

Bayesian Two-way Clustering for Gene Expression Data

Graeme Ambler

July 12, 2003

MackKay and Miskin's model

- MackKay and Miskin (2001) proposed the following model for gene expression data:

$$y_{gs} = \sum_{h=1}^H a_{sh} b_{gh} + \varepsilon_{gs}.$$

where g = gene, s = sample, and ε is a noise term.

- We think that this model is rather too general to use as-is.
- We propose some simplifications of the model which (we hope) will lead to more interpretable results.

A single sample model

- We first consider a simple model for the single-sample case:

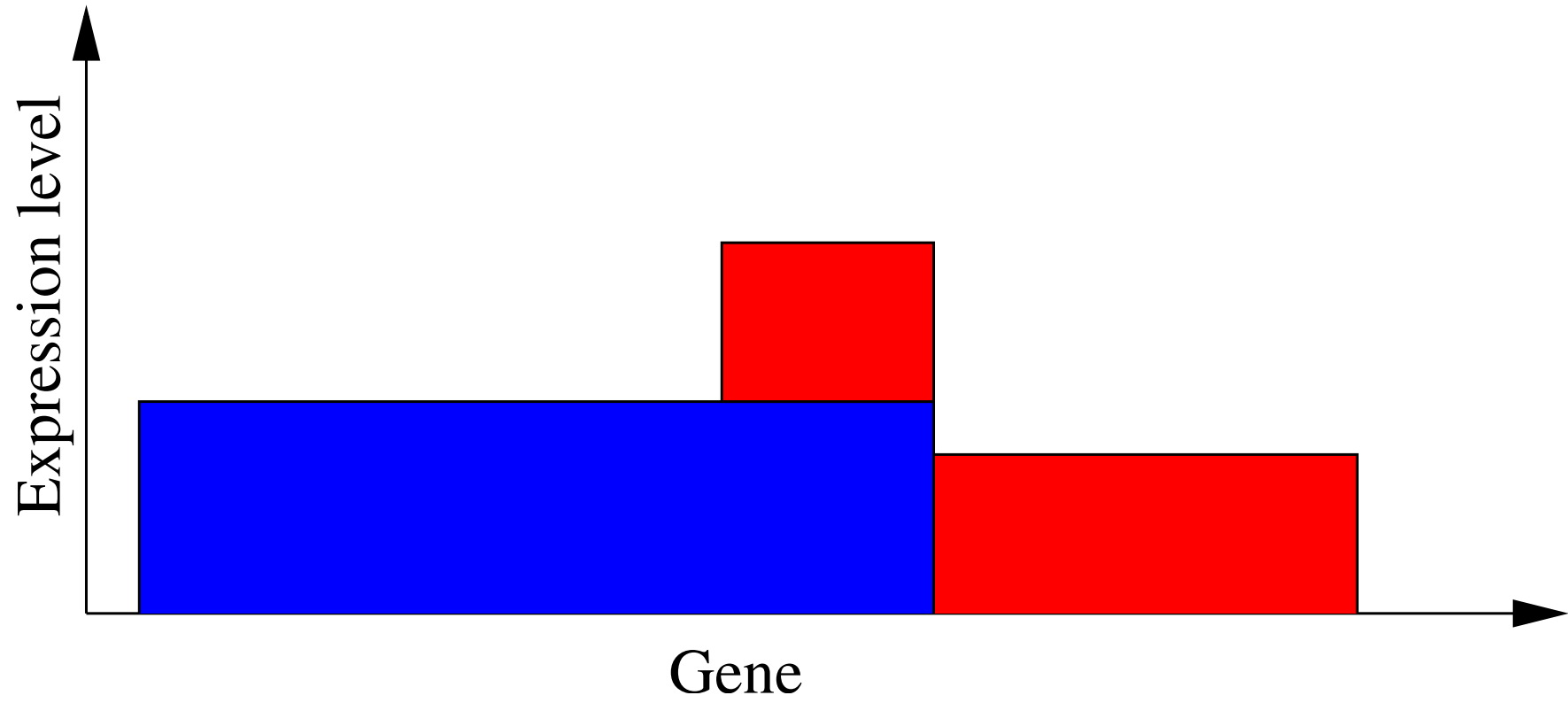
$$y_g = \sum_{h=1}^H z_{gh} b_h + \varepsilon_g,$$

where $z_{gh} \in \{0, 1\}$ is an allocation variable and $b_h \in \mathbb{R}$ is a constant level for layer h .

- This is very similar to the standard mixture model with constant variance:

$$y_g = \sum_{h=1}^H z_{gh} \mu_h + \varepsilon_g$$

except that the allocations for a mixture model are constrained to give $\sum_h z_{gh} \equiv 1 \forall g$.



Prior specification

We propose a fully Bayesian MCMC implementation using the following priors:

$$z_{gh} \sim \text{Bernoulli}\left(\frac{q}{H+1}\right)$$

$$b_h \sim N(0, \tau^2)$$

$$H \sim \text{Poisson}(\lambda)$$

$$\varepsilon \sim N(0, \sigma^2).$$

We also use a conjugate hyperprior on σ :

$$\sigma^{-2} \sim \Gamma(\alpha, \beta).$$

All of the other hyperparameters are held constant.

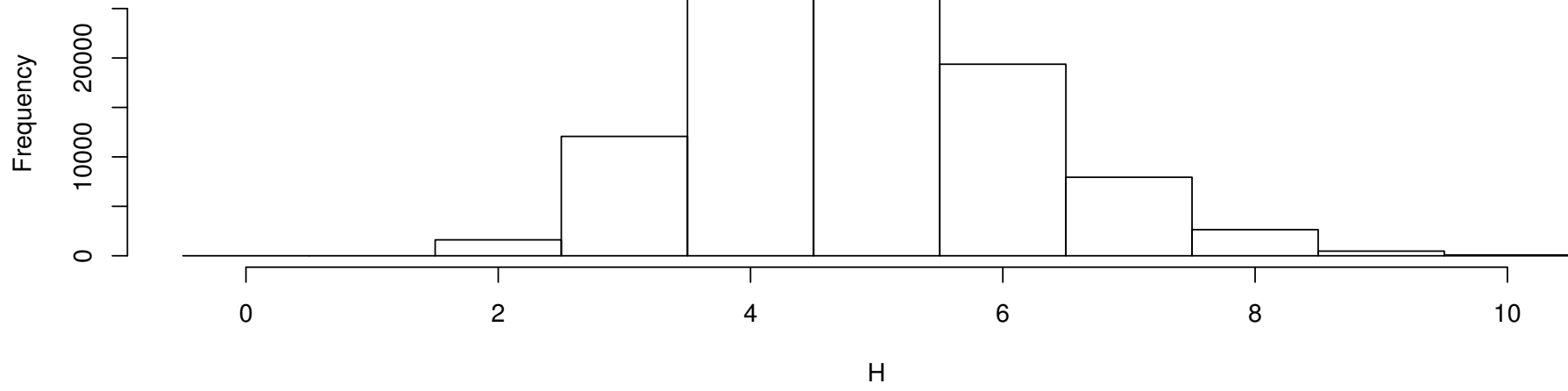
Implementation

- I have implemented an MCMC sampler of this model in C++.
- The sampler performs updates of z , b and σ using Gibbs sampler steps.
- The number (and composition) of layers is updated using both birth-death and split-merge moves.
- There are two different split proposals: one to symmetric values of b_h , and a “bud” move, where one component retains the old value of b_h from the pre-split layer.

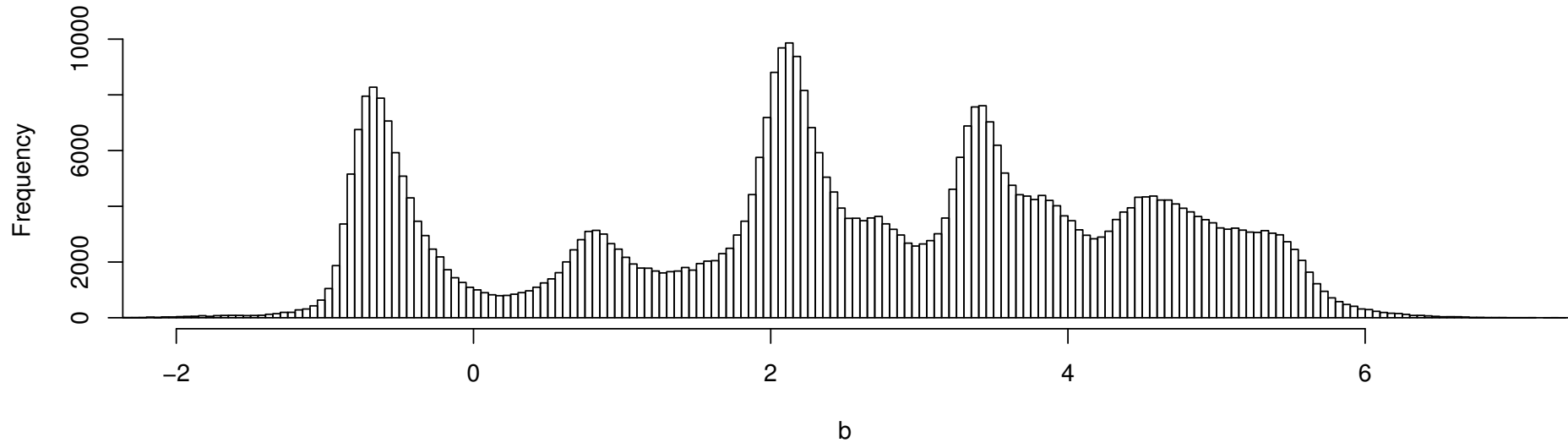
Example: synthetic data

- The test data set consisted of 500 genes with layers at 2.2, 3.4 and 4.7.
- There were also composite values at 5.6 ($= 2.2 + 3.4$), 6.9 ($= 2.2 + 4.7$), 8.1 ($= 3.4 + 4.7$) and 10.3 ($= 2.2 + 3.4 + 4.7$), as well as some values at 0.0.
- Under this model, the number of layers was over-estimated.

Histogram of H

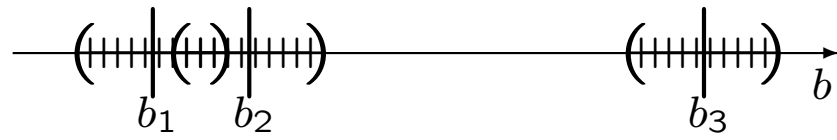


Histogram of b



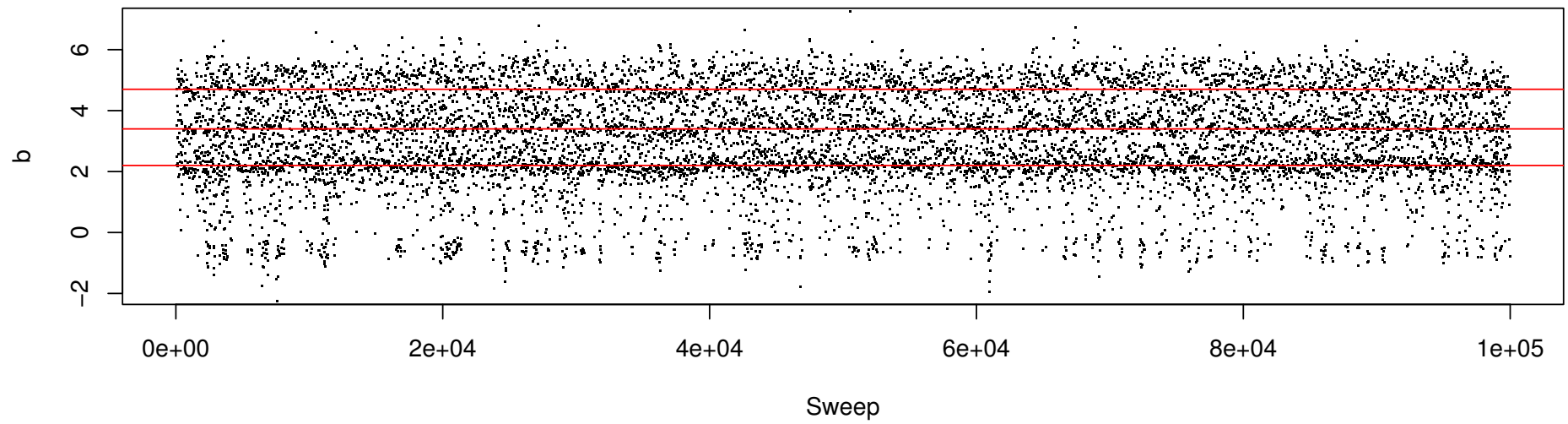
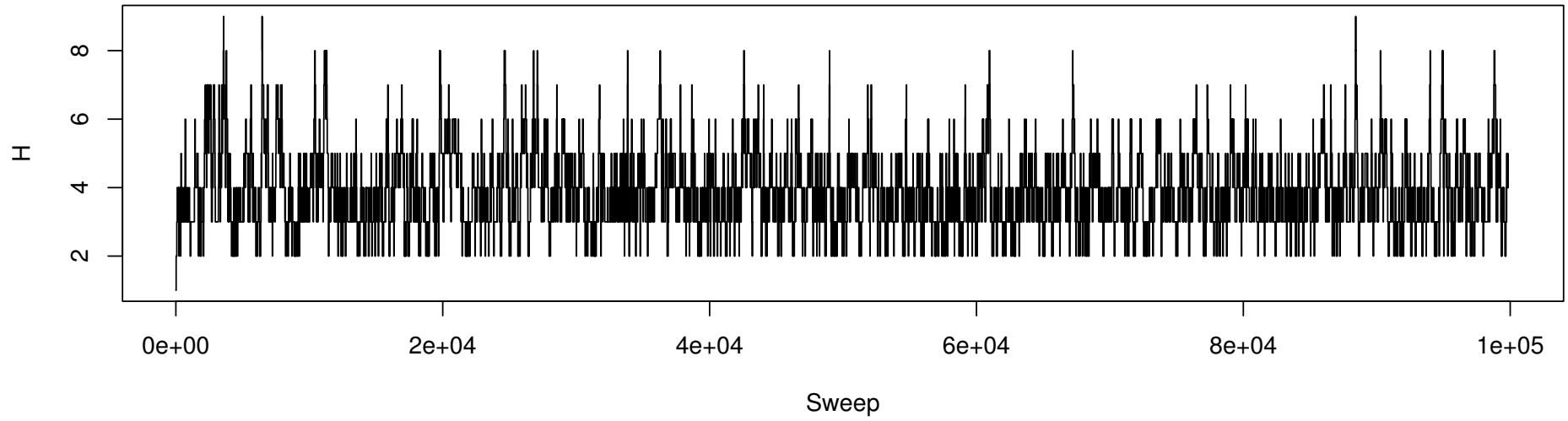
Repulsive b 's(!)

- We wanted to eliminate the possibility of getting several layers with very similar values of b allocated to distinct groups of genes.
- To solve this problem we propose to use some sort of repulsion between values of b .
- Since I have worked with area-interaction point processes before (Baddeley and van Lieshout 1995), it seemed natural to me to use a 1-d area-interaction process for this purpose.

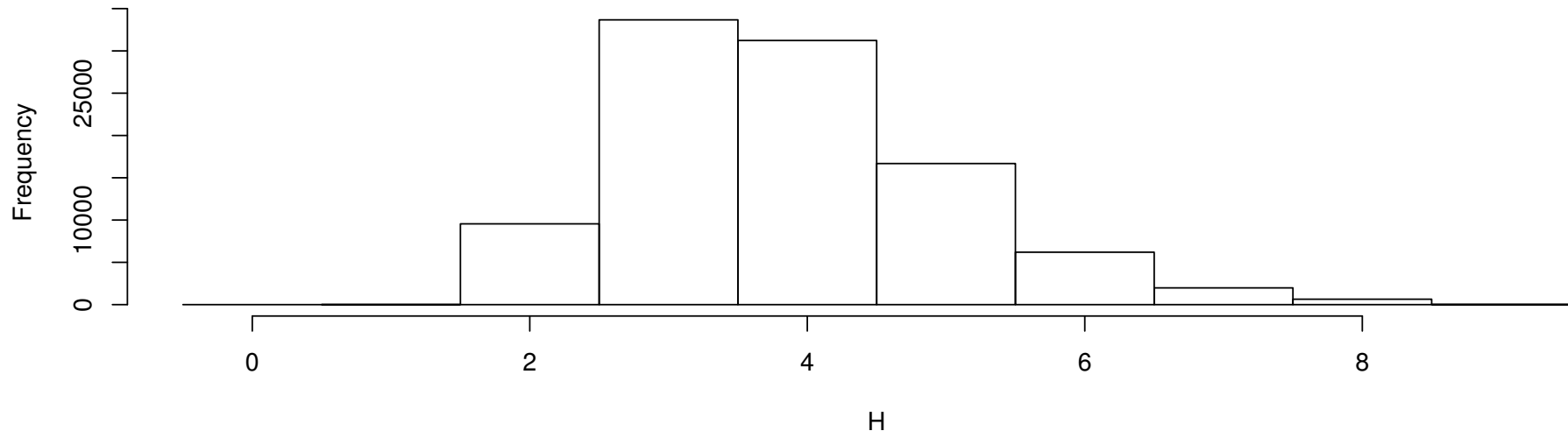


More programming...

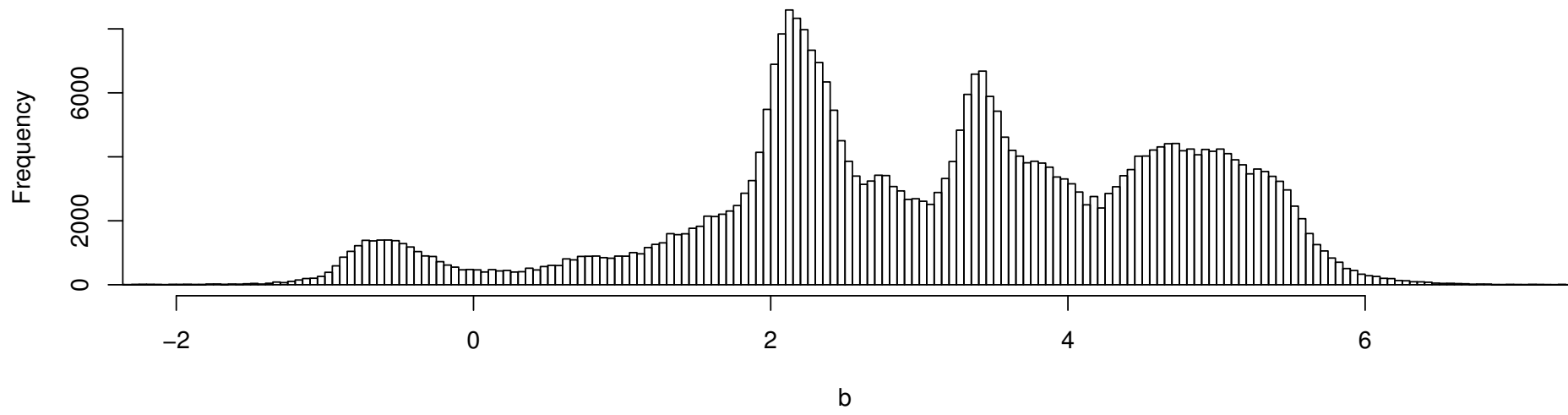
- I have also implemented this additional feature.
- The b -updates are now ‘vanilla’ random-walk Metropolis moves.
- The sampler now finds the correct number and composition of layers with high probability.
- Between-model moves are now accepted less frequently.



Histogram of H



Histogram of b



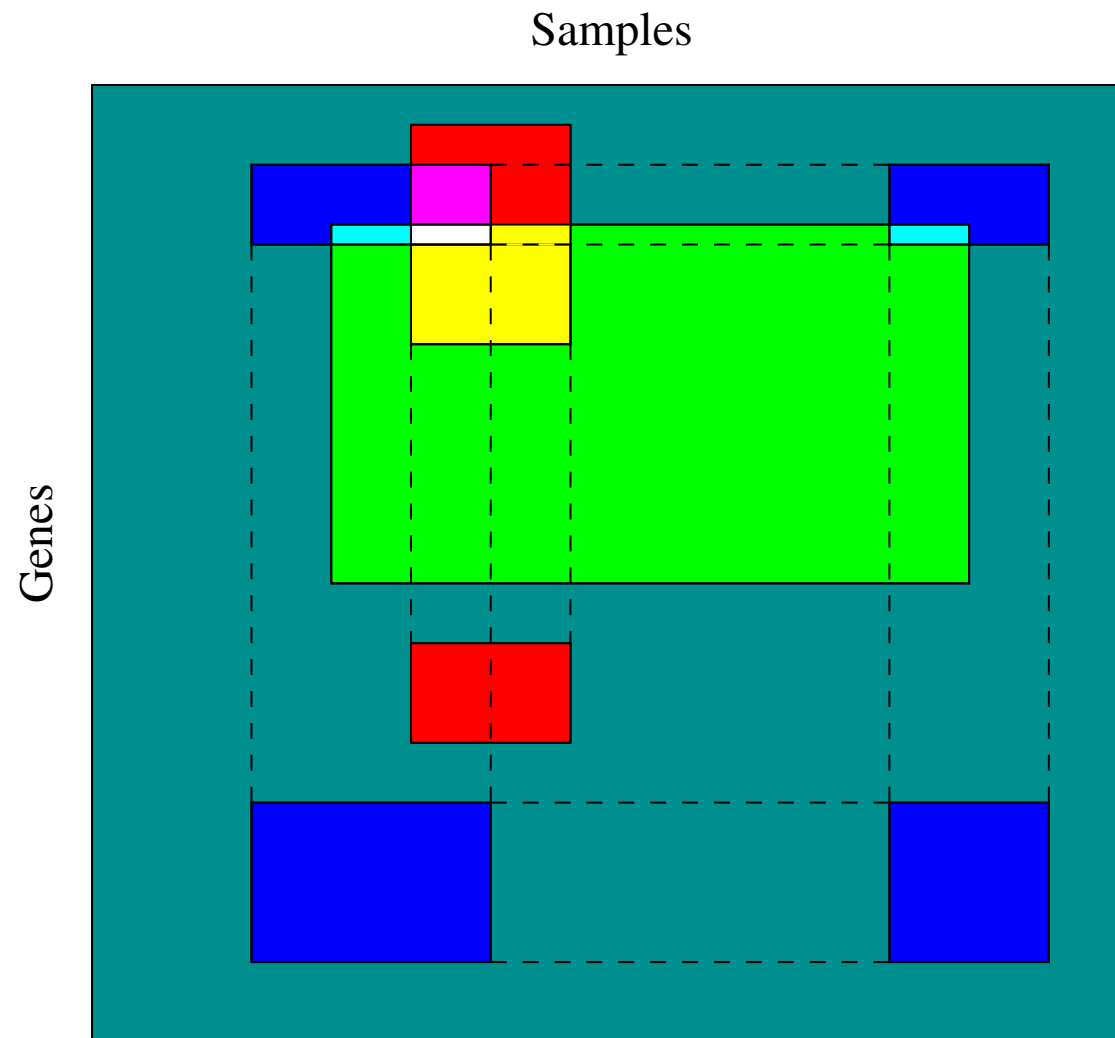
The full model

- Having considered the single-sample case, we now extend the model to the multiple-sample case by adding per-sample allocations:

$$y_{gs} = \sum_{h=1}^H v_{sh} z_{gh} b_h + \varepsilon_{gs},$$

where s =sample and $v_{sh} \in \{0, 1\}$ is an allocation variable with a simple Bernoulli prior for the v 's.

- The layers are constant-height “rectangles”.
- This model is capable of describing basic differential expression of groups of genes between samples.

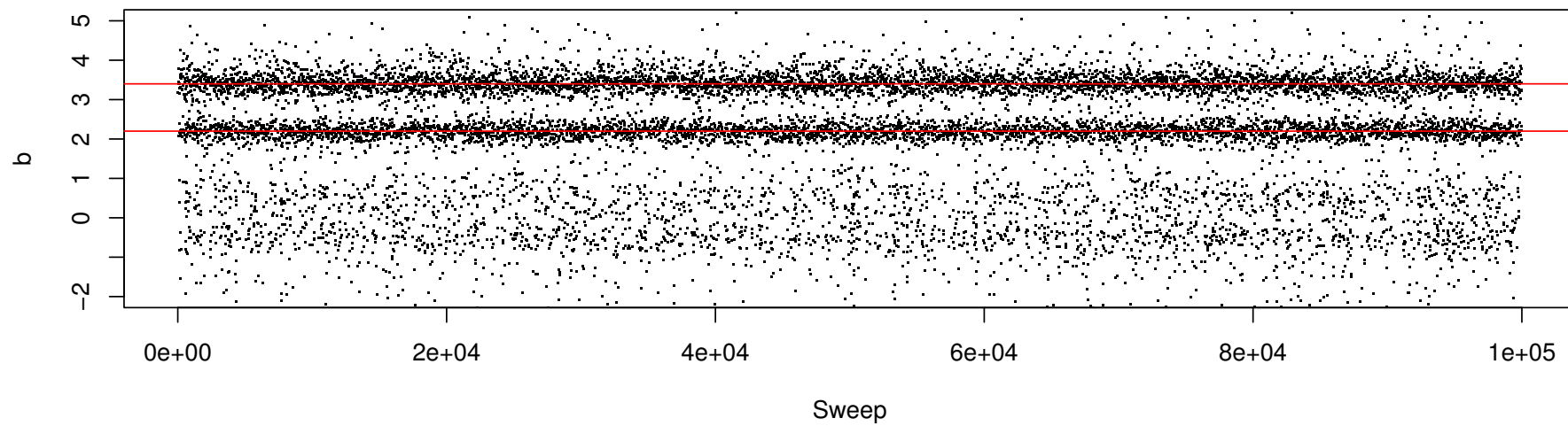
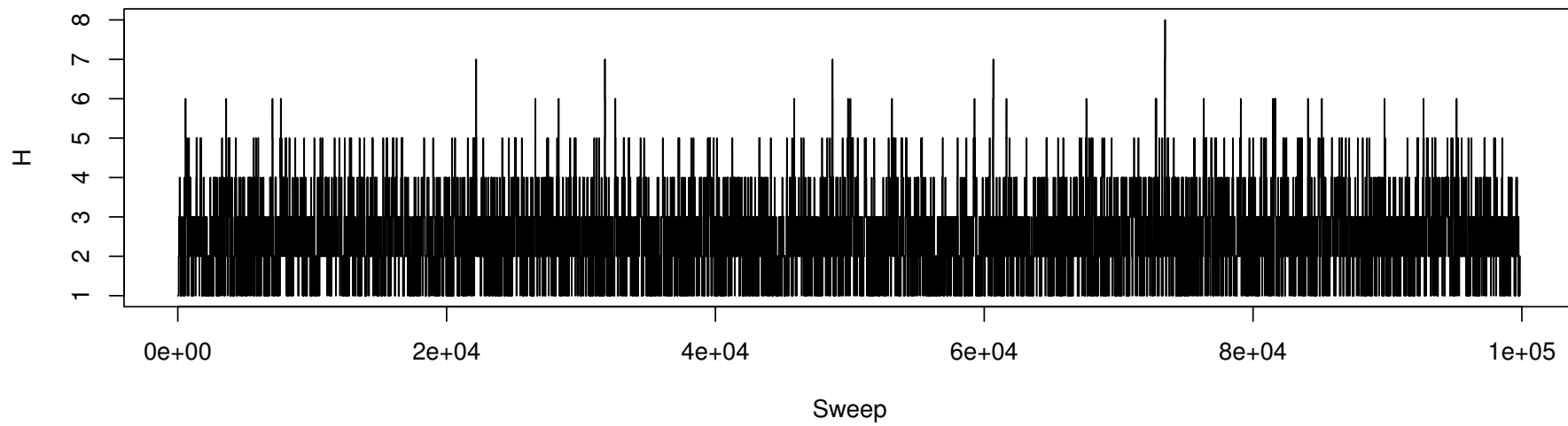


A toy example of the full model

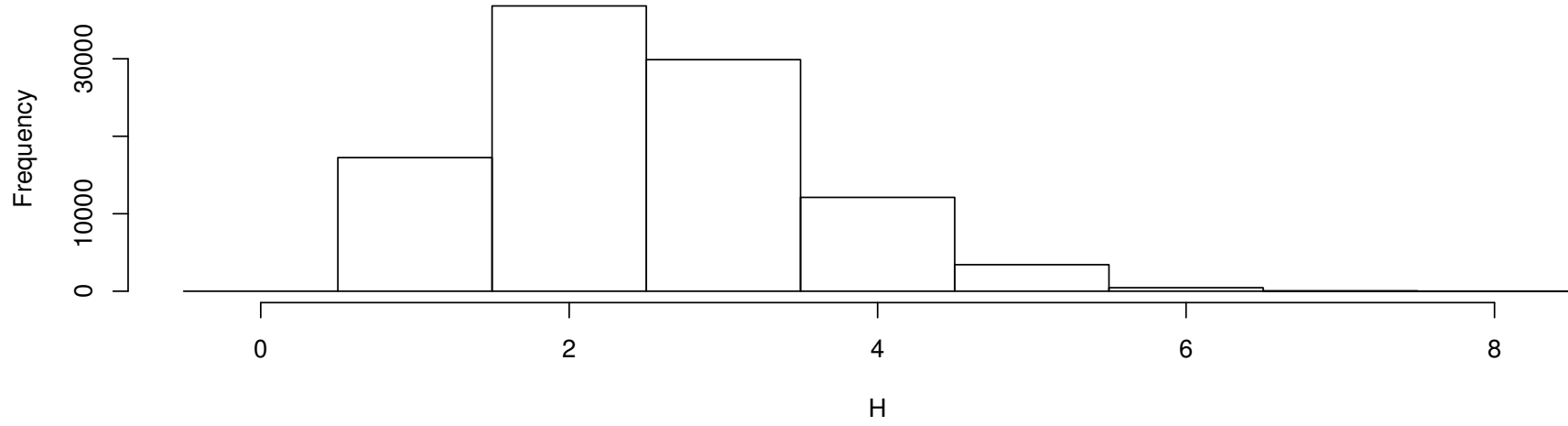
- We next simulated a dataset containing 10 genes for each of 5 samples with 2 layers at 2.2 and 3.4:

2.2	2.2	0	0	2.2
2.2	5.6	3.4	3.4	2.2
0	3.4	3.4	3.4	0
0	3.4	3.4	3.4	0
0	0	0	0	0
0	0	0	0	0
0	0	0	0	0
0	3.4	3.4	3.4	0
2.2	5.6	3.4	3.4	2.2
2.2	2.2	0	0	2.2

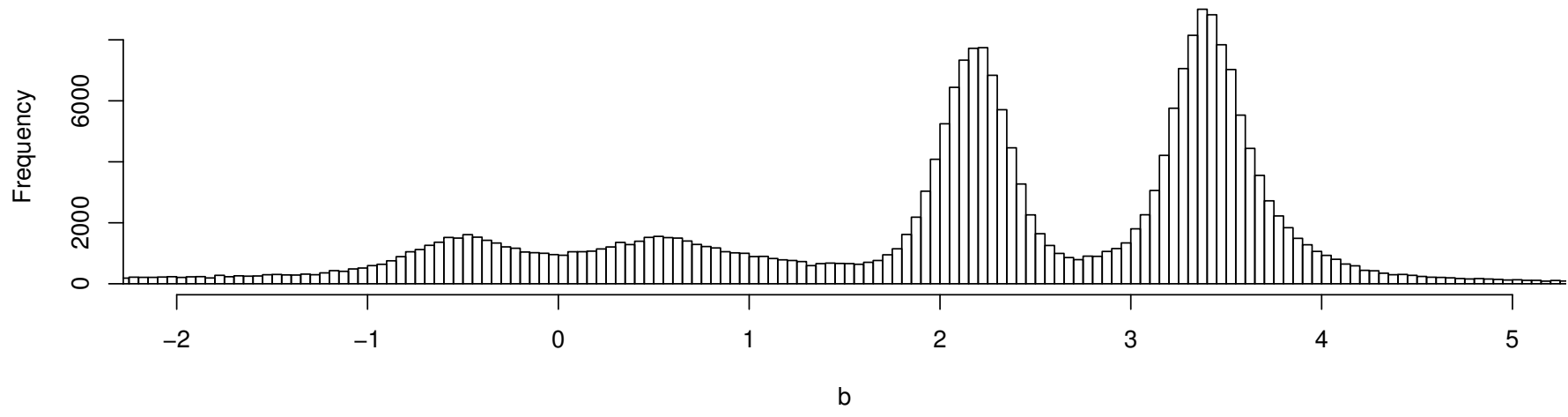
- We added $N(0, (\frac{1}{2})^2)$ noise to this dataset as well.
- The algorithm found the correct number and composition of the layers without difficulty.



Histogram of H



Histogram of b



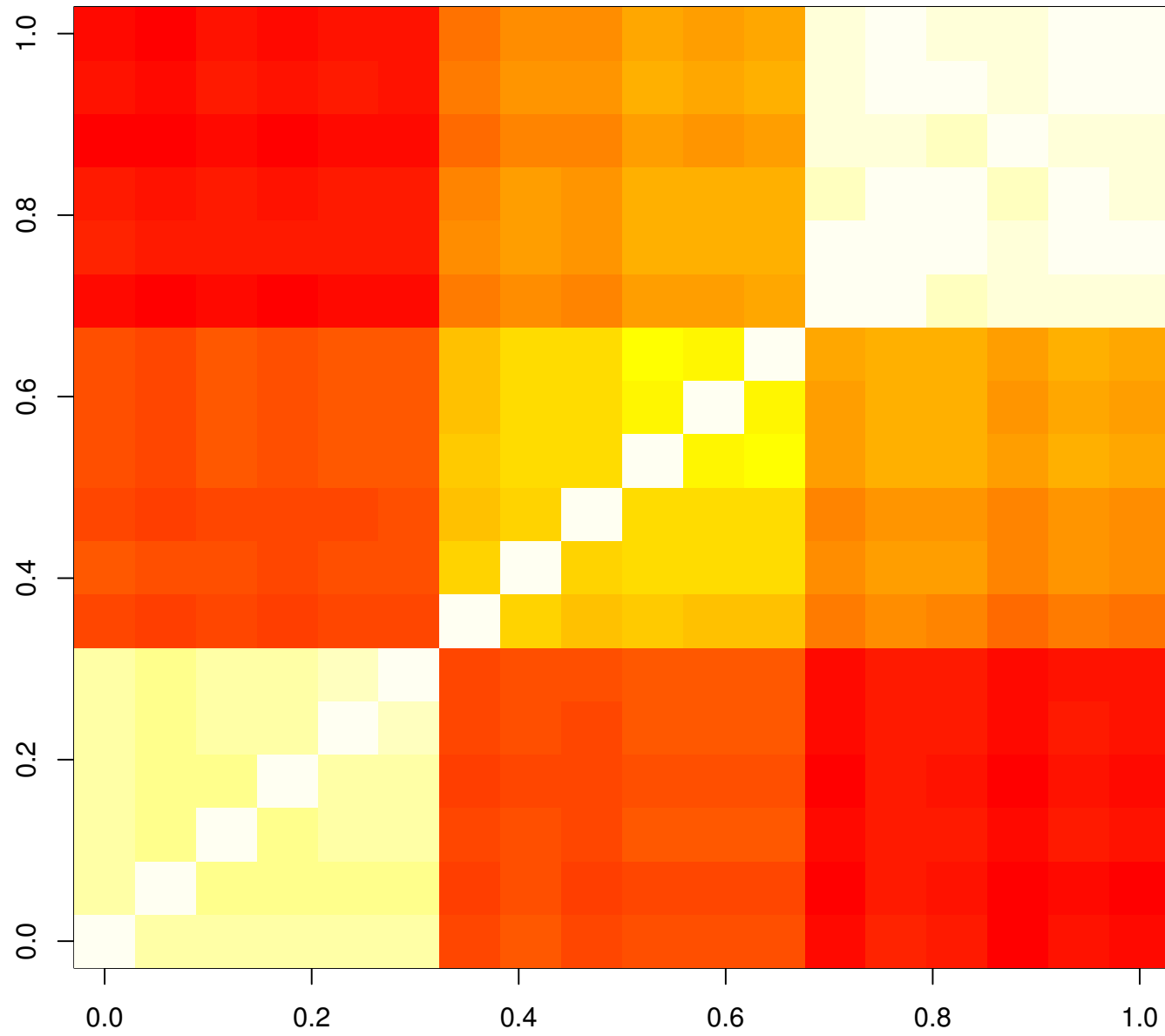
A real dataset

- In order to demonstrate our method with a real dataset, we looked at the human fibroblast data introduced by Lemon et al. (2002).
- This data set consists of 18 samples split into 3 categories: serum starved, serum stimulated and a 50:50 mix of starved/stimulated.

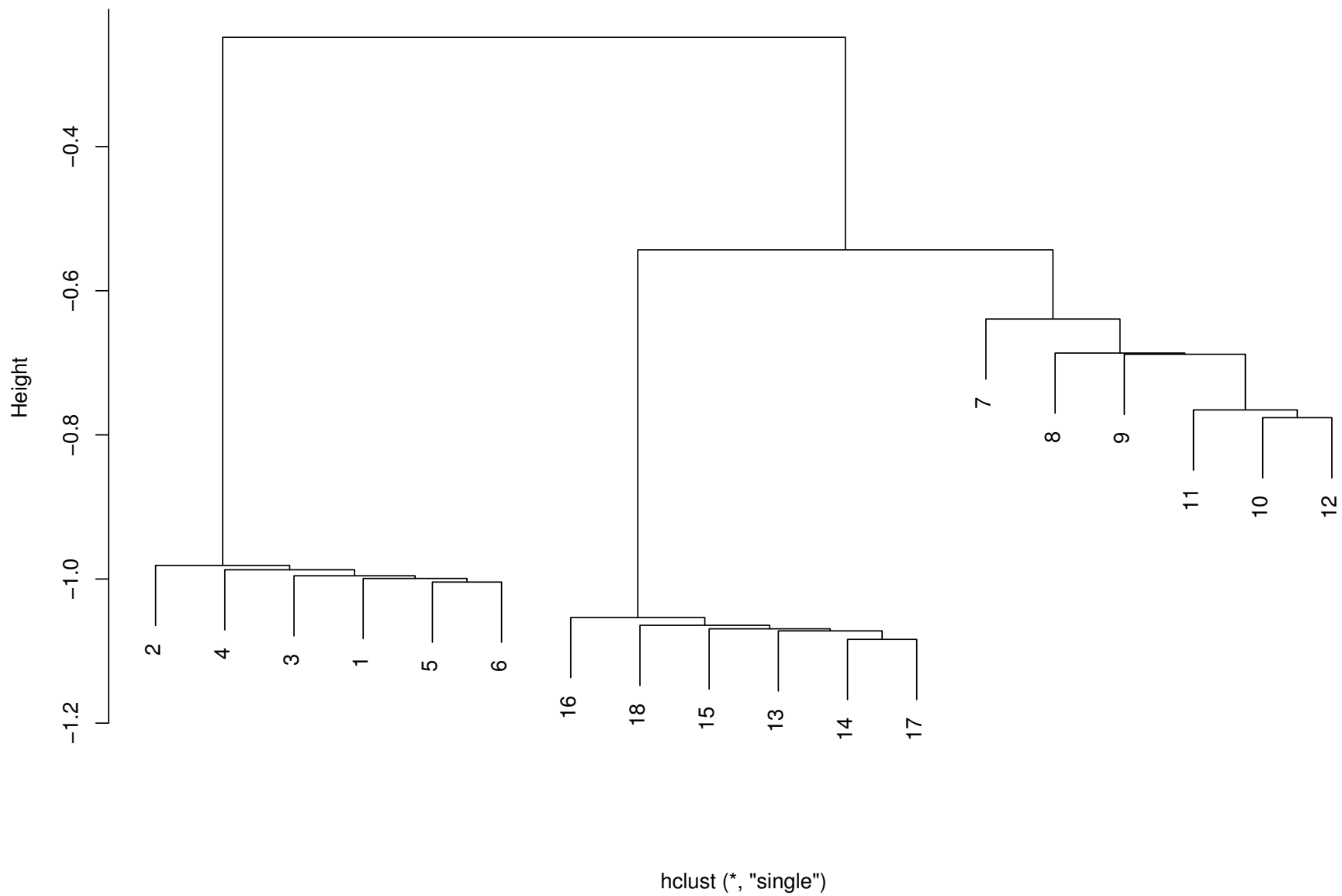
Pre-processing

- We used the natural logarithm of Lemon et al.'s calculated LWF values as our measure of expression and subtracted gene and sample mean levels.
- We then selected the 100 most variable genes across all 18 samples and used this 18×100 array as the input to our sampler.
- The following plots use the expected number of common layers that a pair of samples (genes) are in to cluster the samples (genes).

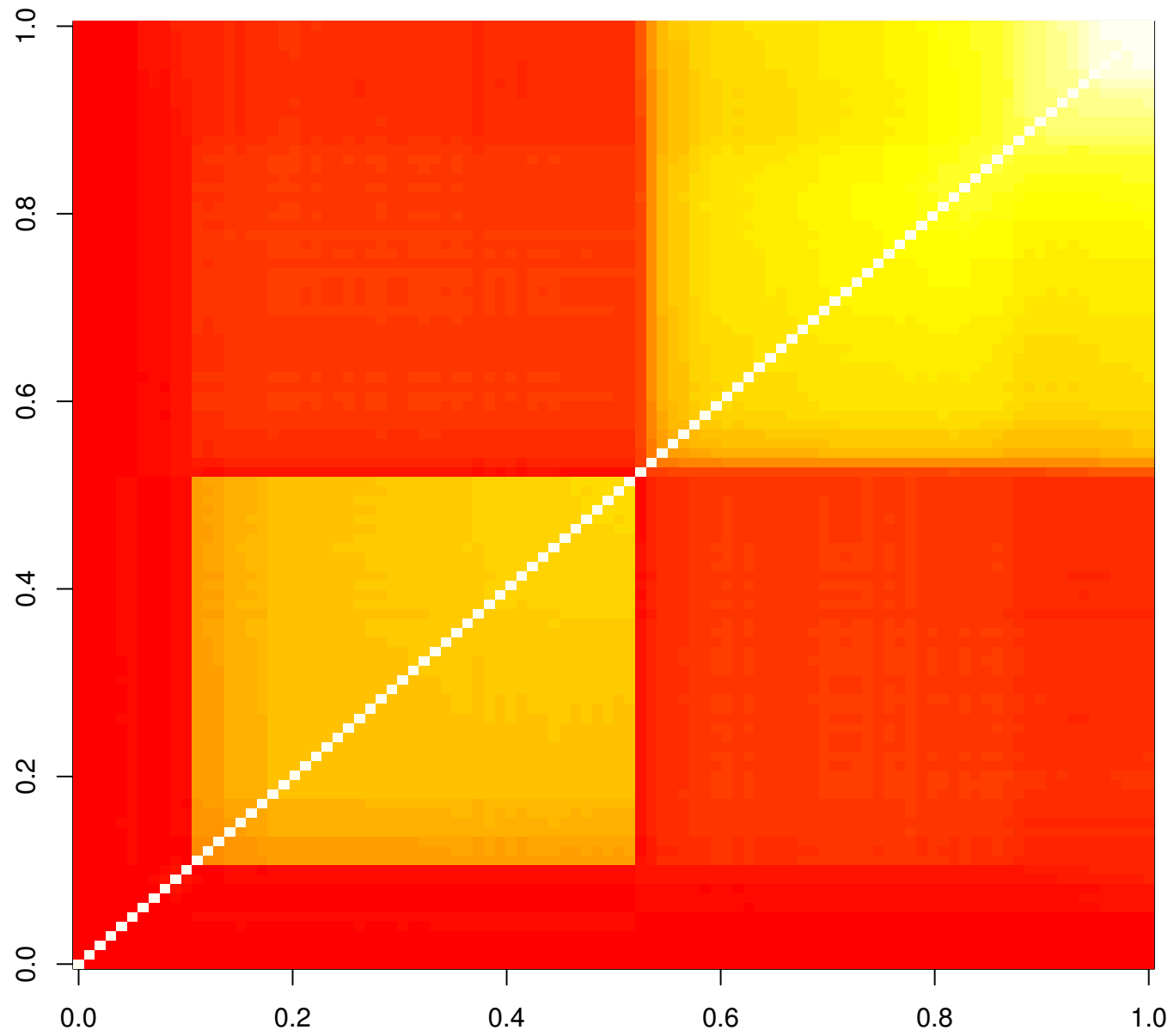
Sample similarities



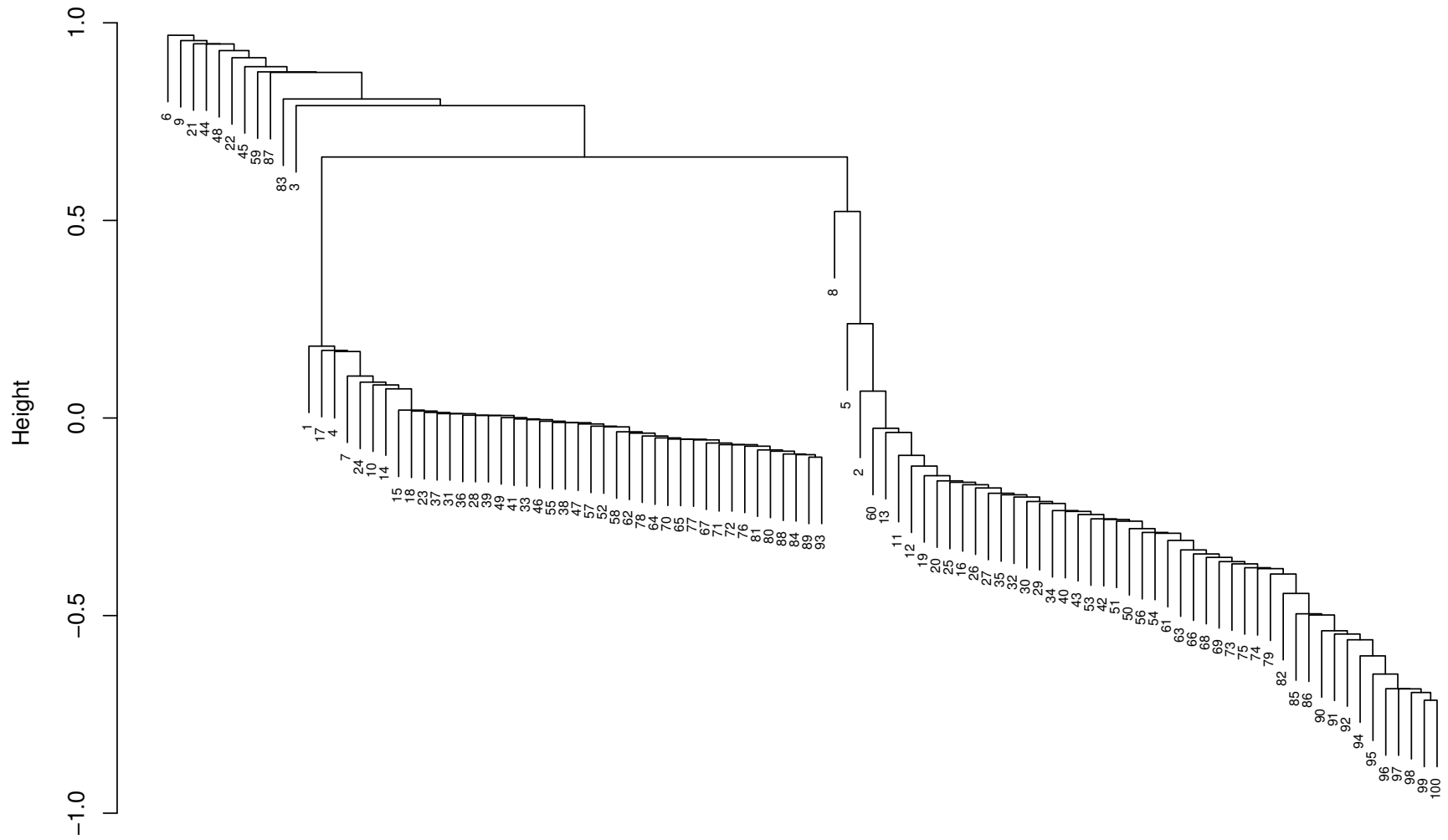
Cluster Dendrogram for Samples



Sorted (using hclust) gene similarities

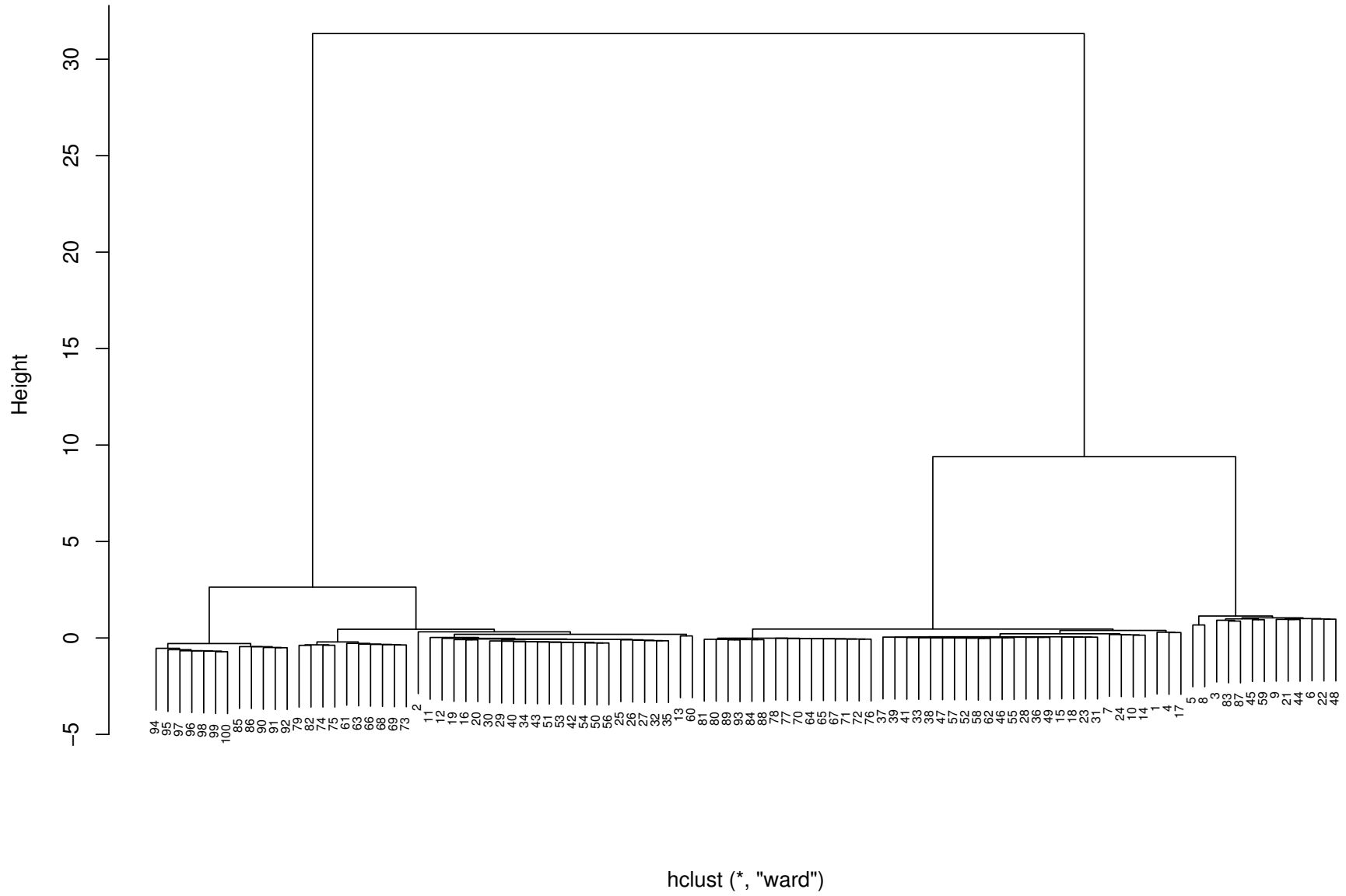


Cluster Dendrogram for Genes

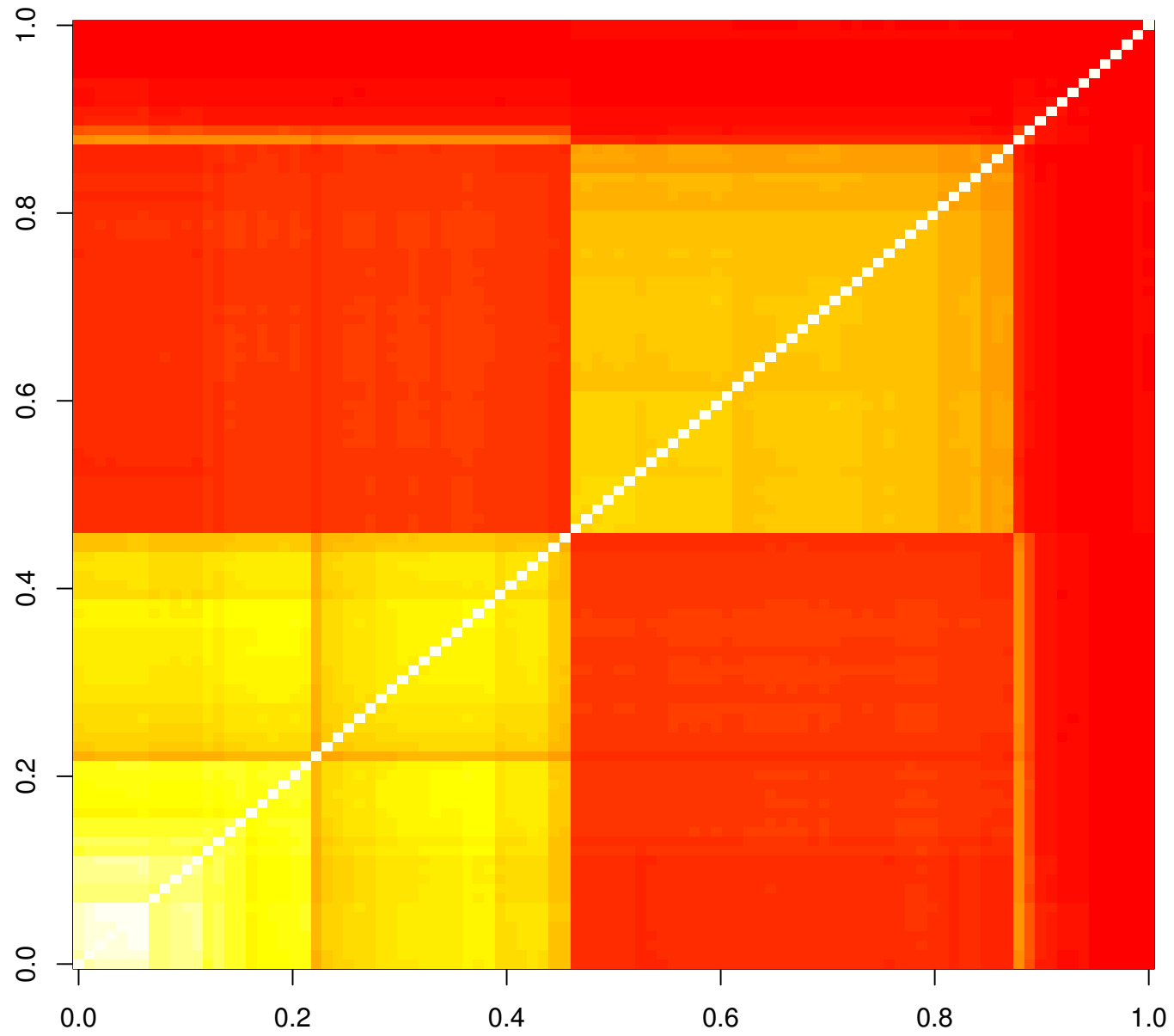


hclust (*, "single")

Cluster Dendrogram for Genes



Sorted (using hclust) gene similarities



Future work

- Work on intelligent methods for choosing the parameters.
- Investigate methods for displaying the posterior in a visually-informative way.

References

Baddeley, A. J. and M. N. M. van Lieshout (1995). Area-interaction point processes. *Annals of the Institute for Statistical Mathematics* 47, 601–619.

Lemon, W. J., J. J. T. Palatini, R. Krahe, and F. A. Wright (2002). Theoretical and experimental comparisons of gene expression indexes for oligonucleotide arrays. *Bioinformatics* 18, 1470–1476.

MacKay, D. J. C. and J. Miskin (2001). Latent variable models for gene expression data. Technical report, Cavendish Laboratory, Madingley Road, Cambridge CB3 0HE. United Kingdom.