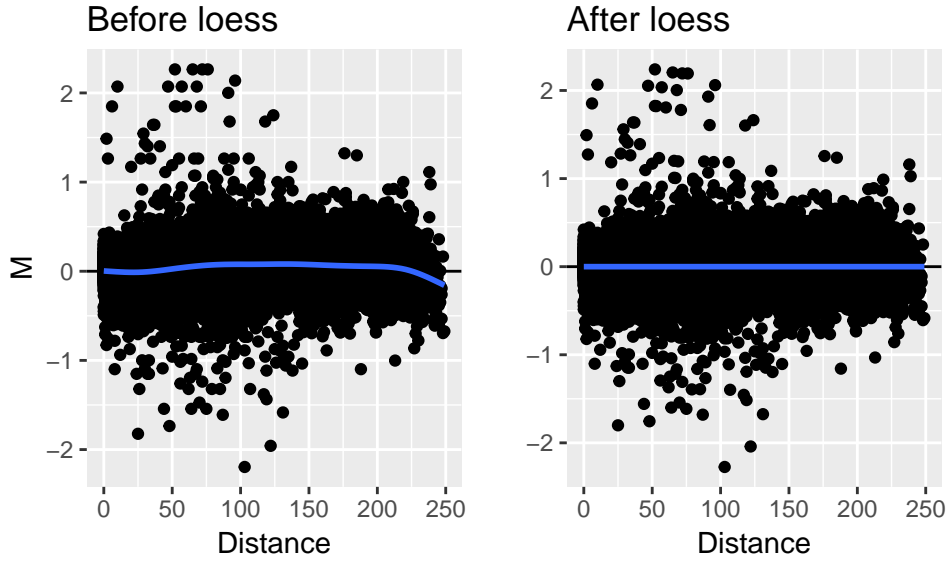
### Summary

As can be seen from the above, the MD plots for the single matrix normalization methods do not all succeed at rescaling the data and thus the main cloud of points are not centered around  $M = 0$ . **Loess** however, was able to take care of rescaling the data and centered the MD plot around 0. Global scaling is recommended for any comparisons between Hi-C datasets.

## The effect of loess normalization on removing chromosome-specific biases

### Common cutting enzyme

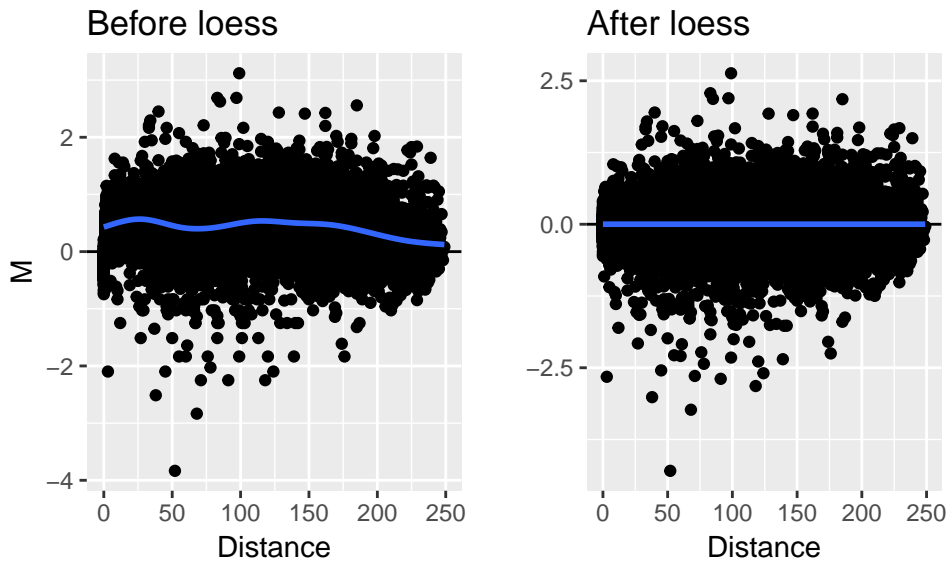
The MD plot below displays data before and after joint **loess** normalization from GM12878 at 1MB resolution, chromosome 1, that were obtained as replicates using the same cutting enzyme. Since the data here is replicate data it is expected that there will not be many differences between the datasets. Any differences found are assumed to be due to bias in the sequencing procedures.



As can be seen by the loess fit on the “Before loess” MD plot there is not a large amount of bias between the two datasets.

### Different cutting enzymes

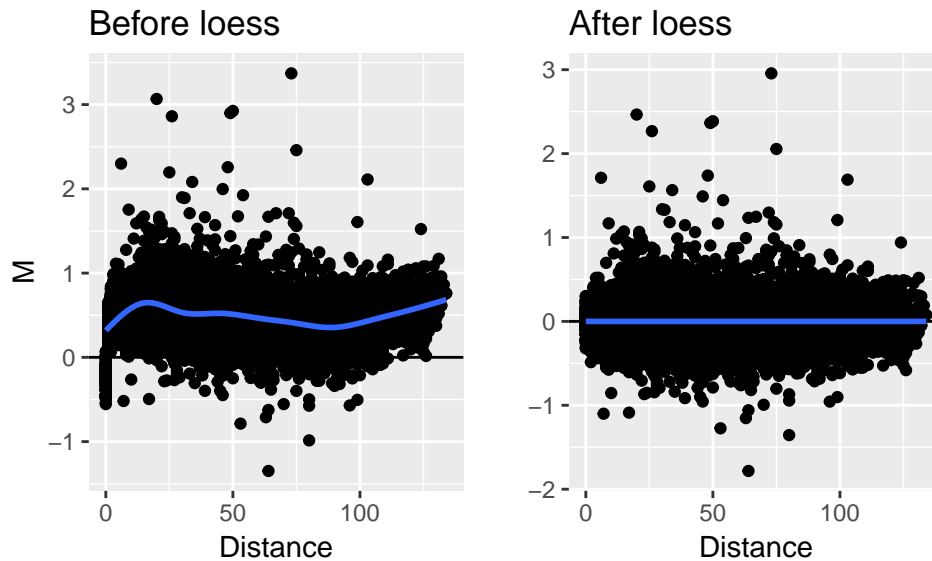
The Hi-C datasets here are from GM12878 cell lines at 1MB resolution, chromosome 1. One dataset was cut using MboI and the other using DpnII. Since different cutting enzymes are used it is expected that there will be some differences in the data due to enzyme choice. Biases between the datasets are successfully removed with `loess` normalization as can be seen in the following MD plots displaying the data before and after joint `loess` normalization. It can also be seen that biases between the datasets differ between each chromosome and dataset. The differences do not appear to follow a trend which makes a non-parametric approach to normalization better suited to the task.



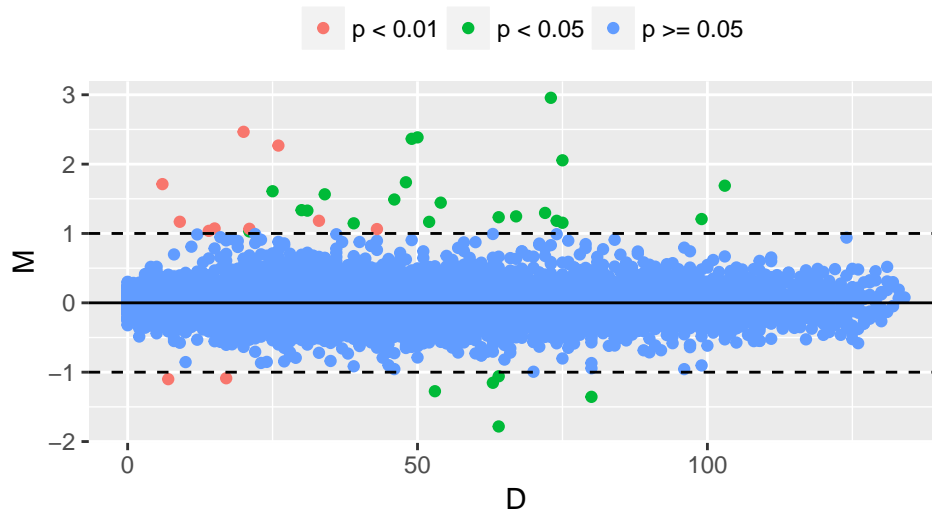
# The effect of normalization methods on detecting differential chromatin interactions

To look at differences between different normalization methods we use data from GM12878 at 1MB resolution on chr 11 generated using two different cutting enzymes, MboI and DpnII. The data is unscaled. For each method tested below, we also test for differences between the two datasets. No artificial changes were added to the datasets. Any differences detected by the method will be examples of existing differences between the replicated Hi-C data on the same cell line when cut by different enzymes. Since the datasets are for the same chromosome and the same cell line we should expect few differences to be detected.

## loess



## MD Plot

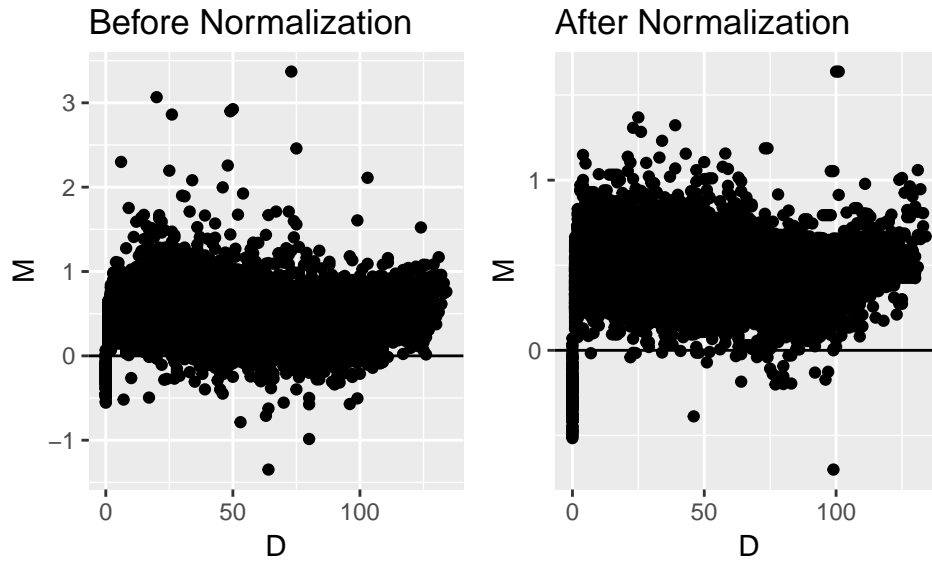


[1] "37 differences found between the datasets"

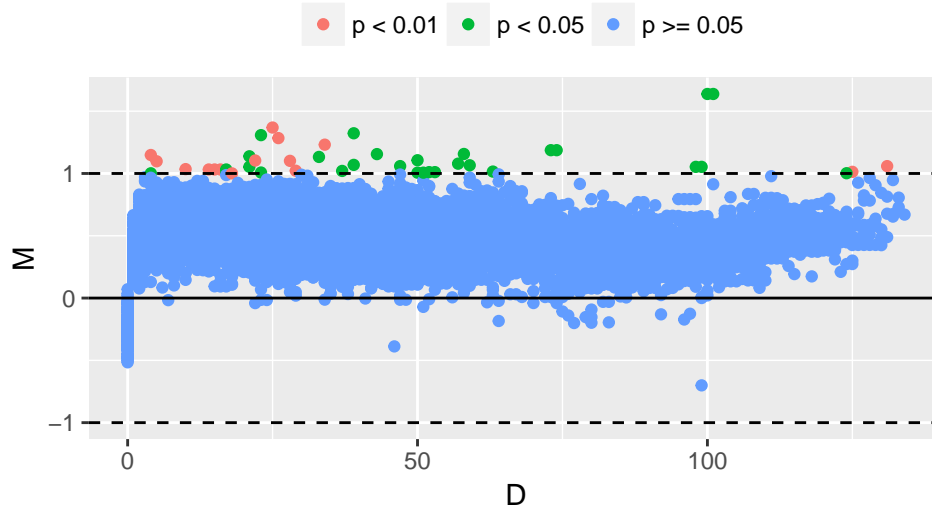
The MD plot above serves as a reference to show that Loess can successfully normalize the data and

removes bias between the two datasets. The following MD plots display the data after the specified individual normalization method has been applied to each matrix. The dotted lines on the second MD plot represent the difference threshold where points are considered significant outside of the two dotted lines. The difference threshold used here is  $\theta = 1$ , see `?hic_diff` and the methods section of the paper for a full explanation on the difference threshold.

### ChromoR

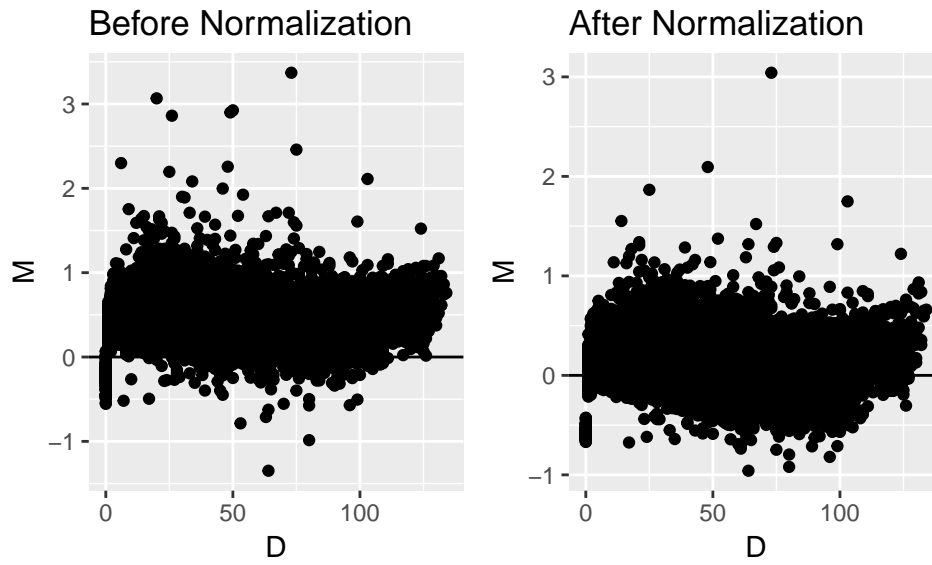


### MD Plot

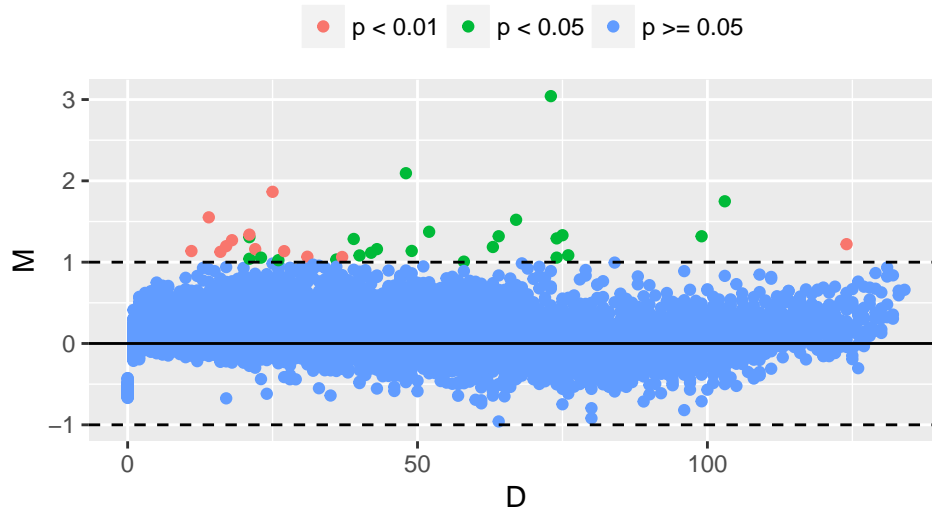


[1] "44 differences found between the datasets"

ICE



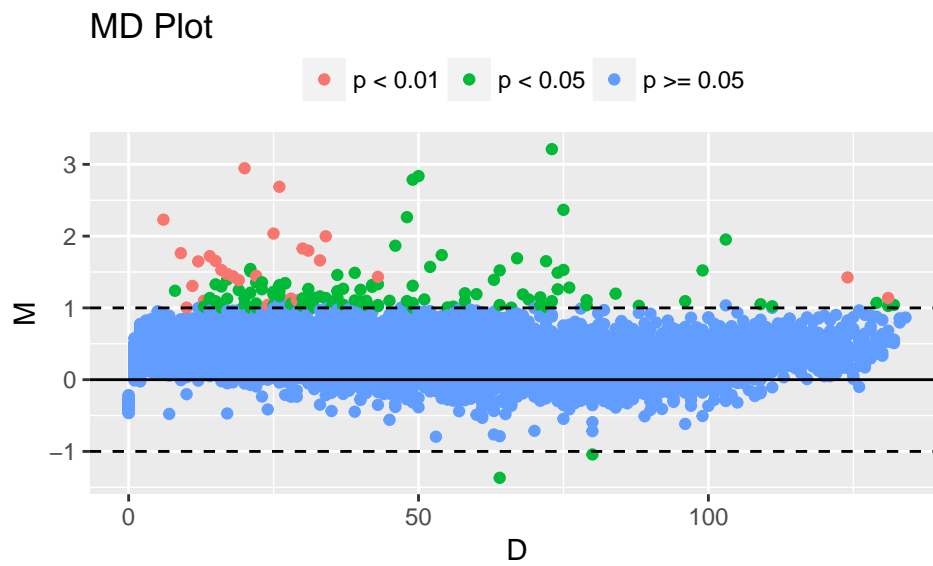
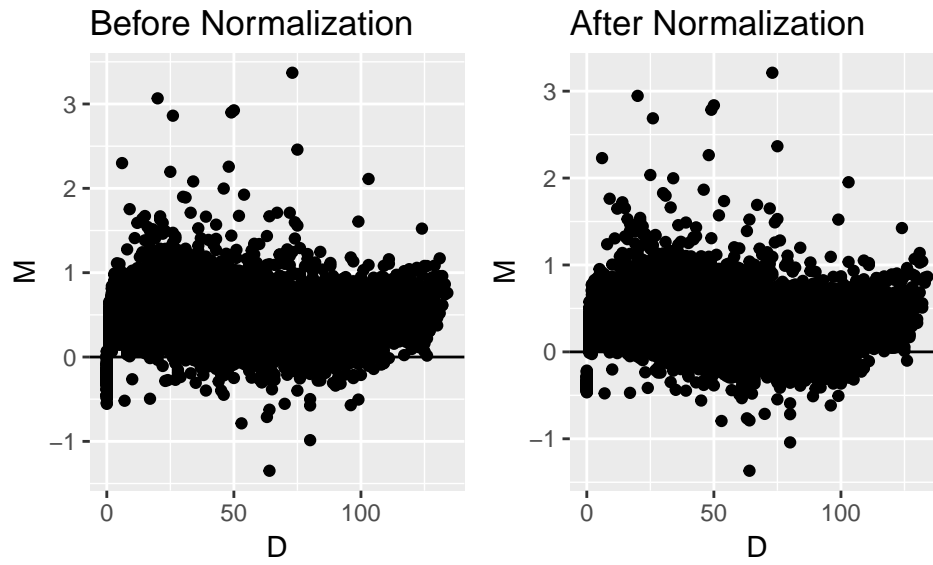
MD Plot



[1] "35 differences found between the datasets"

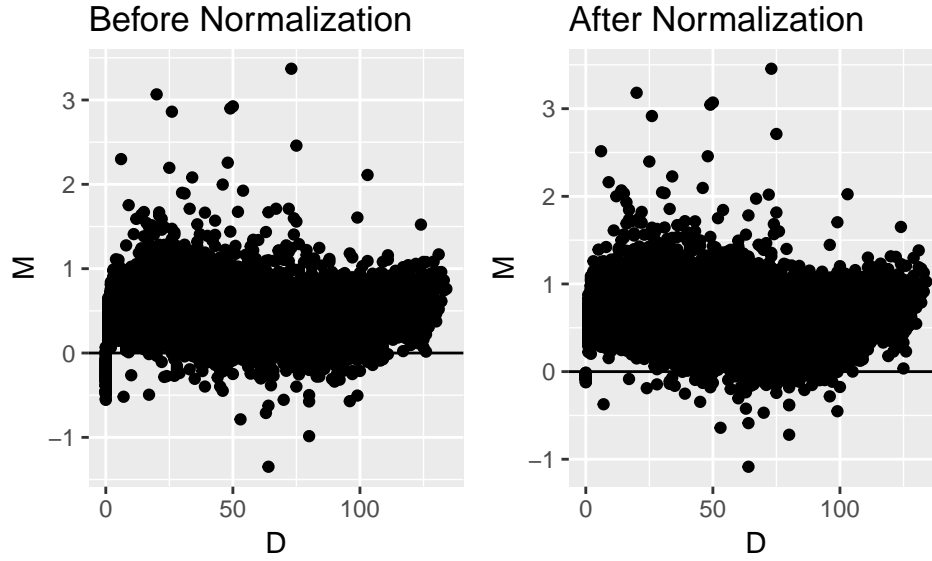


KR

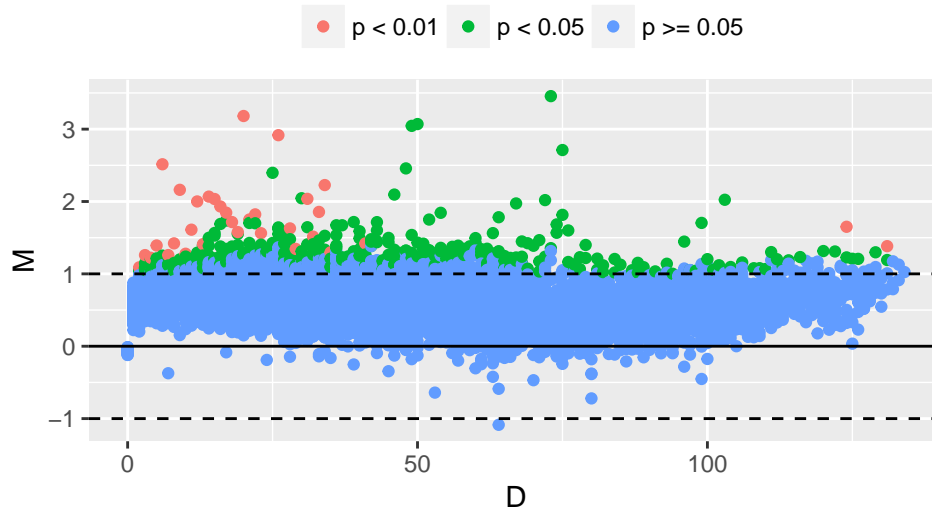


[1] "144 differences found between the datasets"

SCN



MD Plot



[1] "350 differences found between the datasets"

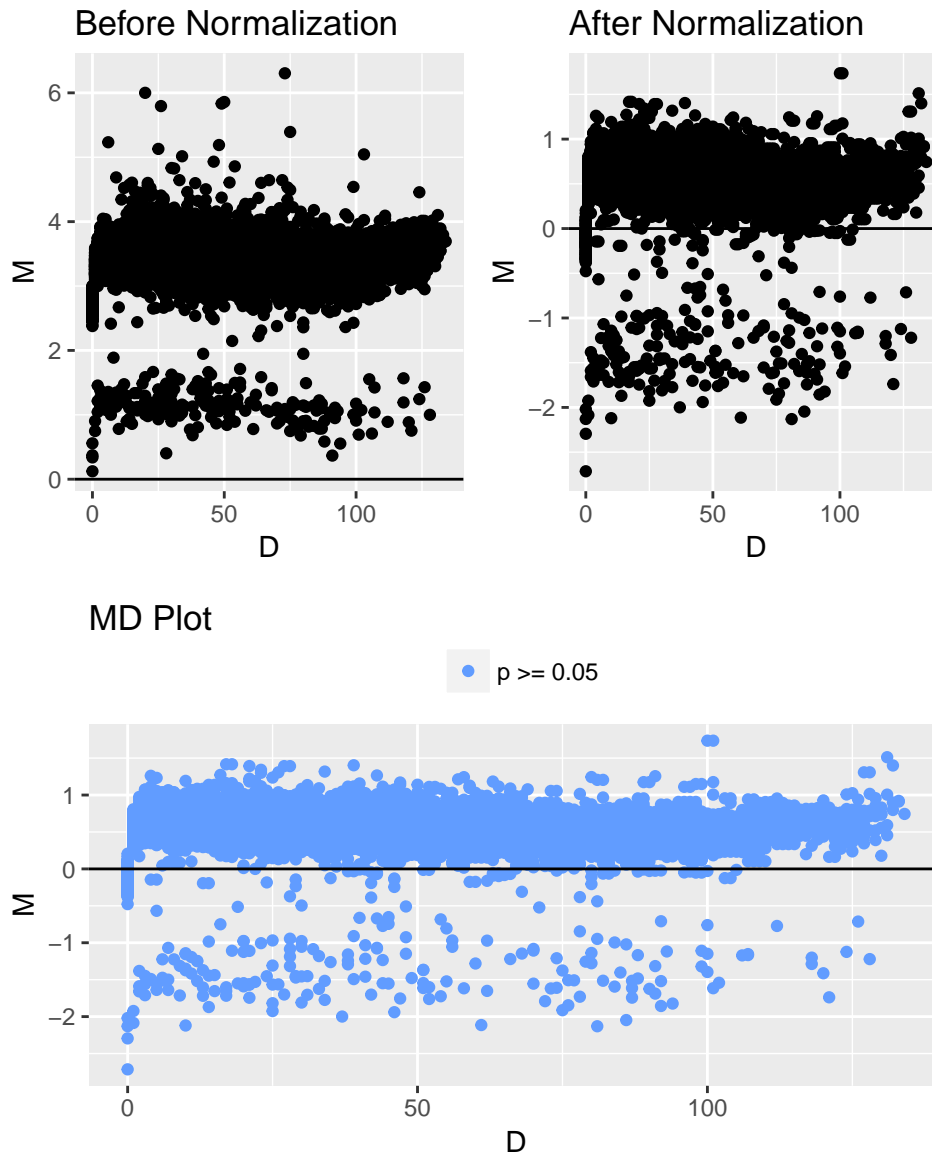
### Summary

Loess is the only method that can successfully remove the bias between the two datasets. The individual normalization techniques fail to remove biases between the datasets though they may be effective at removing bias within a single dataset. KR normalization appears to be second to the loess normalization in removing global and local biases.

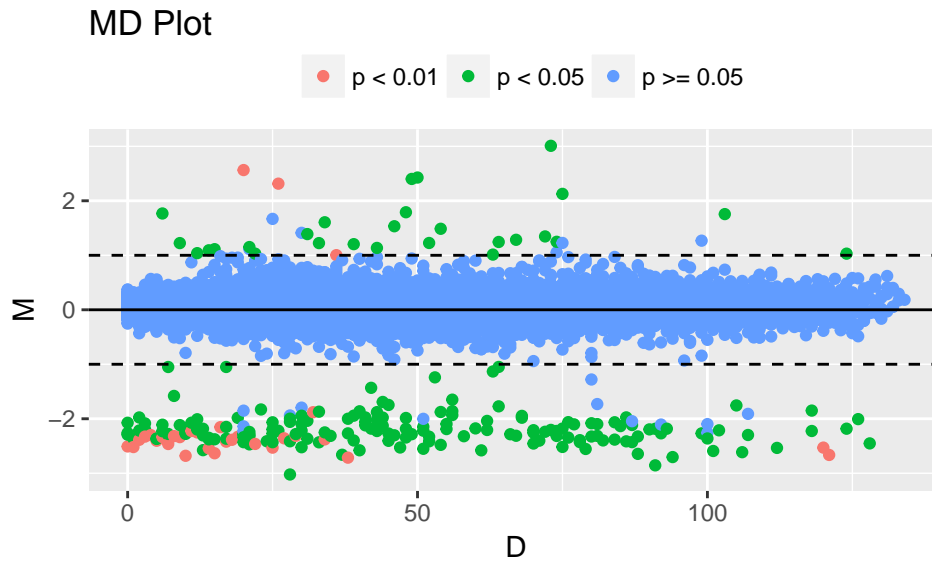
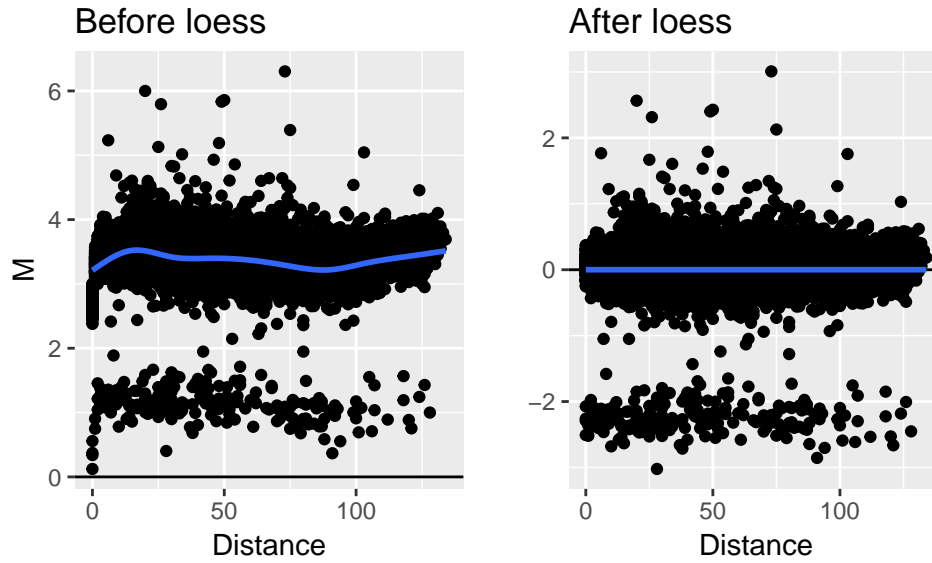
## Comparison of HiCcompare vs. ChromoR in detecting differential chromatin interactions

ChromoR includes a function for detecting differences between two Hi-C datasets. Using the data for chromosome 11 from GM12878 as used in the above normalization comparison we add 200 *a priori* known differences to the matrix at a 5 fold change and attempt to detect them using ChromoR and loess

The MD plot of the ChromoR normalized matrices:



ChromoR found 0 differences between the two matrices. Compared to `hic_loess` below:



Loess found 224 differences between the matrices.

ChromoR's normalization technique fails to remove bias between Hi-C datasets and its difference detection method also fails to detect any differences when true differences are added at a 5 fold change. Loess was capable of normalizing these datasets and detecting the majority of the true differences added to the matrices.