

# Example

Jinlu Liu

15/08/2022

We use the following example (same as Simulation 1 in the main article) to demonstrate the use of functions in the normHDP package. Suppose the dataset we have is named `Y_linear`, to run a short markov chain, we use a burn-in of 1,000, thinning of 5 and total iteration of 3,000:

```
# Load Y_linear from data file
load('data/Y_linear.RData')

# Obtain posterior samples
case1_mcmc <- normHDP_mcmc(Y = Y_linear,
                          J = 3,
                          number_iter = 3000,
                          thinning = 5,
                          empirical = TRUE,
                          burn_in = 1000,
                          quadratic = FALSE,
                          iter_update = 100,
                          beta.mean = 0.72,
                          alpha_mu_2 = 20,
                          print_Z = TRUE,
                          num.cores = 4)
```

If we want to continue to run for more iterations, i.e.: to extend the total number of iterations from 3,000 to 5,000:

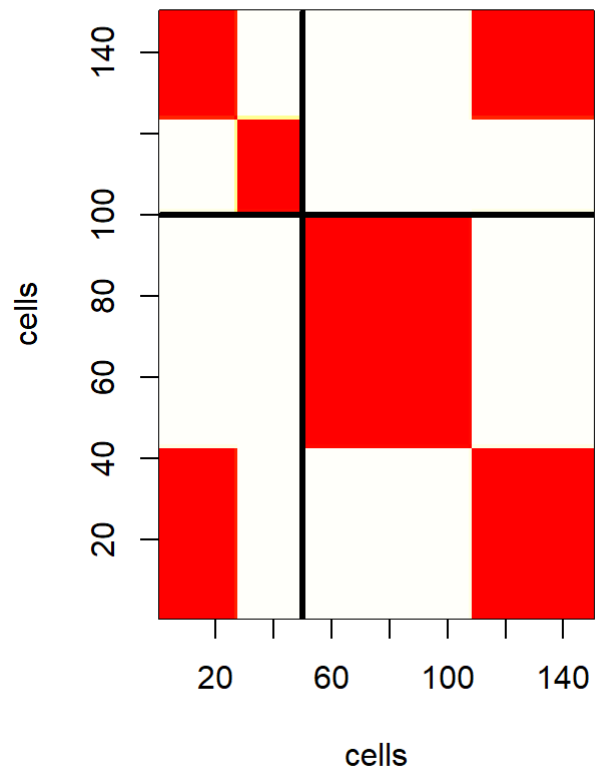
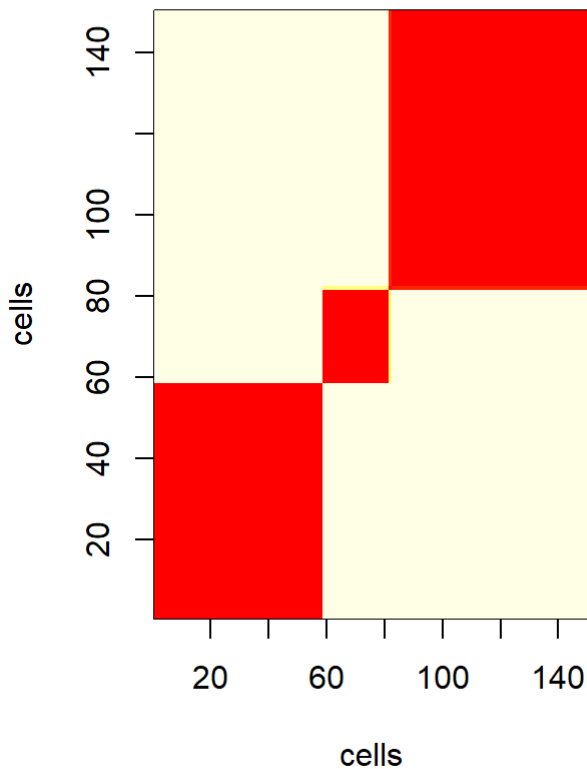
```
normHDP_mcmc_2(normHDP_output = case1_mcmc,
               number_iter = 10000)
```

To obtain the posterior similarity matrix:

```
case1_psm <- similarity_matrix(normHDP_output = case1_mcmc)
```

The output from `similarity_matrix()` contains a list of 2 items: the first item is the posterior similarity matrix between cells from 2 datasets or within the same dataset. The second item is the posterior similarity matrix between cells from all datasets. In plot below, the left sub-plot shows posterior similarity by combining cells from all datasets together; the right sub-plot shows posterior similarity between cells from the same dataset and between (any of the) 2 datasets.

```
plotpsm(psm.ind = case1_psm$psm.within,
        psm.tot = case1_psm$psm.combined,
        method="complete",
        xlab = 'cells',
        ylab = 'cells')
```



Next, we can summarize the posterior estimates of clustering using a single point estimate:

```
# Optimal Clustering
case1_Z_estimate <- opt.clust(normHDP_output = case1_mcmc,
                             psm_output = case1_psm)

# cluster Summary
cluster_summary(Z_point = case1_Z_estimate)
```

Since we have the label-switching problem for the previous MCMC, we cannot carry out analysis on the unique parameters (mean expressions and dispersion) by simply taking the posterior mean or median, instead, we analyze the posterior estimates of mean expression and dispersion by running a separate MCMC chain for the unique parameters based on the posterior estimates of allocations and capture efficiencies:

```

# MCMC of mu and phi
case1_mcmc_post <- normHDP_mcmc_post(normHDP_output = case1_mcmc,
                                   burn_in = 1000,
                                   thinning = 5,
                                   number_iter = 3000,
                                   Z = case1_Z_estimate,
                                   Y = Y_linear,
                                   iter_update = 100)

# Chain Length
length.mcmc <- length(case1_mcmc_post$mu_star_1_J_output)

# Posterior mean of mean expression with fixed allocation
case1.post.mean.mu <- Reduce("+", case1_mcmc_post$mu_star_1_J_output)/length.mcmc

# Posterior mean of dispersion with fixed allocation
case1.post.mean.phi <- Reduce("+", case1_mcmc_post$phi_star_1_J_output)/length.mcmc

# Posterior mean of capture efficiencies
case1.post.mean.beta <- case1_mcmc_post$Beta_posterior_mean

```

Similar to the `normHDP_mcmc` function(), we can extend the total number of iterations using:

```

normHDP_mcmc_post2(normHDP_post_output = case1_mcmc_post,
                   number_iter = 5000)

```

The structure of the output from `normHDP_mcmc_post2()` and `normHDP_mcmc_post()` are the same.

## Global Marker genes

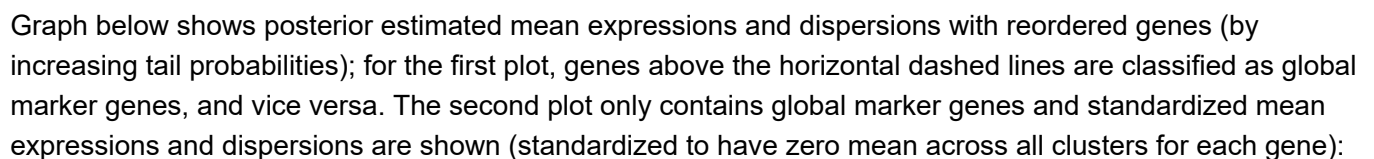
To find global marker genes, we use the function `global_marker_genes()` where we use the output from `normHDP_post_output()` as an input and pre-specify a set of threshold for determining global marker genes, example below uses 1.5 for both mean expressions and dispersions. We can also pre-specify `alpha_M` and `alpha_D` (example below). By default, `alpha_M` and `alpha_D` are set to control the expected false discovery rate (EFDR) to 5 percent. The output of the function includes four items:

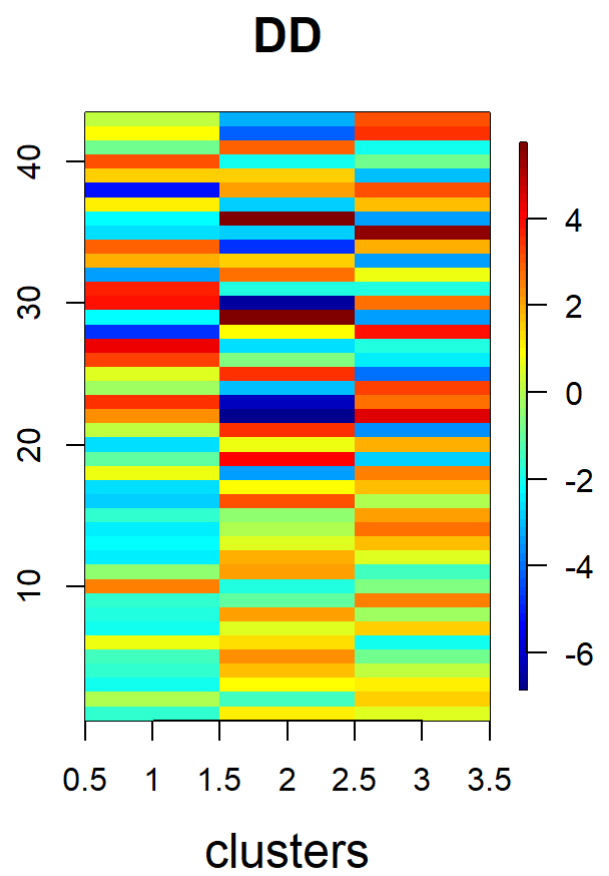
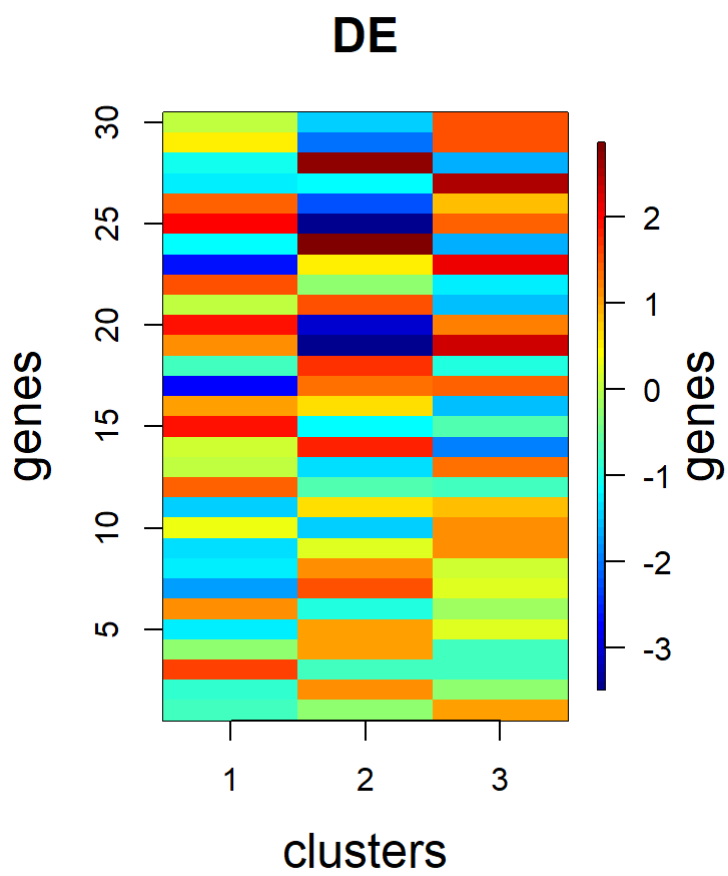
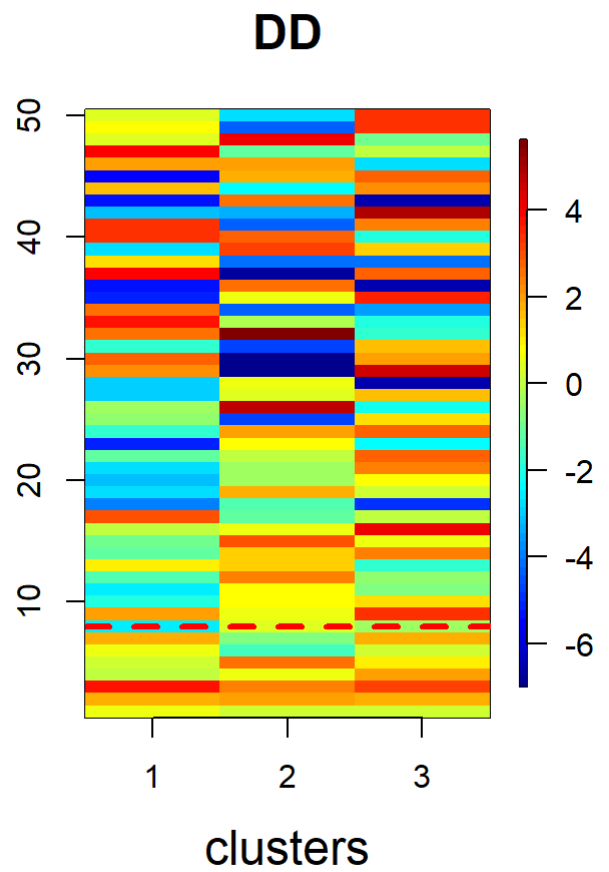
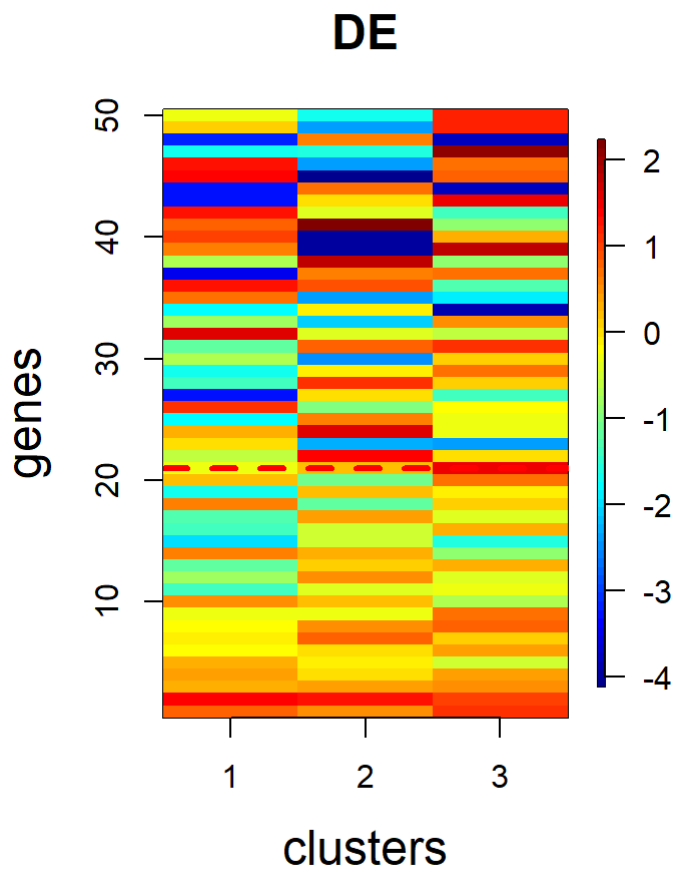
- `marker_DE`: a data frame with four variables: the gene index, the absolute log-fold change ( $\mu$ ) of the gene, tail-probability ( $\mu$ ) of the gene, a binary variable to indicate whether the gene is globally differentially expressed (DE).
- `marker_DD`: a data frame with four variables: the gene index, the absolute log-fold change ( $\phi$ ) of the gene, tail-probability ( $\phi$ ) of the gene, a binary variable to indicate whether the gene is globally differentially dispersed (DD).
- `alpha_M`: threshold for the tail probabilities of  $\mu$ . Genes with tail probabilities greater than `alpha_M` are classified as global DE genes.
- `alpha_D`: threshold for the tail probabilities of  $\phi$ . Genes with tail probabilities greater than `alpha_D` are classified as global DD genes.

The summary plot contains three sub-figures. The sub-figure on the left compares the tail probability against the mean absolute log-fold change for the mean expressions. The sub-figure in the middle compares the tail probability against the mean absolute log-fold change for dispersion. The sub-figure on the right summarizes the total number of genes with these global characteristics.

```
case1_global_marker <- global_marker_genes(normHDP_post_output = case1_mcmc_post,
      threshold = list('mu' = 1.5, 'phi' = 1.5),
      alpha_M = 0.9,
      alpha_D = 0.9)
```

```
global_marker_genes_plot(gmg_output = case1_global_marker)
```

[illegible]



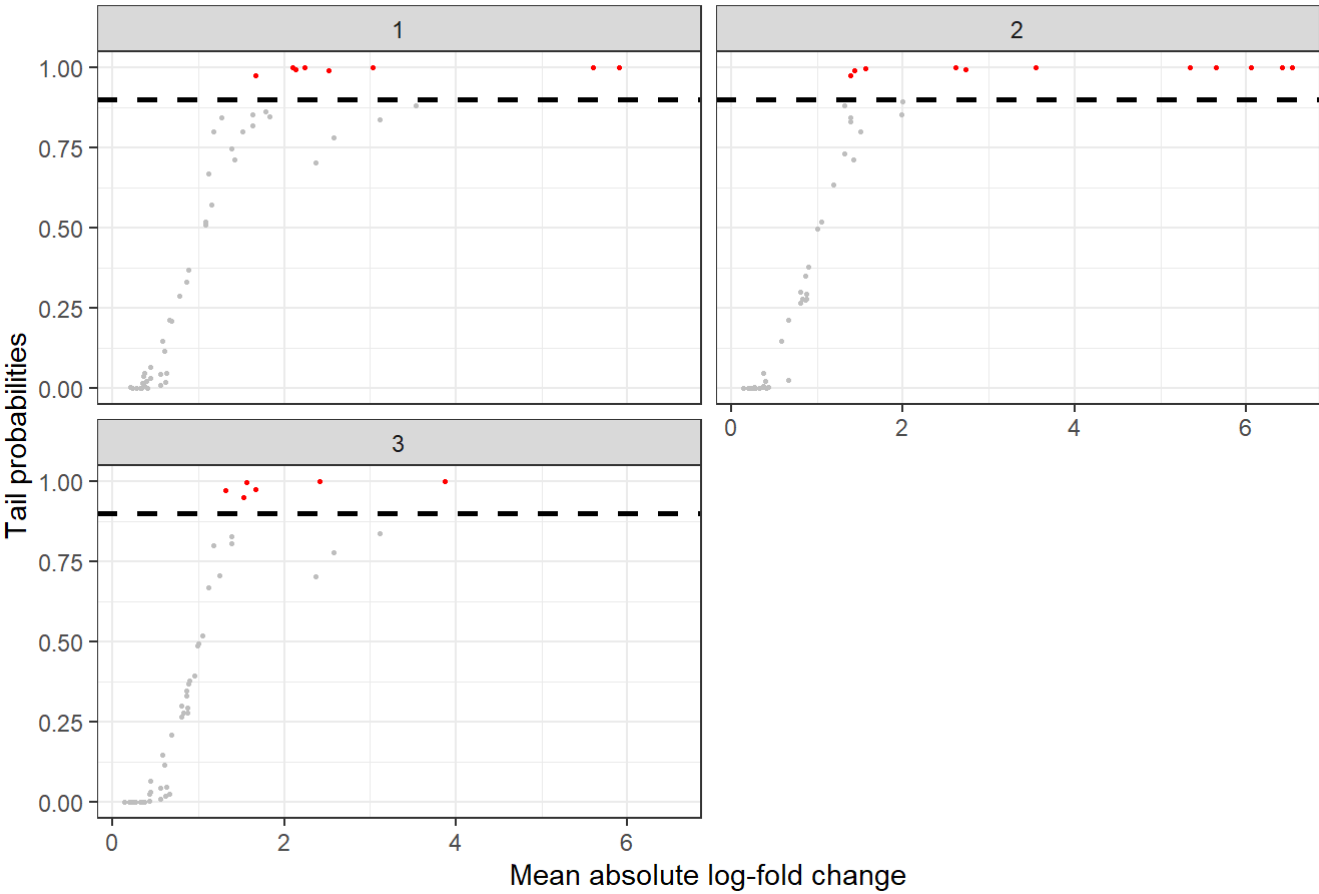
Local marker genes

To find local marker genes, we use the function `local_marker_genes()` where the input structure is the same as for function `global_marker_genes()`. It also returns 4 items as output, but the structure of `marker_DE` and `marker_DD` are slightly different; both of these data frames contain one extra variable to label the cluster. Note that when referring to local marker genes, we need to specify the cluster as well.

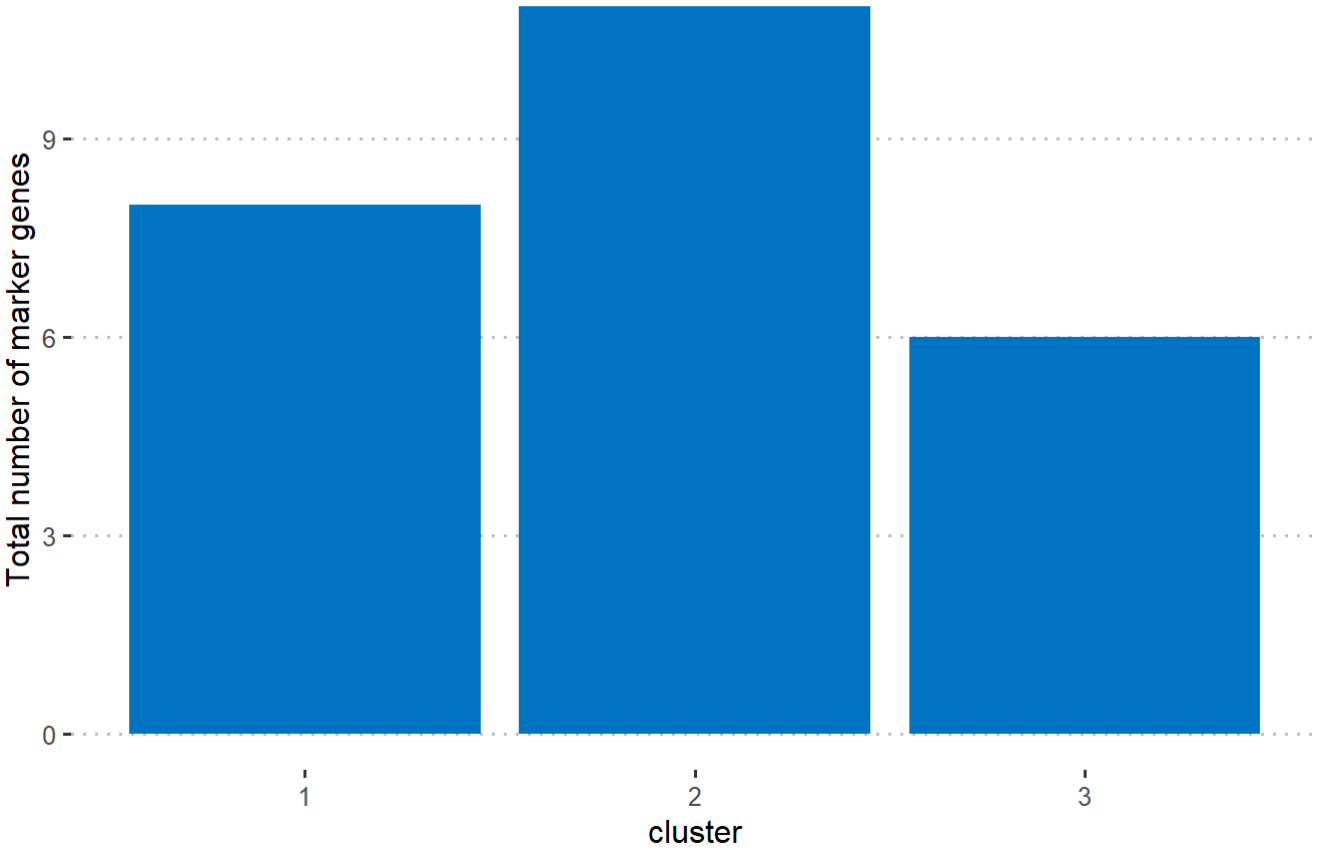
```
case1_local_marker <- local_marker_genes(normHDP_post_output = case1_mcmc_post,  
                                         threshold = list('mu' = 1, 'phi' = 1),  
                                         alpha_M = 0.9,  
                                         alpha_D = 0.9)
```

```
# Relationship between absolute log-fold change and tail probabilities  
local_marker_genes_plot(lmg_output = case1_local_marker)
```

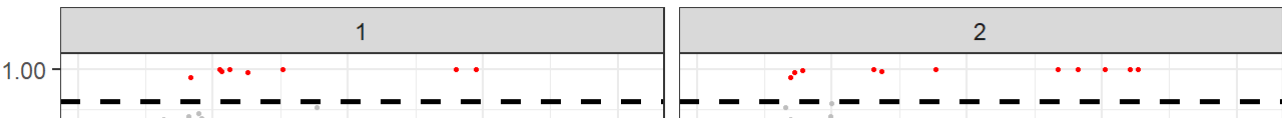
Local Marker Genes (DE)

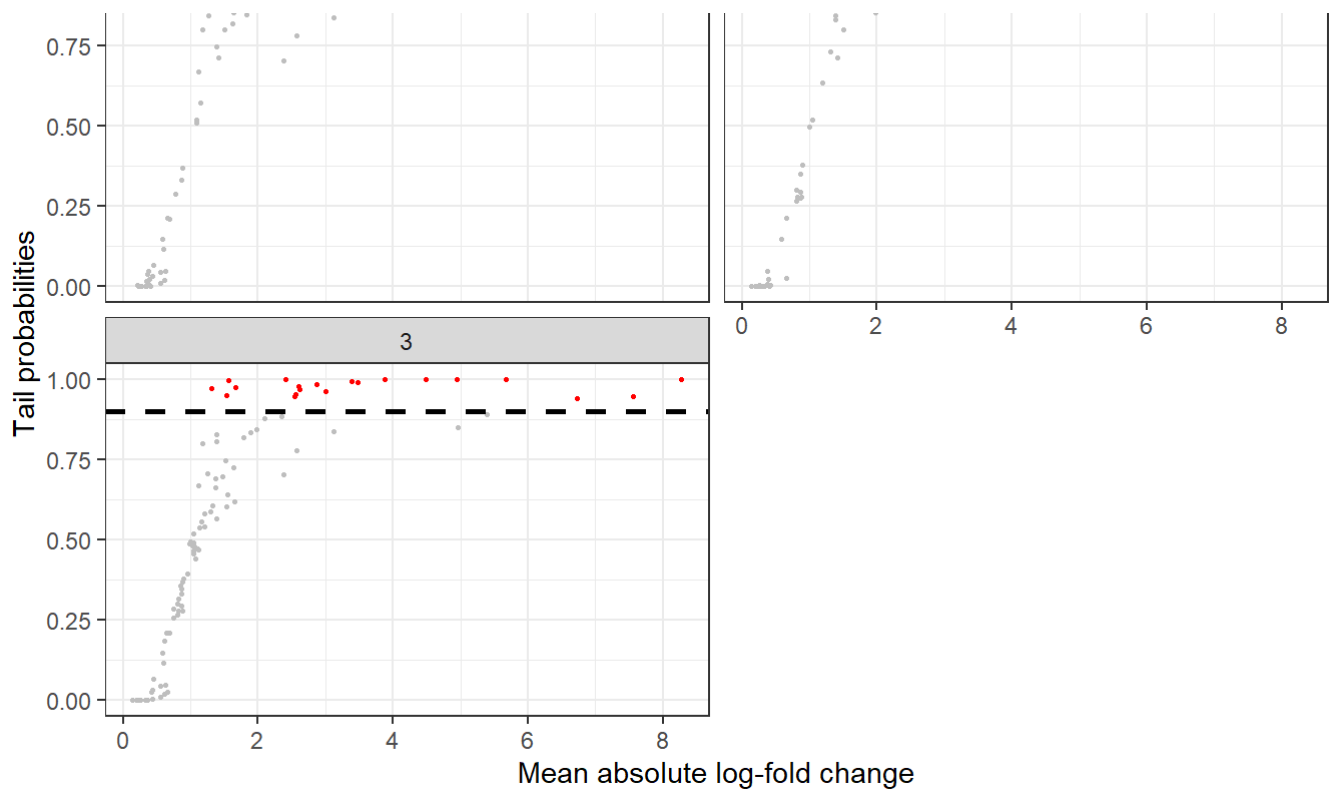


Local Marker Genes (DE)

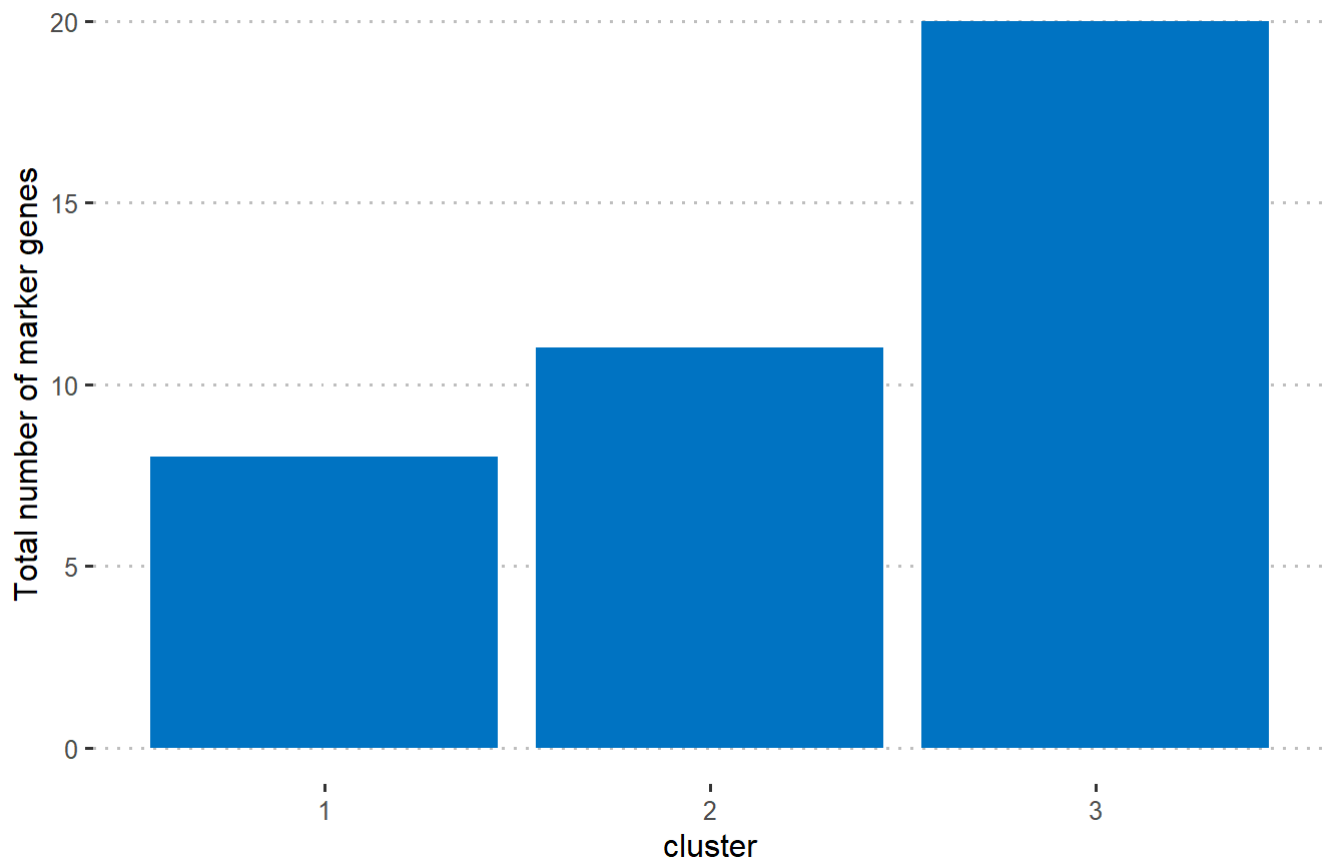


Local Marker Genes (DD)



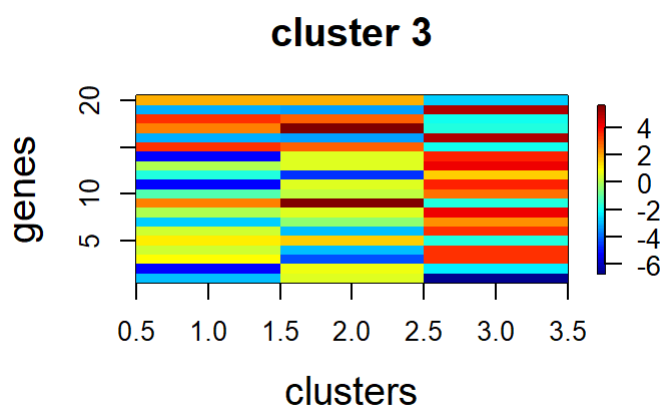
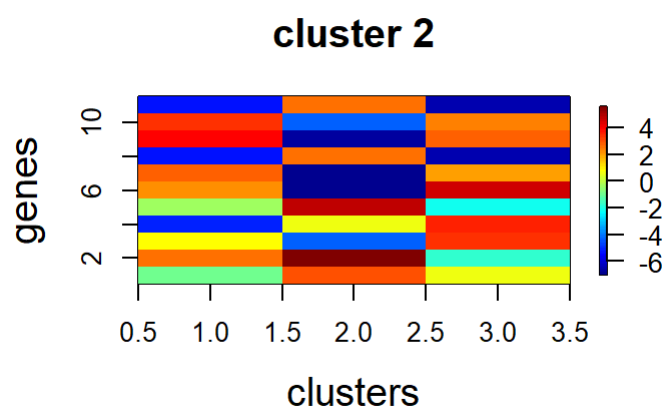
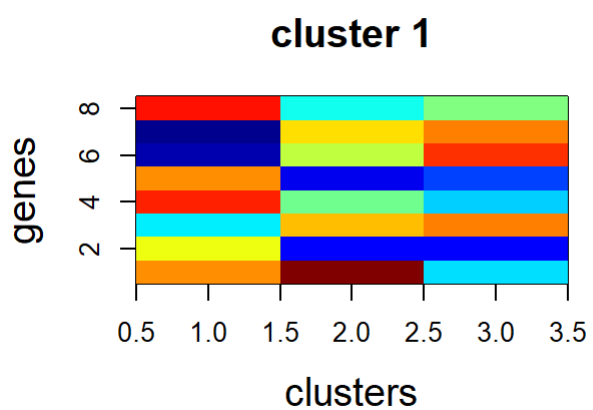
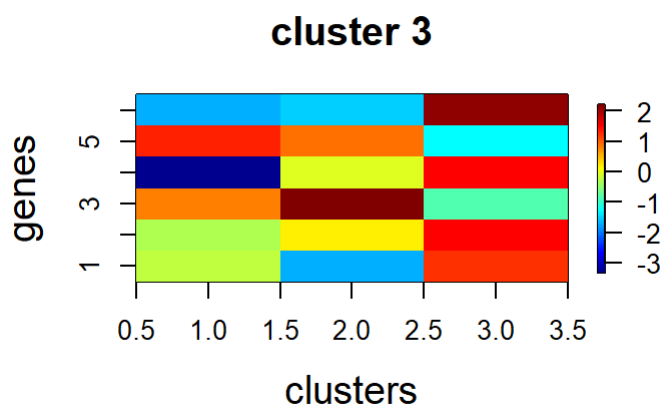
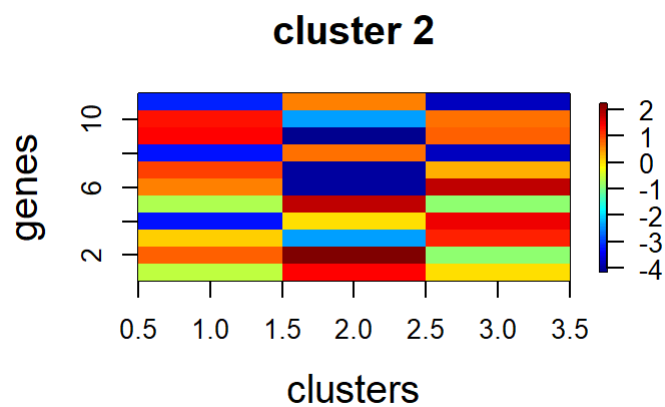
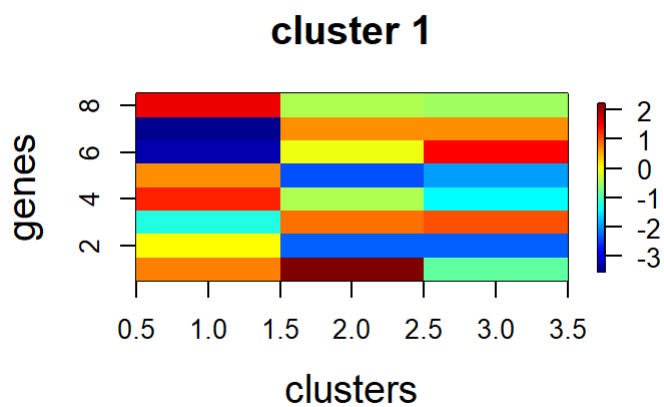


### Local Marker Genes (DD)



```
# Heat map
local_marker_genes_heatmap(lmg_output = case1_local_marker,
                           normHDP_post_output = case1_mcmc_post)
```





Latent counts and observed counts

We use the `latent_counts()` function to calculate and reorder latent counts; By default, only cells are reordered unless output from `gmg_output()` is provided. Specifically, we have a list of 3 items:

- `matrix.output`: a list of D items, same structure as input Y.
- `index.solid`: index to plot on the x-axis to separate cells from different clusters.
- `index.dashed`: index to plot on the x-axis to separate cells from different datasets.

If output from `gmg_output()` is provided, then the output will be a list with 6 items:

- `Y_latent_DE`: a list of D items, same structure as input Y. Cells are reordered by clustering and genes are reordered by the tail probabilities corresponding to mean expressions for each dataset.
- `Y_latent_DD`: a list of D items, same structure as input Y. Cells are reordered by clustering and genes are reordered by the tail probabilities corresponding to dispersions for each dataset.
- `DE_number`: Number of global differentially expressed genes.
- `DD_number`: Number of global differentially dispersed genes.
- `index.solid` and `index.dashed` are also provided.

```
# Estimated latent counts
case1_estimated_latent <- latent_counts(mu_JG = case1.post.mean.mu,
                                       phi_JG = case1.post.mean.phi,
                                       Z_CD = case1_Z_estimate,
                                       beta_CD = case1.post.mean.beta,
                                       Y = Y_linear,
                                       gmg_output = case1_global_marker)
```

We can use the `observed_counts()` function to reorder cells and genes of the observed counts; by default, the function only reorder cells unless output from `gmg_output()` is provided. Under default, the `observed_counts()` function outputs a list of 3 items:

- `matrix.output`: a list of D items, same structure as input Y.
- `index.solid`: index to plot on the x-axis to separate cells from different clusters.
- `index.dashed`: index to plot on the x-axis to separate cells from different datasets.

If output from `gmg_output()` is provided, then the `observed_counts()` function outputs a list of 6 items:

- `Y_DE`: a list of D items, same structure as input Y. Cells are reordered by clustering and genes are reordered by the tail probabilities corresponding to mean expressions.
- `Y_DD`: a list of D items, same structure as input Y. Cells are reordered by clustering and genes are reordered by the tail probabilities corresponding to dispersions.
- `DE_number`: Number of global differentially expressed genes.
- `DD_number`: Number of global differentially dispersed genes.
- `index.solid` and `index.dashed` are also provided.

```
# Observed counts
case1_observed_counts <- observed_counts(Z_CD = case1_Z_estimate,
                                       Y = Y_linear,
                                       gmg_output = case1_global_marker)
```

## Posterior Predictive Checks

Section below demonstrate the use of posterior predictive checks (ppc). We repeat steps below twice to show variation. For each posterior predictive check below, we compare 2 relationships:

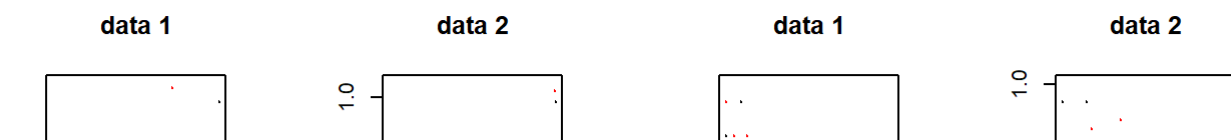
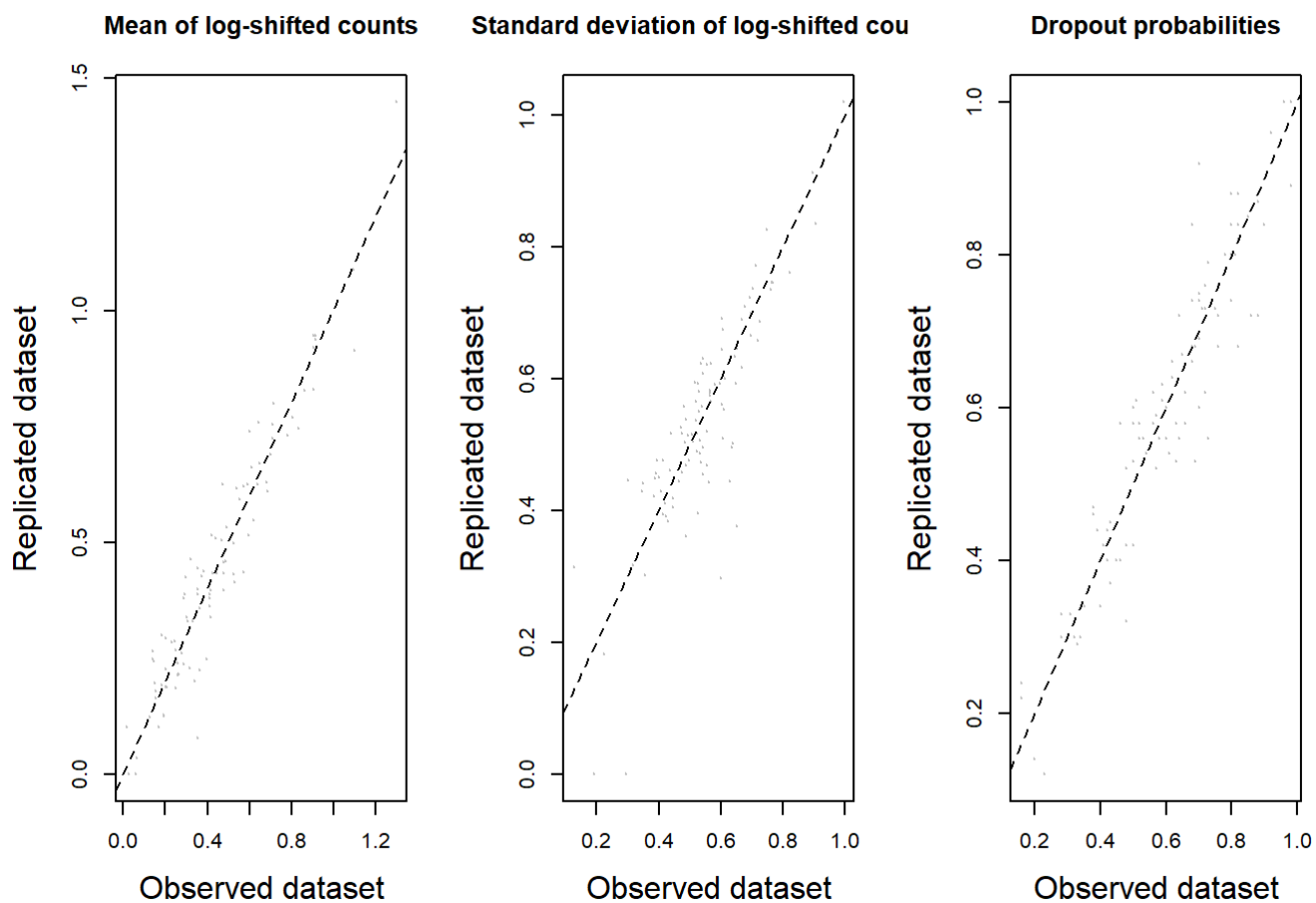
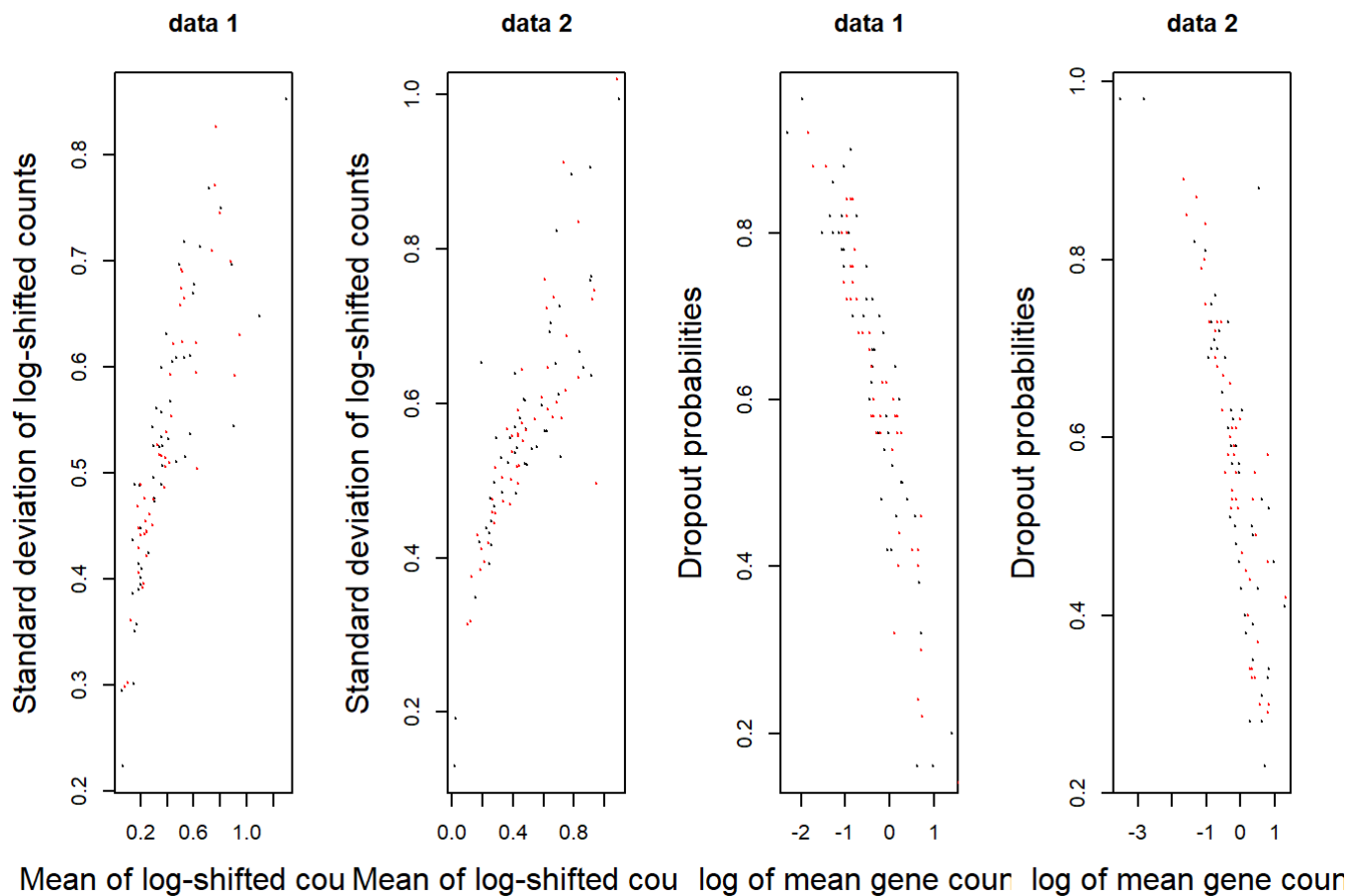
- Relationship between mean of log-shifted counts and standard deviation of log-shifted counts.

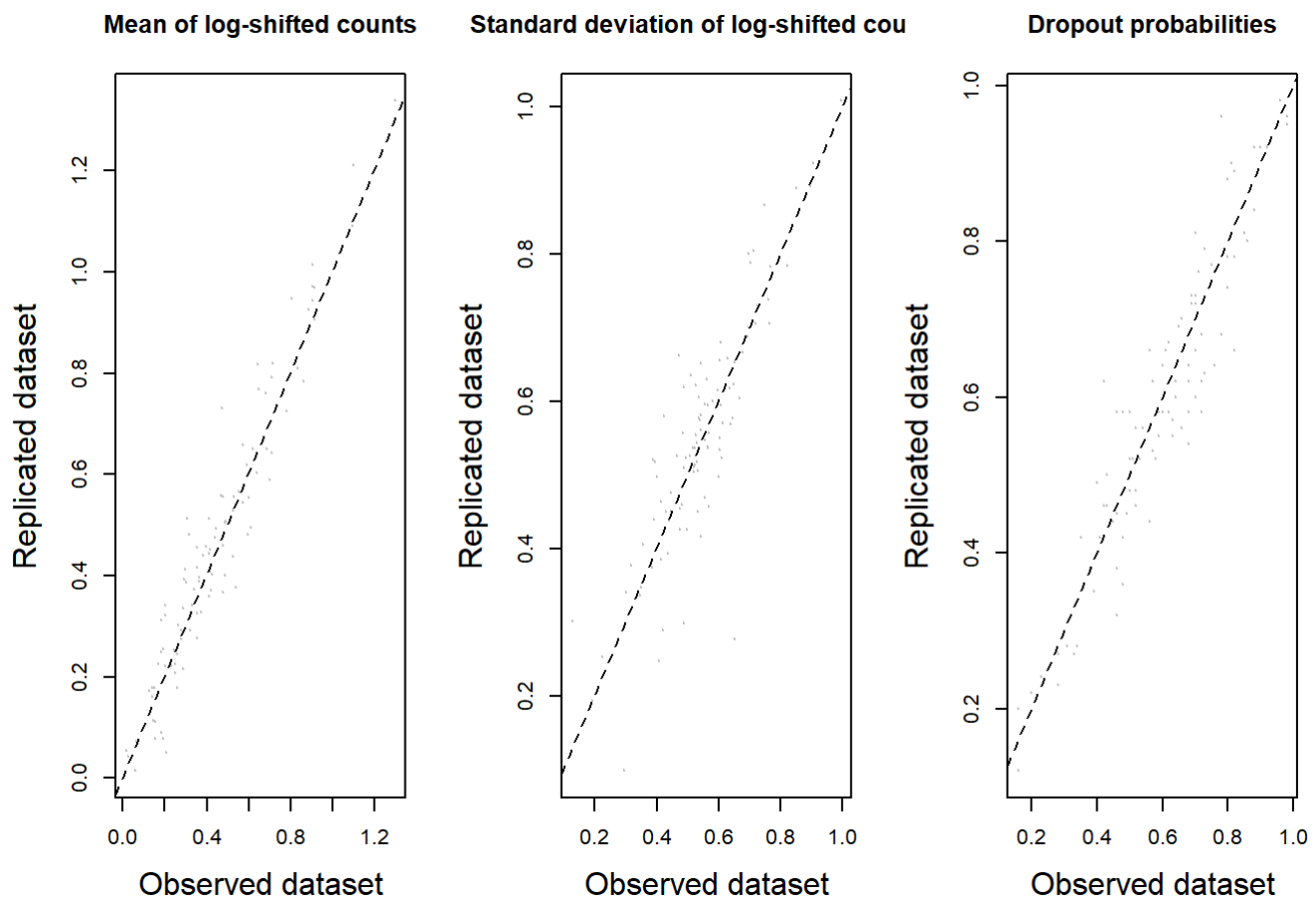
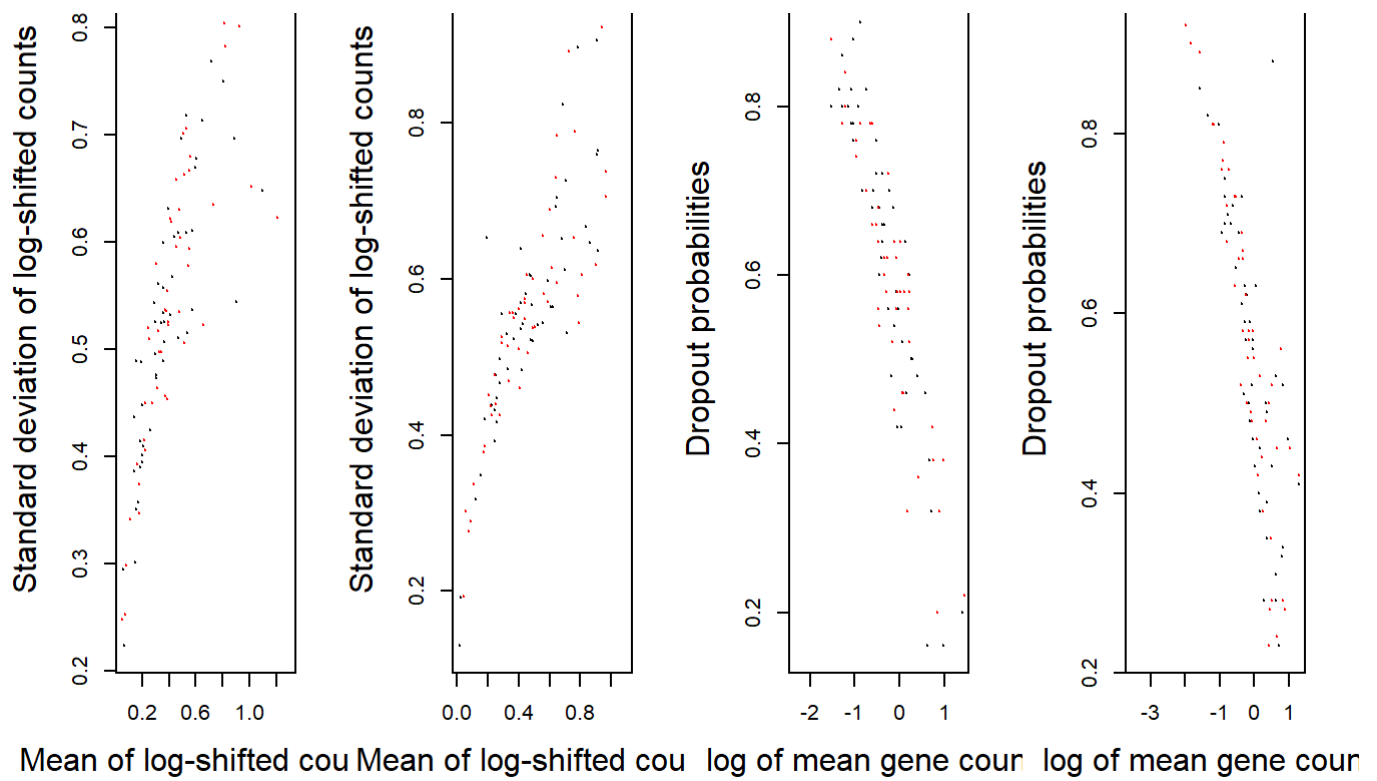
- Relationship between log of mean counts and dropout probabilities (the proportion of zero counts for each gene).

```
# Posterior predictive checks
for(i in 1:2){

  ppc_single(normHDP_output = case1_mcmc,
             Y = Y_linear,
             data.name = c('data 1', 'data 2'))

}
```





In addition, there is also an option to check the recovery of these above relationships with multiple replicated datasets, and to show each relationship for each dataset on one plot. We use function `ppc()` to obtain multiple replicates; example below uses 30 replicates. The `ppc()` function outputs 2 items:

- Statistics of the replicated datasets. The index of the replicated datasets are labelled with variable `t` in the data frame.
- Statistics of the observed datasets.

```
# Posterior predictive checks multiple
case1_ppc_multiple <- ppc(normHDP_output = case1_mcmc,
  Y = Y_linear,
  number_rep = 30)
```

Unlike ppc with single replicate, ppc with multiple replicates ignores the exact statistic for each gene and only consider the general fitted relationship.

```
ppc_plot(ppc_output = case1_ppc_multiple,
  data.name = c('data1', 'data2'))
```

