

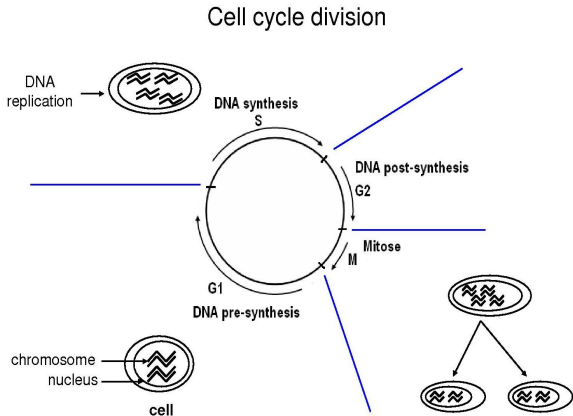
Time Series Clustering: application on cell cycle genes expression profiles

Ahlame Douzal (Ahlame.Douzal@imag.fr)

AMA, LIG, Université Joseph Fourier

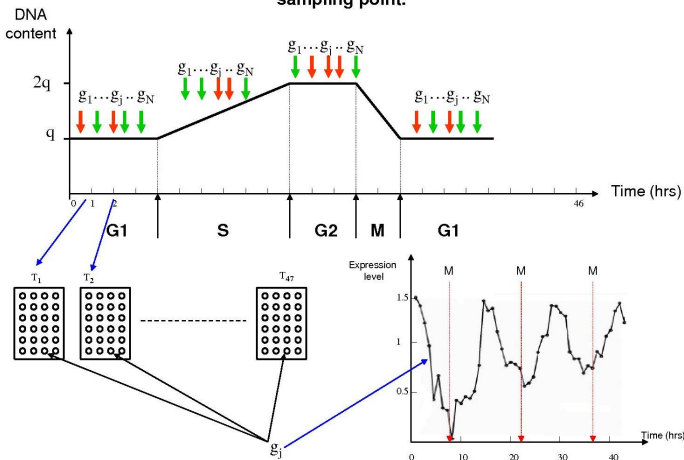
Master 2R - MOSIG (2011)

Problem statement

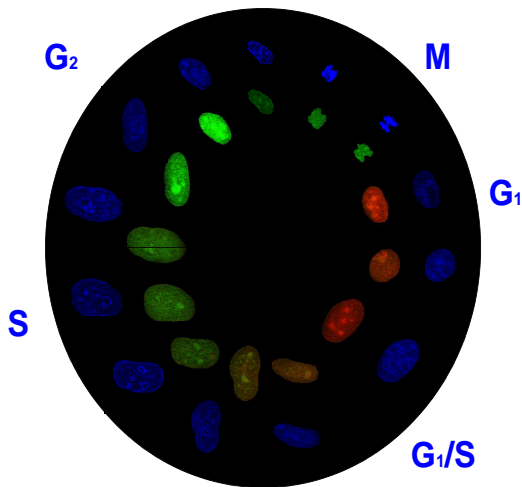


Gene expression data and cell cycle division

In gene expression time-series, each microarray corresponds to a sampling point.



Problem statement



Objectives

- Identification of the cell cycle expressed genes
- Determine of differentially expressed genes:
 - which genes are involved in different type of cells (cancer versus healthy cells) ?

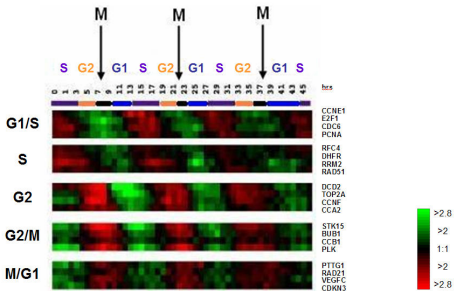
Clustering and classifying have proven helpful to:

- Extract main groups of behaviour expressions
- Identify new genes (unlabeled, unknown)
- Identify new genes relationships: co-expressed genes, co-regulated genes, understand genes functions,

Conventional Approach

Hela data (Whitfield et al. 2002)

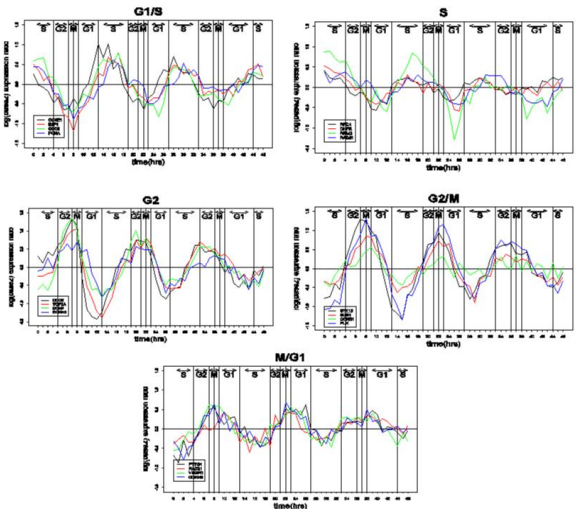
- Studied genes are preprocessed: non-cyclic expression, unexpressed genes and noisy expression
- Experimentally, a set of reference genes are selected



- Each measured gene is assigned to one phase by its peak similarity to the reference genes: Pearson correlation.

The expression profiles of the 20 reference genes

- The expression profiles of the 20 reference genes, illustrating their peak expression at one phase during three cell division cycles. The double arrowed lines delimit the time duration for four cell cycle phases: G_1 , S , G_2 and M .



Aims

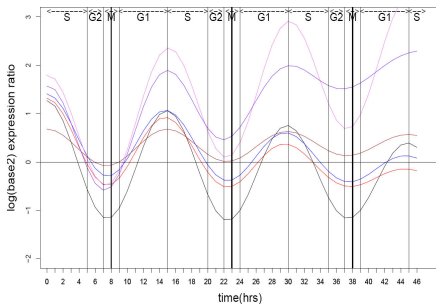
Limitations

- No consensus on the well-characterized genes, experimentation bias
- Similarity ignores the temporal structure, fixed a priori and unjustified

Aims

- Learn the best dissimilarity measure to classify or cluster genes expression profiles.
- Propose a well-founded set of reference genes.

Gene expression progression during the cell division process



- periodic profiles
- variation on: cell-cycle duration, initial amplitude
- amplitude attenuation
- tendency and drifts effects

Which metric for clustering and classifying genes expression profiles ?

- **Values-based metrics ?**
 - DTW, Euclidean distance, Manhattan distance, Fréchet distance, LCSS,...
- **Behavior-based metrics ?**
 - Pearson correlation coefficient, Qualitative distance (slope, derivative comparison), Kendall ratio, temporal correlation coefficient, etc...
- **Values and Behavior-based metrics ?**

Proximity measure specifications

- Genes expression are cyclic profiles
- Time of peak expression determines the cell-cycle phase assignment,
- Genes expression data may include tendency effects, amplitude attenuation,...

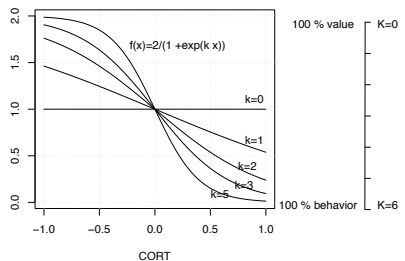
Induced constraints ...

- r should not include time warping
- δ_E is considered for values proximity measure
- *cort* is considered for behavior proximity measure

Values & Behavior-based metrics

$$D_k(S_1, S_2) = f(\text{CORT}(S_1, S_2)) \cdot \delta_E(S_1, S_2) \text{ with } f(x) = \frac{2}{1 + \exp(k x)} \quad , \quad k \geq 0$$

k : the contribution of values and of behavior to D



k learned during the classification or the clustering processes

Adaptive clustering: Partitioning around medoids approach

Motivations

- Use the PAM (Partitioning Around Medoids) algorithm to partition the set of genes into 5 clusters (5 cell cycle phases)
- More robust than k-means faced to outliers,
- Provide a more detailed analysis of the obtained partition
- Indicating for each object if it is (width silhouette):
 - well classified (i.e., genes well characterizing a cell cycle phase)
 - or lying on the boundary (genes involved in an inter-phase transition)

Algorithm steps

Learning D_k

- Perform the PAM algorithm based on D_k and for several values of (n, k) . Let $P_{n,k}$ be the obtained partition.
- note P_{n^*,k^*} the optimal partition according to two goodness criteria (asw, wb ratio)

Identification of the well-characterized genes

- Extract a kernel set of the p first genes maximizing the silhouette width,
- Identify the cell cycle phase of the kernel set,
- Assign each cluster to the cell cycle phase of it's kernel set.

Classification of the expressed cell cycle genes

- Assign each gene to the cell cycle phase of the cluster it belongs in.

Application specifications

- Analysis of experimental transcriptomic data from Human cancer cell line (Whitfield et al. 2002)
- Third experimentation
- The expression of 1099 periodically expressed genes through 48 instants covering 3 cell division cycles

<http://genome-www.stanford.edu/Human-CellCycle/Hela/>

Adaptive Clustering: Identification of the expressed cell cycle genes

Learning the most appropriate D_k

- Perform the PAM algorithm based on D_k for k varying in $[0, 6]$ and n in $[4, 10]$, $P_{n,k}$ the obtained partition for a given (n, k) .
- $(n^*, k^*) = \operatorname{argmax}_{n,k}(\operatorname{avgSil}(P_{n,k}))$, D_{k^*} is the most appropriate dissimilarity.

Identification of the well-characterized genes

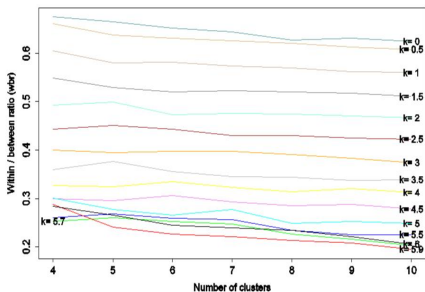
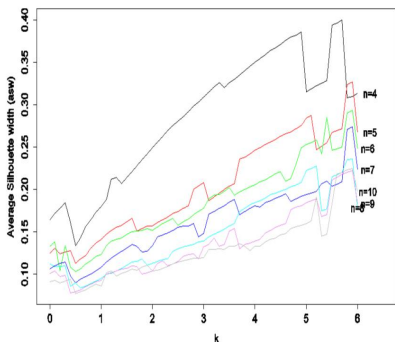
- Extract a kernel set of the N first genes maximizing the silhouette width,
- Identify the cell cycle phase of the kernel set,
- Assign each cluster to the cell cycle phase of it's kernel set.

Classification of the expressed cell cycle genes

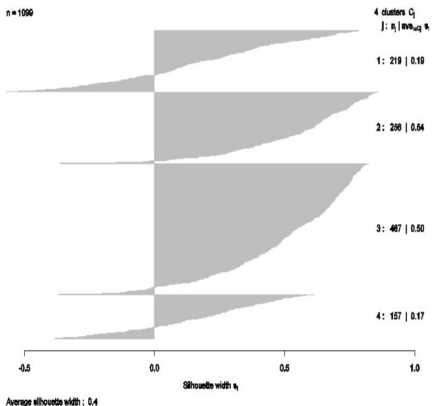
- Assign each gene to the cell cycle phase of the cluster it belongs in.

Learning the dissimilarity D_k through the partitioning process

- Average silhouette width (left plot) and the within/between ratio (right plot) of $P_{n,k}$, n from 4 to 10 and k from 0 to 6



Average silhouette width of the genes profiles partition



$n^* = 4, k^* = 5.7$: essentially the behavior separate well genes in 4 clusters
 Average Silhouette width=0.4: the clustering structure is reasonable

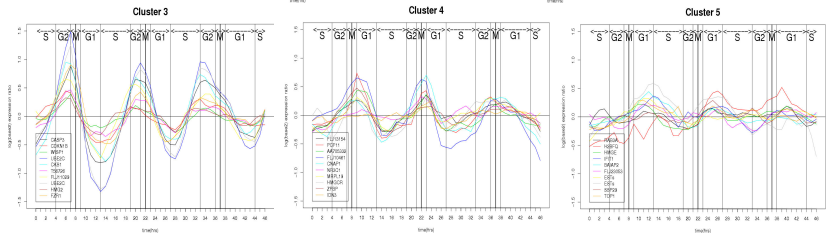
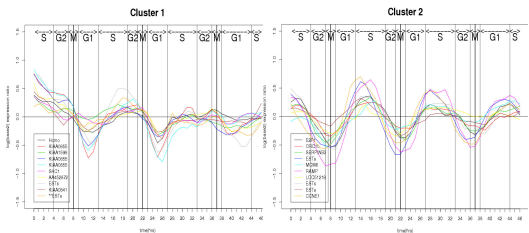
Identification of the well-characterized genes

-Extract a kernel set of the $n = 10$ first genes maximizing the silhouette width

Cluster Number	Gene Name	Whitfield Assignment	Gene Type	Neighbor Cluster	Silhouette width (sw)	High peaked phase
1	Homo	S	K	2	0.806	S
	KIAA0855	S	K	3	0.697	
	KIAA1598	S	K	2	0.688	
	KIAA0855	S	K	2	0.686	
	KIAA0855	S	K	3	0.681	
	SHC1	S	K	2	0.677	
	AA452872	S	K	3	0.674	
	ESTs	S	K	3	0.665	
	KIAA0841	S	K	2	0.658	
	**ESTs	S	K	3	0.635	
	RRM2	S	R	2	0.586	
	DHFR	S	R	2	0.315	
	RAD51	S	R	3	0.238	
2	E2F1*	G_1/S	K	1	0.832	G_1/S
	ORC1L	G_1/S	K	1	0.820	
	SERPINB3	G_1/S	K	1	0.82	
	ESTs	G_1/S	K	1	0.812	
	MCM6	G_1/S	K	1	0.812	
	RAMP	G_1/S	K	1	0.812	
	LOC51218	G_1/S	K	1	0.802	
	ESTs	G_1/S	K	1	0.794	
	ESTs	G_1/S	K	1	0.794	
	CCNE1	G_1/S	K/R	5	0.786	
	E2F1	G_1/S	R	1	0.775	
	CDC6	G_1/S	R	1	0.682	
	PCNA	G_1/S	R	1	0.625	
RFC4	S	R	1	0.526		
3	CASP3	G_2	K	4	0.811	G_2/M
	CDKN1B	G_2	K	4	0.807	
	WISP1	G_2	K	4	0.799	
	UBE2C	G_2	K	4	0.788	
	CKS1	G_2	K	4	0.784	
	T56726	G_2	K	4	0.770	
	FLJ11029	G_2	K	1	0.770	
	UBE2C	G_2	K	4	0.770	
	HMG2	G_2	K	4	0.768	
	FZR1	G_2	K	4	0.765	
	CCNF	G_2	R	4	0.757	
	TOP2A	G_2	R	4	0.669	
	CDC2	G_2	R	1	0.618	
STK15	G_2/M	R	4	0.478		
CCNA2	G_2	R	4	0.458		

4	FLJ13154	M/G_1	K	3	0.737	M/G_1
	PCF11	M/G_1	K	5	0.717	
	AA705332	G_2/M	K	5	0.695	
	FLJ10461	G_2/M	K	3	0.651	
	CNAP1	G_2/M	K	3	0.599	
	NR3C1	G_2	K	3	0.593	
	MRPL19	M/G_1	K	3	0.585	
	HMGCR	M/G_1	K	3	0.579	
	ZPBP	M/G_1	K	3	0.578	
	IDN3	G_2	K	3	0.576	
	RAD21	M/G_1	R	3	0.433	
	CDKN3	M/G_1	R	3	0.320	
	PTTG1	M/G_1	R	5	0.282	
BUB1	G_2/M	R	3	0.184		
VEGFC	M/G_1	R	3	0.148		
CCNB1	G_2/M	R	3	0.095		
PLK	G_2/M	R	3	0.003		
5	RAB3A	M/G_1	K	2	0.561	G_1
	H2BFQ	M/G_1	K	2	0.502	
	HMGE	M/G_1	K	4	0.489	
	IFIT1	M/G_1	K	2	0.484	
	BALAP2	G_1/S	K	2	0.478	
	FLJ23053	G_1/S	K	2	0.475	
	ESTs	M/G_1	K	4	0.429	
	ESTs	G_1/S	K	2	0.407	
SSP29	G_2/M	K	4	0.398		
TOP1	M/G_1	K	4	0.394		

Identify the cell cycle phase of each kernel set



$k^*=5.9$, (S, G₁/S, G₂/M, M/G₁, G₁)

Assignment of boundary genes

The kernel sets of the partition $P_{n^*=4, k=5, 7}$

Name	Ada-Assl	Neig	Sw	Final-Assl
MZF1	<i>S</i>	G_1/S	0.049	begin of <i>S</i>
TncRNA	<i>S</i>	G_1/S	0.049	begin of <i>S</i>
CAPS	<i>S</i>	G_1/S	0.041	begin of <i>S</i>
AURKE	<i>S</i>	G_2/M	0.039	G_2
ZFX	<i>S</i>	G_2/M	0.038	G_2
KATNA1	<i>S</i>	G_2/M	0.028	G_2
KBTBD2	<i>S</i>	G_2/M	0.026	G_2
CDKL5	<i>S</i>	G_2/M	0.02	G_2
TTC31	<i>S</i>	G_1/S	0.013	begin of <i>S</i>
LOC134121	<i>S</i>	G_2/M	0.012	G_2
UEL3	<i>S</i>	G_1	0.011	G_1/S
CDKN2C	<i>S</i>	G_2/M	0	G_2
REEF1	<i>S</i>	G_1/S	-0.012	begin of <i>S</i>
TOP2A	<i>S</i>	G_2/M	-0.023	G_2
MICA/HCP5	<i>S</i>	G_2/M	-0.039	G_2
CDH24	<i>S</i>	G_1/S	-0.041	begin of <i>S</i>
ABCC5	<i>S</i>	G_1/S	-0.044	begin of <i>S</i>
RECQL4	G_1/S	<i>S</i>	0.047	begin of <i>S</i>
SLC9A3	G_1/S	<i>S</i>	0.046	begin of <i>S</i>
FLJ13231	G_1/S	<i>S</i>	0.028	begin of <i>S</i>
ESTs	G_1/S	<i>S</i>	0.001	begin of <i>S</i>
EST	G_1/S	G_1	-0.019	end of G_1
HIST1H2AM	G_1/S	<i>S</i>	-0.045	begin of <i>S</i>
BAIAP2	G_1/S	G_1	-0.045	end of G_1
CRLF3	G_1/S	<i>S</i>	-0.05	begin of <i>S</i>

Name	Ada-Assl	Neig	Sw	Final-Assl
NR5A2	G_2/M	G_1	0.05	<i>M</i>
HERPUD2	G_2/M	G_1	0.045	<i>M</i>
AMD1	G_2/M	G_1	0.035	<i>M</i>
NIPBL	G_2/M	G_1	0.012	<i>M</i>
NFIC	G_2/M	G_1	0.008	<i>M</i>
ESTs	G_2/M	G_1	0.006	<i>M</i>
ChGn	G_2/M	G_1	0.003	<i>M</i>
ECLAF1	G_2/M	<i>S</i>	0.001	G_2
WWC1	G_2/M	G_1	-0.003	<i>M</i>
HLA-DOA	G_2/M	G_1	-0.012	<i>M</i>
AGPAT3	G_2/M	G_1	-0.015	<i>M</i>
C20orf199	G_2/M	G_1	-0.017	<i>M</i>
SLC39A10	G_2/M	G_1	-0.02	<i>M</i>
LARP1	G_2/M	G_1	-0.024	<i>M</i>
ANP32B	G_2/M	G_1	-0.026	<i>M</i>
ABHD10	G_2/M	<i>S</i>	-0.029	G_2
FXR1	G_2/M	G_1	-0.032	<i>M</i>
ZNF207	G_1	G_2/M	0.05	<i>M</i>
HSPA2	G_1	G_2/M	0.048	<i>M</i>
PPP2CA	G_1	G_2/M	0.044	<i>M</i>
CEP350	G_1	G_2/M	0.017	<i>M</i>
OC146517	G_1	G_2/M	0.013	<i>M</i>
SAP30BP	G_1	<i>S</i>	0.009	G_1/S
DR1	G_1	G_2/M	0.007	<i>M</i>
TMEM132A	G_1	G_2/M	0.002	<i>M</i>
W85890	G_1	G_2/M	-0.007	<i>M</i>
PCF11	G_1	G_2/M	-0.021	<i>M</i>
DNAJA1	G_1	G_2/M	-0.022	<i>M</i>
TSC22	G_1	G_2/M	-0.023	<i>M</i>
EST	G_1	G_2/M	-0.024	<i>M</i>
EST	G_1	G_2/M	-0.024	<i>M</i>
EST	G_1	G_2/M	-0.024	<i>M</i>

Metrics efficiency comparison: Random-Periods model for periodically expressed genes

The sinusoid function characterizing the expected periodic expression of a cell-cycle gene g (Liu et al. (2004)):

$$f(t, \theta_g) = a_g + b_g t + \frac{K_g}{\sqrt{2\pi}} \int_{-\infty}^{+\infty} \cos\left(\frac{2\pi t}{T \exp(\sigma z)} + \Phi_g\right) \exp\left(-\frac{z^2}{2}\right) dz,$$

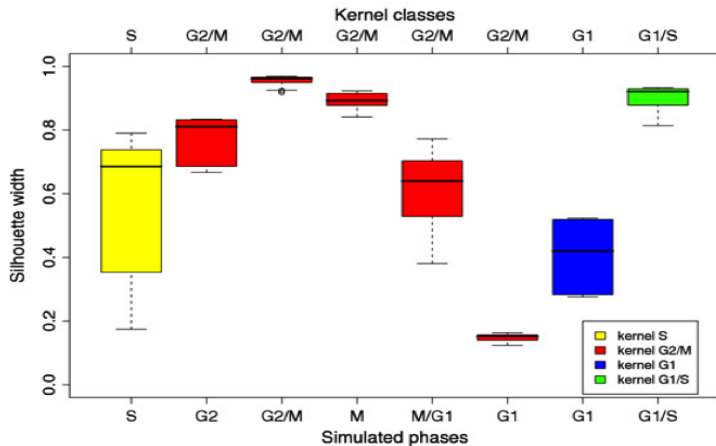
where θ_g is explicitly $(K_g, T, \sigma, \Phi_g, a_g, b_g)$, specific to each gene g

- K_g : initial amplitude of the periodic expression pattern
- T : cell-cycle duration
- σ : governs the rate of attenuation in amplitude
- Φ_g : corresponds to the cell-cycle phase during which the gene undergoes its peak level of transcription
- a_g and b_g : account for any drift (intercepts and slopes, respectively) in a gene's background expression level

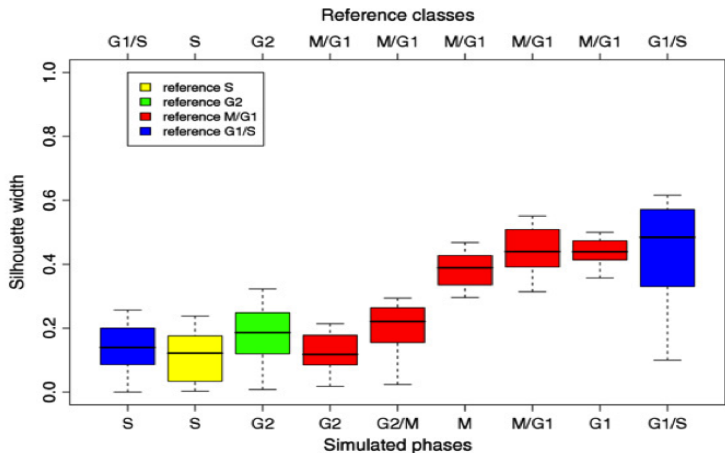
Dissimilarity efficiency for classifying gene expression profiles

Simulated classes	Kernel classes				Reference classes				
	S	G ₂ /M	G ₁	G ₁ /S	S	G ₂	G ₂ /M	M/G ₁	G ₁ /S
S	100	0	0	0	29	0	0	0	71
G ₂	0	100	0	0	0	65	0	35	0
G ₂ /M	0	100	0	0	0	0	0	100	0
M	0	100	0	0	0	0	0	100	0
M/G ₁	0	100	0	0	0	0	0	100	0
G ₁	0	27	73	0	0	0	0	100	0
G ₁ /S	0	0	0	100	0	0	0	0	100

Dissimilarity efficiency for classifying gene expression profiles



Dissimilarity efficiency for classifying gene expression profiles

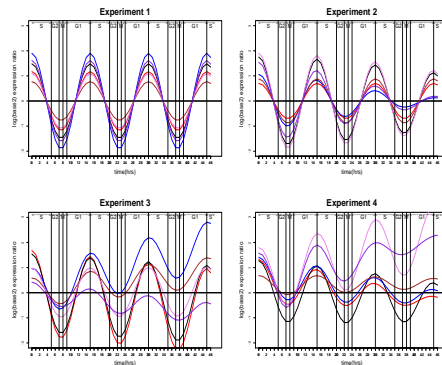


Validation Protocol

- A simulation study based on 20000 genes expression profiles,
- genes are equally generated from 5 classes (five cell-cycle phases)
- each gene expression is observed through 3 cell-cycles on 47 instants,
- Four experiments are simulated (500 genes /experiment)
- 10 samples are generated for each experiment

$$T=15, \Phi_g=(0, 5.190, 3.823, 3.278, 2.459)$$

Experiment number	K_g	σ	b_g	a_g
1	[0.34, 1.33]	0	0	0
2	[0.34, 1.33]	[0, 0.115]	0	0
3	[0.34, 1.33]	0	[-0.05, 0.05]	[0, 0.8]
4	[0.34, 1.33]	[0, 0.115]	[-0.05, 0.05]	[0, 0.8]



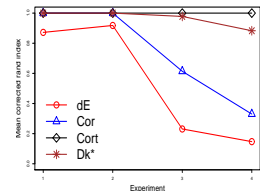
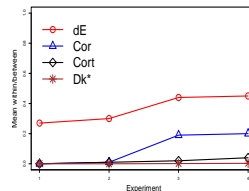
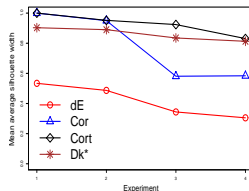
Proximity measures efficiency for clustering genes expression profiles

For each experiment $j \in \{1, 2, 3, 4\}$ and for each measure δ_E , COR, and CORT

- a PAM algorithm is used to partition each sample S_{ij} , $i \in \{1, \dots, 10\}$ into 5 clusters (5 cell cycle phases)
- Three goodness criteria: the average silhouette width (asw), the within/between ratio (wbr), and the corrected Rand index (RI)

For the adaptive dissimilarity D_k

- a PAM algorithm is performed for k varying in $[0, 6]$,
- select $P_{k^*}^{ij}$, with $k^* = \operatorname{argmax}_k(\operatorname{avgSil}(P_k^{ij}))$,



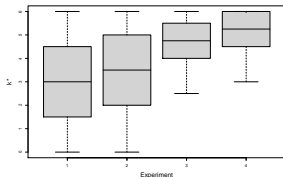
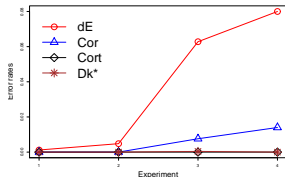
Proximity measure efficiency for classifying gene expression profiles

For each experiment $j \in \{1, 2, 3, 4\}$ and for each measure δ_E , COR, and CORT

- A 10-NN algorithm is performed to classify each sample S_{ij} , $i \in \{1, \dots, 10\}$
- The misclassification error rate is evaluated

For the adaptive dissimilarity D_k

- the 10-NN algorithm is performed for k varying in $[0, 6]$,
- $C_{k^*}^{ij}$ is selected with $k^* = \operatorname{argmin}_k(\operatorname{ErrorRate}(C_k^{ij}))$,



References

- Liu, D., Umbach, D.M., Peddada, S.D., Li, L., Crockett, P.W., Weinberg, C.R., 2004. A random-periods model for expression of cell-cycle genes. *Proceedings of the National Academy of Sciences of the USA* 101, 72407245.
- Liu, X., Lee, S., Casella, G., Peter, G.F., 2008. Assessing agreement of clustering methods with gene expression microarray data. *Computational Statistics and Data Analysis* 52 (12), 53565366.
- Maharaj, E.A., 2000. Cluster of time series. *Journal of Classification* 17, 297314.
- Oates, T., Firoiou, L., Cohen, P.R., 1999. Clustering time series with Hidden Markov models and dynamic time warping. In: *Proceedings of the 6th IJCAI-99, Workshop on Neural, Symbolic and Reinforcement Learning Methods for Sequence Learning*, Stockholm, pp. 1721.
- Peddada, S.D., Lobenhofer, L., Li, L., Afshari, C., Weinberg, C., Umbach, D., 2003. Gene selection and clustering for time-course and dose-response microarray experiments using order-restricted inference. *Bioinformatics* 19, 834841.
- Pihur, V., Datta, S., Datta, S., 2007. Weighted rank aggregation of cluster validation measures: A Monte Carlo cross-entropy approach. *Bioinformatics* 23 (13), 16071615.
- Whitfield, M.L., Sherlock, G., Murray, J.I., Ball, C.A., Alexander, K.A., Matese, J.C., Perou, C.M., Hurt, M.M., Brown, P.O., Botstein, D., 2002. Identification of genes periodically expressed in the human cell cycle and their expression in tumors molecular. *Biology of the Cell* 13, 19772000.