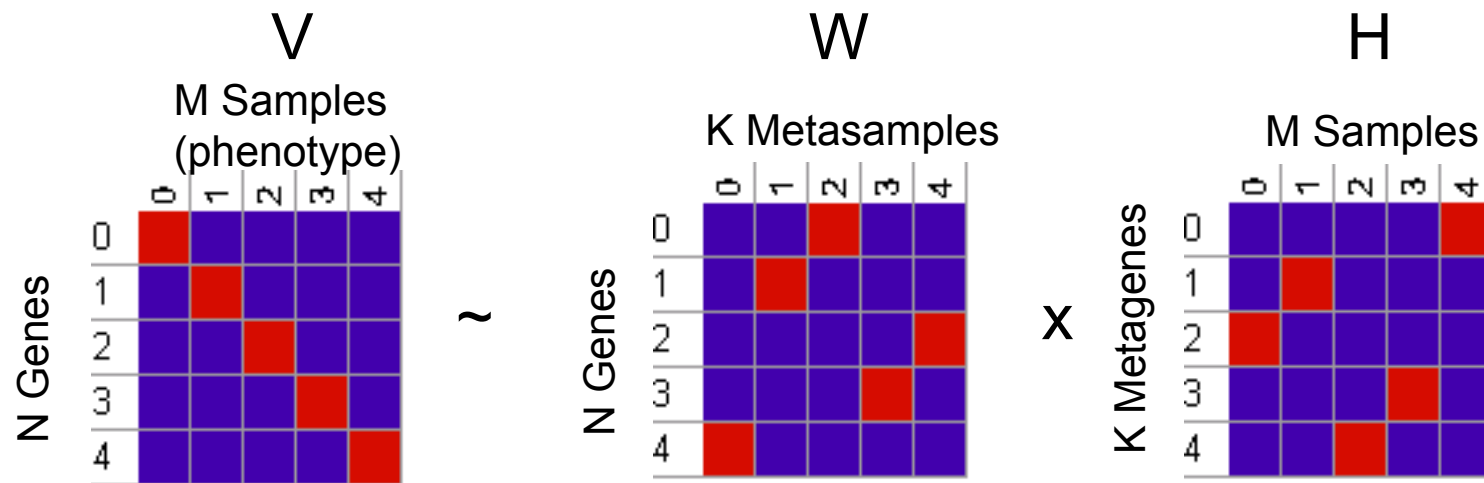


Non-Negative Factorization Examples and Analysis of Encoding and Clustering Tendencies (as a way to interpret the meaning of “metasamples” and “metagenes” in the context of “discovered” patterns, clusters and elements of a biological “by parts” representation).

Definitions:



In a real microarray example the “cells” may represent a group of several genes or samples with similar expression.



V

W

H

	0	1	2	3	4
0	Red	Blue	Blue	Blue	Blue
1	Red	Red	Blue	Blue	Blue
2	Blue	Blue	Red	Red	Red
3	Blue	Blue	Blue	Red	Red
4	Blue	Blue	Red	Blue	Red

~

	0	1
0	Blue	Red
1	Blue	Red
2	Red	Blue
3	Red	Blue
4	Red	Blue

X

	0	1	2	3	4
0	Blue	Blue	Red	Red	Red
1	Red	Red	Blue	Blue	Blue

The phenotypes have mainly expression in genes 0-1 and 2-4. It recognizes that and separates them into to groups of metasamples. The same clustering takes place with the metagenes.

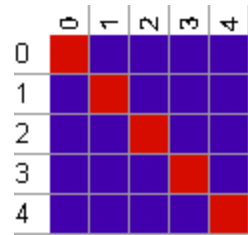
[illegible]

	k0	k1	k2
gene5_k0	red	blue	blue
gene7_k0	red	blue	blue
gene1_k0	red	blue	blue
gene8_k0	red	blue	blue
gene2_k1	blue	red	blue
gene9_k1	blue	red	blue
gene4_k2	blue	blue	red
gene0_k2	blue	blue	red
gene3_k2	blue	blue	red
gene6_k2	blue	blue	red

[illegible][illegible]

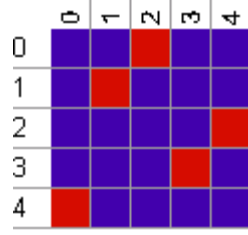
W H

V



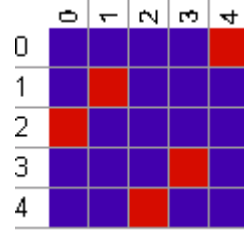
~

W



X

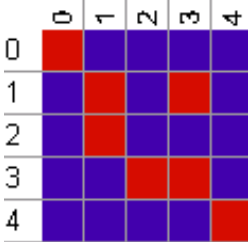
H



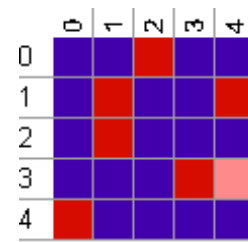
Encoding Tendency

Same number of factors as genes and samples

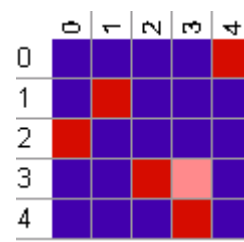
Each phenotype and gene is unique and orthogonal and this is preserved by the encoding. As the data is already sparse the encoding doesn't produce new patterns.



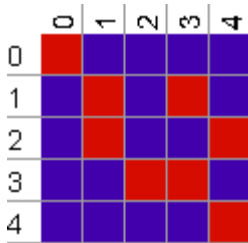
~



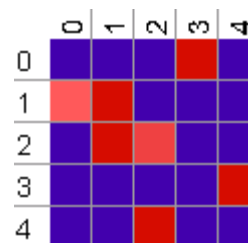
X



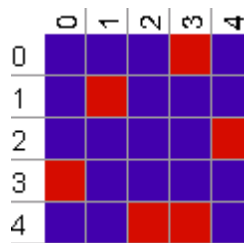
W and H are similar to V (with different column and row order) with one notable difference: gene 1 is decoupled and produced by combining metagenes 1 and 4. Metagene 4 is a new pattern that helps to build the profiles for genes 1 and 3. Notice the small increase in sparseness in matrix H (6 red cells) compared with V (7 red cells). This is an example of the encoding "by parts." As the V matrix is symmetric the same encoding trick could have been applied to the samples (e.g. making sample 3 and 1 with a metasample that has only gene 3 expressed. This break of symmetry is probably due to the fact that W is updated before H.



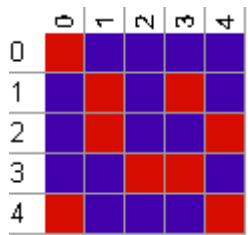
~



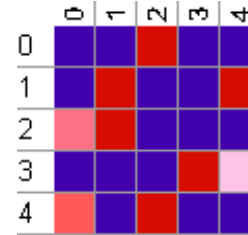
X



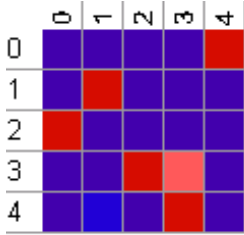
Similar metagenes as previous example. Previous metasamples 1, 2 and 3 are now 1, 3 and 4. The additional expression of gene 2 in sample 4 changes the encoding of previous metasamples 0 and 4. The new metasamples 0 and 2 allow to reproduce a V with higher density (8 red cells) with the same density in W (7 red cells).



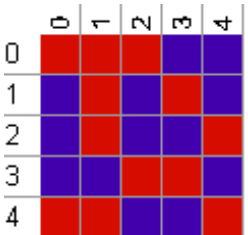
~



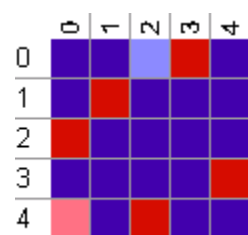
X



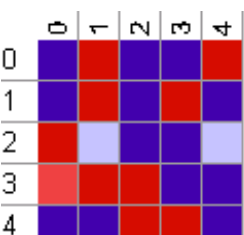
The metasamples are very similar to the samples. The metagenes decompose the gene patterns in an efficient way using almost orthogonal rows. This allows to reproduce the V matrix with 9 red cells with only six red cells in H.



~



X

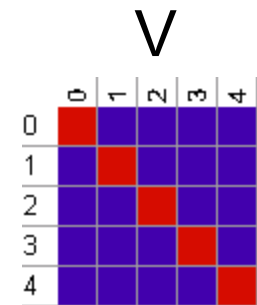


Here the "by parts" decomposition works better for W than H. The gene patterns are harder to decouple and the metagenes look closer to the original genes (no increase sparseness. In contrast W is more sparse with 6-7 red cells instead of 11 for V). There is more sample "overlap" in this case and then for example samples 1 is decomposed in metasamples 3, 2, 1 and 0.

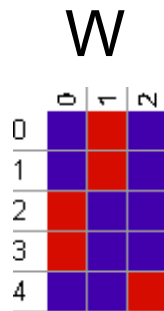
Clustering Tendency

Same number of factors as genes and samples

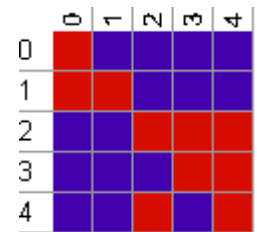
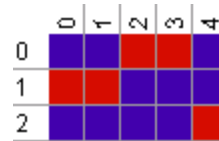
Five orthogonal phenotypes and only three metasamples. It groups them in two groups of two and then one apart.



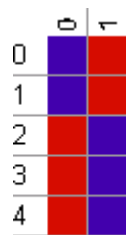
~



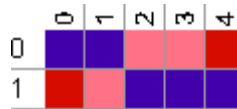
X



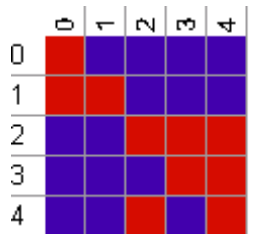
~



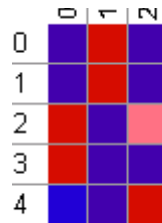
X



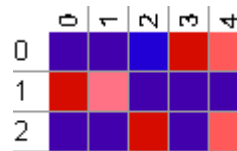
The phenotypes have mainly expression in genes 0-1 and 2-4. It recognizes that and separates them into to groups of metasamples. The same clustering takes place with the metagenes.



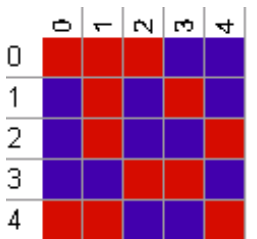
~



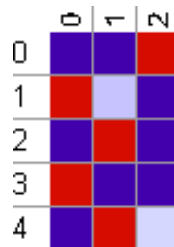
X



Same as before but there is one extra metasample that picks up the phenotype of sample 2. Metagene 0 from the previous example now decouples in two. One (2) now accounts for the new metasample (2).



~

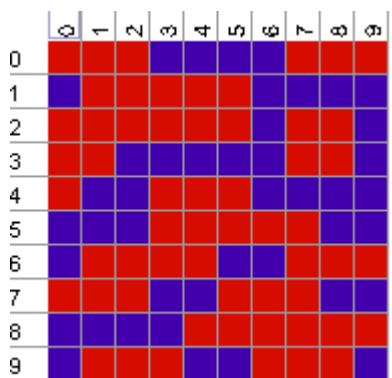


X



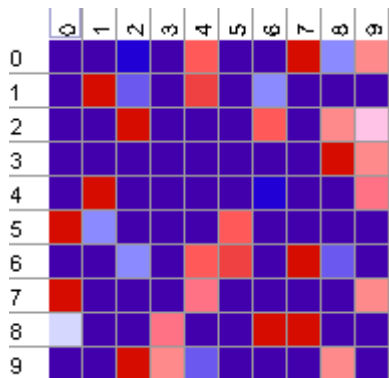
The samples so not easily cluster in less than 5 groups. Here we see a combination of clustering and encoding. It separates the phenotypes of samples 0 and 3 in metasamples 2 and 0 but encodes samples 1, 2 and 4 into metasample 1 and the partially in metasamples 0 and 2. This is illustration of the encoding by parts: sample 1 is mainly made of metagenes 1 and 2 which combine metasamples 1 and 2.

V



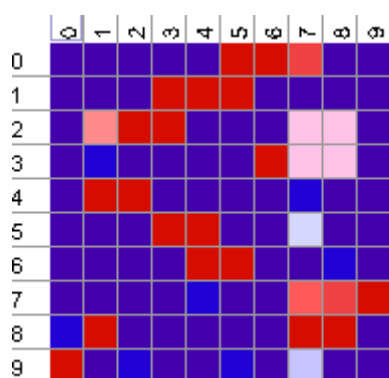
~

W



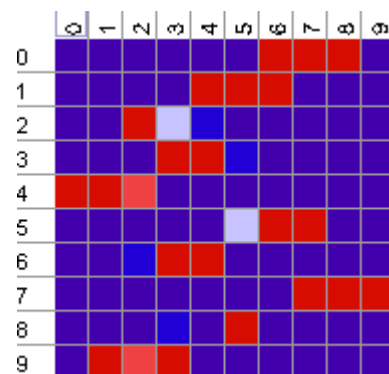
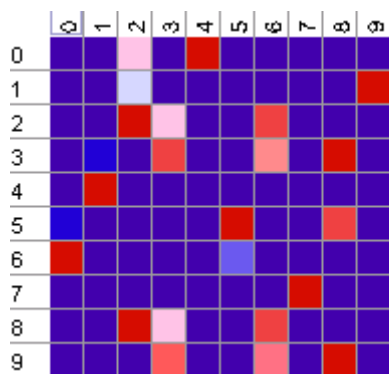
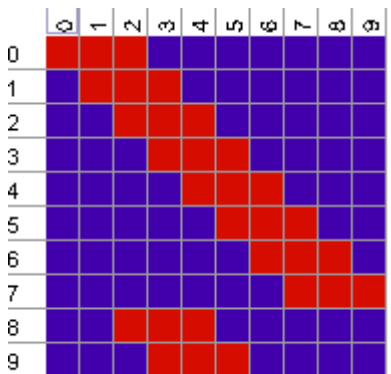
X

H



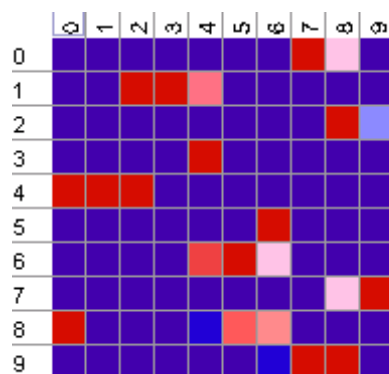
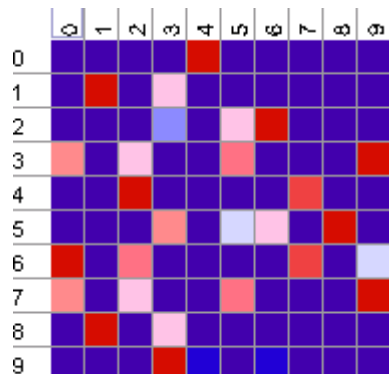
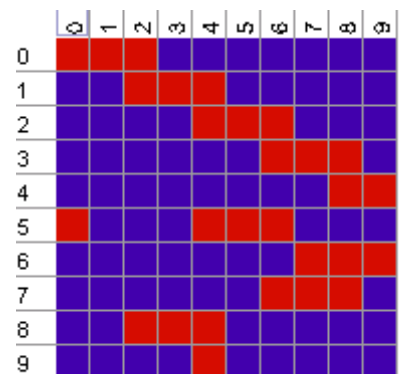
~

X



~

X

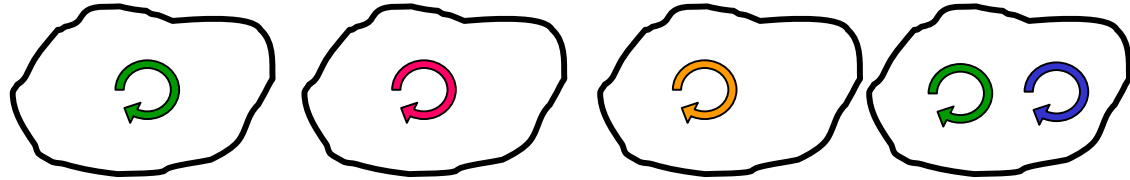


Biological Model and example

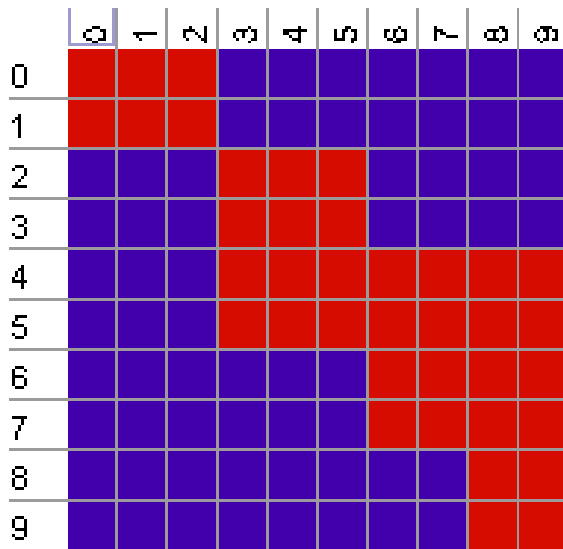
4 metagenes



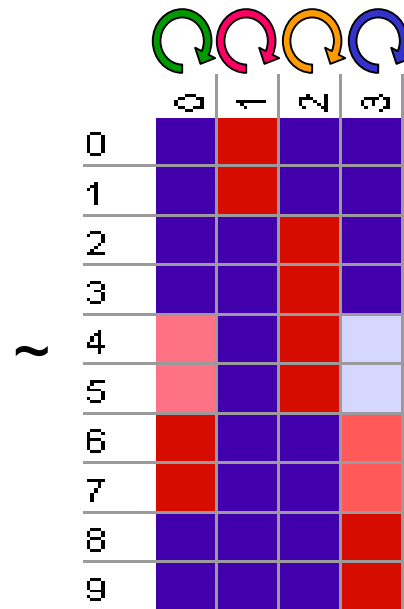
4 metasamples



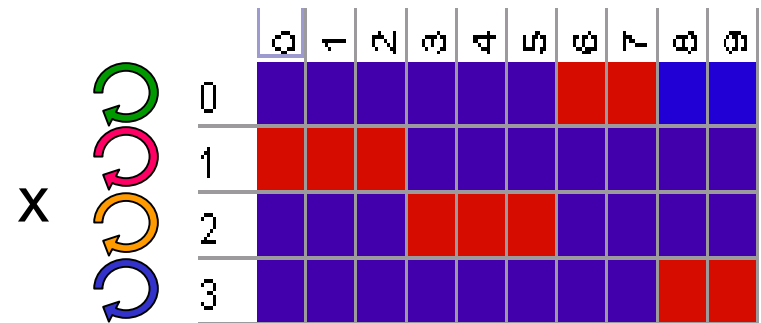
V



W



H



[illegible]

	0	1	2	3
0	0	1	2	3
1	0	0	0	0
2	0	0	0	0
3	0	0	0	0
4	0	0	0	0
5	0	0	0	0
6	0	0	0	0
7	0	0	0	0
8	0	0	0	0
9	0	0	0	0
10	0	0	0	0
11	0	0	0	0
12	0	0	0	0
13	0	0	0	0
14	0	0	0	0

	0	1	2	3	4	5	6	7	8	9
0										
1										
2										
3										

 \sim [illegible]

	0	1	2	3
0				
1				
2				
3				
4				
5				
6				
7				
8				
9				
10				
11				
12				
13				
14				
15				
16				
17				
18				

~

	0	1	2	3	4	5	6	7	8	9
0	0	0	0	1	1	1	0	0	1	1
1	0	0	0	0	0	0	1	1	0	0
2	0	0	1	1	1	1	1	0	0	0
3	1	1	0	0	0	0	0	0	1	0

Summary of trends:

Very similar patterns of genes or samples get cluster together. This is a consequence of the number of factors being less than the number of genes or samples but also of the trend of increase sparseness in the “by parts” decomposition.

A sparse W (or H) as a result of a significant decomposition by parts may imply a less sparse H (or W) because more parts are needed to recreate the samples (or genes).

If V mainly contains orthogonal samples and gene patterns overlap significantly (i.e. V has “row patterns”) then the genes get decomposed in parts and H will be sparse.

If V mainly contains orthogonal genes and samples patterns overlap significantly (i.e. V has “column patterns”) then the samples get decomposed in parts and W will be sparse.