A Comparison of Clustering Techniques for Gene Expression Data

Michiel de Hoon, Seiya Imoto, Satoru Miyano Human Genome Center, University of Tokyo





The 10th International Conference on Intelligent Systems for Molecular Biology August 3 – 7, 2002, Edmonton, Canada



ISMB 2002 2002.08.04

Genes / microarrays can be clustered based on the similarity between their expression profiles



What do we mean by similarity?

Several distance measures are commonly used.

Unnormalized distance functions:

Minkowski metrics: Euclidean distance, city-block distance, ...

Mahalanobis distance

Requires some normalization of the expression data

Normalized distance functions:

Usually based on some correlation coefficient *r* via d = 1 - rPearson correlation, Spearman rank correlation, Kendall's τ Can be applied to normalized and unnormalized data Information about the magnitude of changes in the gene expression

is ignored

These distances do not satisfy the triangle inequality!

Intermediate:

Angle between expression data vectors ("uncentered correlation" *)

* M. B. Eisen, P. T. Spellman, P. O. Brown, D. Botstein, "Cluster analysis and display of genome-wide expression patterns", *Proc. Natl. Acad. Sci. USA*, Vol. 95, pp. 14863-14868, 1998

Which distance function should be used?





東京大学



Pearson:

The Pearson correlation coefficient is used most commonly: Alon *et al.*, PNAS **96**, 6745-6750 (1999); Heyer *et al.*, Genome Research **9**, 1106-1115 (1999); Ewing & Claverie, Proc. PSB **5**, 430-441 (2000); Zhu & Zhang, Proc. PSB **5**, 479-490 (2000).

The angle between expression data vectors ("uncentered correlation") is also widely used: Eisen *et al.*, PNAS **95**, 14863-14868 (1998).



Euclid:

However, the Euclidean distance has been reported to perform better: Slonim *et al.*, RECOMB, 263-272 (2000).



Three types of clustering algorithms are commonly used



• Hierarchical clustering:

Initially each gene forms one cluster

On each step, join the two closest clusters

Continue until only one cluster remains

Draw a tree diagram

Several types of hierarchical clustering exist:

Simple linkage, maximum linkage, centroid linkage, average linkage

• *k*-means clustering:

Create *k* clusters; assign genes randomly to a cluster Iterate: reassign genes to the cluster whose centroid is the closest Continue until no further reassignment takes place This algorithm should be repeated many times to find the optimal solution. The optimal solution has the smallest sum of distances to the cluster centroids.

• Self-Organizing Maps:

Based on a topology of clusters

Which clustering algorithm should be used?



ISMB 2002 2002.08.04

Various clustering programs are available



Cluster/TreeView (Michael Eisen, Berkeley/Stanford):
Widely used, especially by biologists
Mainly focused on hierarchical clustering (simple, maximum, centroid linkage)
but also includes *k*-means clustering and 1D Self-Organizing Maps
Good documentation
Recently, the source code (C++) was (partially) released.
We have modified Cluster/TreeView to include an improved *k*-means clustering algorithm.
Available from our website (http://bonsai.ims.u-tokyo.ac.jp/~mdehoon/software/cluster/index.html)

Tamayo's GeneCluster program:
 Offers 2D Self-Organizing Maps
 Written in Java with nice graphics
 The source code is not available: impossible to check the algorithm or to improve the code

• The C Clustering Library:

Numerical routines for hierarchical and *k*-means clustering, as well as 2D Self-Organizing Maps. Released under the GNU Lesser Public License. Also available for use with Python.

Statistical programs such as S and R



Different clustering methods produce different clustering solutions



We consider a real-life example to compare different clustering methods

Gene expression levels were measured for 66 patients, each afflicted by one of five disease types:

A (14 patients)
B (7 patients)
C (7 patients)
D (10 patients)
E (28 patients)

We only include genes that satisfy the following two conditions:

The gene expression levels deviate by at least one standard deviation;

• Less than 10% of the data values are missing.

This resulted in 1224 remaining genes.

Pairwise centroid linkage clustering with the Euclidean distance produces one large cluster ...





Can we learn anything from this clustering solution?





... while pairwise complete linkage clustering with the Euclidean distance performs slightly better





... and *k*-means clustering with the Euclidean distance performing better still







k-means with the Euclidean distance; 100,000 trials



Pairwise centroid linkage clustering with the uncentered correlation is used most commonly



東京大学



ISMB 2002 2002.08.04

The University of Tokyo, Human Genome Center

Pairwise complete linkage clustering with the uncentered correlation gives a better result ...







... but again *k*-means gives the clearest result



This is similar to the result found with the centered correlation using *k*-means





ISMB 2002 2002.08.04







ISMB 2002 2002.08.04

Our second example considers gene clustering



Cyanobacterium Synechocystis sp PCC 6803

(Cyanobase website; Kazusa DNA Research Institute)

Analysis of Cyanobacterial Gene Expression during Acclimation to High Light Yukako Hihara, Ayako Kamei, Minoru Kanehisa, Aaron Kaplan, Masahiko Ikeuchi The Plant Cell, Vol. 13, 793-806, April 2001

Acclimation from low light to high light:

 Short-term processes take place within several minutes (State transitions, protective energy dissipation, changes in the efficiency of energy transfer from the harvesting complex to photosystem II, formation of nonfunctional PII reaction centers)

 Long-term processes take place in hours or days (Changes in the composition, function, and structure of the photosynthetic apparatus and other photosynthesis-related components)



The time-dependent response of Cyanobacterium to high light was measured



東京大学

Change the light intensity from low to high and examine the response as a function of time



Preprocessing gene expression data



Remove noisy genes

Expression levels close to the background level cannot be distinguished from noise and should therefore not be included in the analysis.

Apply global normalization

Scale the cy5 and cy3 levels by dividing them by their sum over all genes. This avoids errors due to different rates of hybridization.

Remove genes that fail the t-test

Perform a *t*-test for each gene to see if its expression ratio is significantly different from unity. If not, that gene was not significantly affected by the experimental manipulation and should be removed from the analysis.

Fit linear spline functions to the gene expression levels

Using a linear spline function instead of the expression data directly reduces the effect of noise present in the data.

These steps have been implemented in Python/C routines.



ISMB 2002 2002.08.04

Clustering reveals patterns of gene expression



k-means clustering, using Euclidean distance



k-means clustering using Euclidean distance: Cluster I





sll0322	hypF	transcriptional regulatory protein HypF
sll0927	metX	S-adenosylmethionine synthetase
sll1028	ccmK	carbon dioxide concentrating mechanism protein
sll1029	ccmK	carbon dioxide concentrating mechanism protein
sll1031	ccmM	carbon dioxide concentrating mechanism protein
slr1963		hypothetical protein
slr2075	groES	10kD chaperonin
sll0416	groEL-2	60kD chaperonin 2

10,000 repetitions; solution found 111 times



ISMB 2002 2002.08.04

k-means clustering using Euclidean distance: Cluster II





sll0170	dnaK	DnaK protein
sll0430	htpG	heat shock protein
sll0521	ndhG	NADH dehydrogenase subunit 6
slr1280	ndhK	NADH dehydrogenase subunit NdhK
sir0011 sir1350 sir1604 sir1641 sil0814 sir0476 sir1687	rbcX desA ftsH clpB	unknown function fatty acid desaturase cell division protein ClpB protein hypothetical protein hypothetical protein hypothetical protein



The University of Tokyo, Human Genome Center

東京大学

k-means clustering using Euclidean distance: Cluster III





sll1097 rps7 30S ribosomal protein S7
sll1260 rps2 205 ribosomal protoin 52
sll1743 rpl11 50S ribosomal protein L11
sll1745 rpl10 50S ribosomal protein L10
sll1799 rpl3 50S ribosomal protein L3
sll1801 rpl23 50S ribosomal protein L23
sll1802 rpl2 50S ribosomal protein L2
sll1804 rps3 30S ribosomal protein S3
sll1807 rpl24 50S ribosomal protein L24
sll1809 rps8 30S ribosomal protein S8
sll1810 rpl6 50S ribosomal protein L6
sll1816 rps13 30S ribosomal protein S13



ISMB 2002 2002.08.04

k-means clustering using Euclidean distance: Cluster IV



東京大学



010000	alac	ATD demondent Cln protococ regulatory cubunit
SILUUZU	cipe	ATP-dependent Cip protease regulatory subunit
sll0144	pyrH	uridine monophosphate kinase
sll0185		unknown
sll0262	desD(des6)	delta-6 desaturase
sll0414		hypothetical protein
sll0519	ndhA	NADH dehydrogenase subunit 1
sll0587	pykF	pyruvate kinase
sll0680	pstS or pho	phosphate-binding periplasmic protein
sli0854		unknown
sll1002	ycf22	hypothetical protein
sll1327	atpC	ATP synthase g subunit
sll1614	pma1	cation-transporting ATPase (E1-E2 ATPase)
slr0208		unknown
slr0452	ilvD	dihydroxyacid dehydratase
slr0642		integral membrane protein
slr0839	hemH	ferrochelatase
slr1237	codA	cytosine deaminase
slr1469	rnpA	protein subunit of ribonuclease P (RNase P)
slr1513		unknown
slr1557		unknown
slr1718		unknown
00*2554		unknown
5512334		UNKNOWN



k-means clustering using Euclidean distance: Cluster V

. . .





SII0258	pspv	cytochrome c550
sll0617	im30	chloroplast membrane-associated 30 kD protein
sll0851	psbC	photosystem II CP43 protein
sll0901	purE or add	e6 phosphoribosylaminoimidazole carboxylase
sll1091		391 aa (43 kD) bacteriochlorophyll synthase subunit
sll1214	PNIL34 or A	AT103 phytochrome-regulated gene
sll1306		hypothetical protein
sll1471	cpcG	phycobilisome rod-core linker polypeptide CpcG
sll1626	lexA	SOS function regulatory protein
sll1694	hofG or ho	<i>pG</i> general secretion pathway protein G
sll1712		DNA binding protein HU
ssl3093	срсD	phycocyanin associated linker protein
slr0151		hypothetical protein
slr0329	xylR	xylose repressor
slr0374		cell division cycle protein
slr1128		erthyrocyte band 7 integral embrane protein,
		protein 7.2B, stomatin
slr1348	cysE	serine acetyltransferase
slr1459	apcF	phycobilisome core component
slr1545	rpoE	RNA polymerase sigma-E factor
slr1793	talB	transaldolase
slr1853		hypothetical protein
slr1856		anti-sigma B factor antagonist
ssl1533		hypothetical protein
		——————————————————————————————————————
city of T		man Canama Cantar 界以人子

ISMB 2002 2002.08.04

k-means clustering using Euclidean distance: Cluster VI



東京大学



sll1577	срсВ	phycocyanin b subunit
sll1578	срсА	phycocyanin a subunit
sll1579	срсС	phycocyanin associated linker protein
sll1580	срсС	phycocyanin associated linker protein
slr2051	срсG	phycobilisome rod-core linker polypeptide
sll0819 slr0737 slr1655	psaF psaD psaL	photosystem I subunit III photosystem I subunit II photosystem I subunit XI
slr1834 slr1835	psaA psaB	P700 apoprotein subunit la P700 apoprotein subunit lb
slr1855		hypothetical protein
slr1986 slr2067	арсВ арсА	allophycocyanin b chain allophycocyanin a chain

Note: The Euclidean distance distinguishes cluster V and VI.

ISMB 2002 2002.08.04

k-means clustering using Pearson Correlation: Cluster I



sll0262 sll0322 sll0854 sll0927	desD(des6 hypF metX)delta-6 desaturase <mark>(cluster IV)</mark> transcriptional regulatory protein HypF unknown (cluster IV) S-adenosylmethionine synthetase
sll1029 sll1031	ccmK ccmM	carbon dioxide concentrating mechanism protein carbon dioxide concentrating mechanism protein
slr0208 slr1718 slr1963		unknown (cluster IV) unknown (cluster IV) hypothetical protein
slr2075 sll0416	groES groEL-2	10kD chaperonin 60kD chaperonin 2

1,000,000 repetitions; solution found 37 times



ISMB 2002 2002.08.04

k-means clustering using Pearson correlation: Cluster II



s 0020 s 0170 s 0414 s 0430	clpC dnaK htpG	ATP-dependent CIp protease regulatory subunit (cluster IV) DnaK protein hypothetical protein (cluster IV) heat shock protein
sll0519 sll0521 slr1280	ndhA ndhG ndhK	NADH dehydrogenase subunit 1 (cluster IV) NADH dehydrogenase subunit 6 NADH dehydrogenase subunit NdhK
sll1002 sll1028 sll1614 slr0011 slr0642 slr1350 slr1513	ycf22 ccmK pma1 rbcX desA	hypothetical protein (cluster IV) carbon dioxide concentrating mechanism protein (cluster I) cation-transporting ATPase (E1-E2 ATPase) (cluster IV) unknown function integral membrane protein (cluster IV) fatty acid desaturase unknown (cluster IV)
slr1604 slr1641 sll0814 slr0476 slr1687	ftsH clpB	cell division protein ClpB protein hypothetical protein hypothetical protein hypothetical protein



k-means clustering using Pearson correlation: Cluster III



sll0144 sll0587 sll0680 sll1327	pyrH pykF pstS or phoS atpC	uridine monophosphate kinase <mark>(cluster IV)</mark> pyruvate kinase <mark>(cluster IV)</mark> phosphate-binding periplasmic protein precursor (PBP) <mark>(cluster IV</mark> ATP synthase g subunit <mark>(cluster IV)</mark>
SII1096 SII1097 SII1260 SII1743 SII1745 SII1745 SII1799 SII1801 SII1802 SII1804 SII1807 SII1809 SII1810 SII1810	rps12 rps7 rps2 rpl11 rpl10 rpl3 rpl23 rpl23 rpl2 rps3 rpl24 rps8 rpl6 rps13	30S ribosomal protein S12 30S ribosomal protein S7 30S ribosomal protein S2 50S ribosomal protein L11 50S ribosomal protein L10 50S ribosomal protein L3 50S ribosomal protein L23 50S ribosomal protein L2 30S ribosomal protein S3 50S ribosomal protein S3 50S ribosomal protein S8 50S ribosomal protein S8 50S ribosomal protein S8
slr1237 slr1469 slr1557 ssr2554	codA rnpA	cytosine deaminase (cluster IV) protein subunit of ribonuclease P (RNase P) (cluster IV) unknown (cluster IV) unknown (cluster IV)



k-means clustering using Pearson correlation: Cluster IV



sll0185unknownslr0839hemHferrochelatase



k-means clustering using Pearson correlation: Cluster V



im30 psbC	chloroplast membrane-associated 30 kD protein photosystem II CP43 protein
cpcG	phycobilisome rod-core linker polypeptide CpcG
hofG or hopG	general secretion pathway protein G
	cell division cycle protein
psaD	photosystem I subunit II <mark>(cluster VI)</mark>
	erthyrocyte band 7 integral embrane protein, protein 7.2B, stomatin
cysE	serine acetyltransferase
talB	transaldolase
	hypothetical protein
psaA	P700 apoprotein subunit la <mark>(cluster VI)</mark>
psaB	P700 apoprotein subunit Ib (cluster VI)
	im30 psbC cpcG hofG or hopG psaD cysE talB psaA psaB



k-means clustering using Pearson correlation: Cluster VI



sll0258	psbV	cytochrome c550 (cluster IV)
sll0901	purE or ade6	phosphoribosylaminoimidazole carboxylase (cluster IV)
sll1091		391 aa (43 kD) bacteriochlorophyll synthase subunit (cluster IV)
sll1214	<i>PNIL34</i> or <i>AT103</i>	phytochrome-regulated gene (cluster IV)
sll1306		hypothetical protein (cluster IV)
sll1577	срсВ	phycocyanin b subunit
sll1578	срсА	phycocyanin a subunit
sll1579	cpcC	phycocyanin associated linker protein
sll1580	срсС	phycocyanin associated linker protein
ssl3093	срсD	phycocyanin associated linker protein (cluster V)
sll1626	lexA	SOS function regulatory protein (cluster V)
sll1712		DNA binding protein HU <mark>(cluster V)</mark>
slr2051	срсG	phycobilisome rod-core linker polypeptide
sll0819	psaF	photosystem I subunit III
slr1655	psaL	photosystem I subunit XI
slr0151		hypothetical protein (cluster V)
slr0329	xylR	xylose repressor (cluster V)
slr0452	ilvD	dihydroxyacid dehydratase (cluster IV)
slr1459	apcF	phycobilisome core component (cluster V)
slr1545	rpoE	RNA polymerase sigma-E factor (cluster V)
slr1853		hypothetical protein (cluster V)
slr1855		hypothetical protein
slr1856	_	anti-sigma B factor antagonist (cluster V)
sir1986	арсВ	allophycocyanin b chain
sir2067	арсА	aliophycocyanin a chain

ISMB 2002 2002.08.04

k-means clustering using uncentered correlation: Cluster I



sll0262	desD(des6)delta-6 desaturase (cluster IV)		
sil0814 sil0322 sil0854 sil0927	hypF metX	transcriptional regulatory protein unknown (cluster IV) S-adenosylmethionine synthetase	HypF e
sll1028 sll1029 sll1031	ccmK ccmK ccmM	carbon dioxide concentrating med carbon dioxide concentrating med carbon dioxide concentrating med	chanism protein chanism protein chanism protein
slr1963		hypothetical protein	
slr2075 sll0416	groES groEL-2	10kD chaperonin 60kD chaperonin 2	
sir0011 sir0208 sir0642 sir1513 sir1718	rbcX	unknown function (cluster II) unknown (cluster IV) integral membrane protein (cluster unknown (cluster IV) unknown (cluster IV)	er IV)

100,000 repetitions; solution found 131 times



k-means clustering using uncentered correlation: Cluster II



sll0519 sll0521 slr1280	ndhA ndhG ndhK	NADH dehydrogenase subunit 1 <mark>(cluster IV)</mark> NADH dehydrogenase subunit 6 NADH dehydrogenase subunit NdhK
slr1350 slr1604 slr1641 slr0476 slr1687	desA ftsH clpB	fatty acid desaturase cell division protein ClpB protein hypothetical protein hypothetical protein
sll0020	clpC	ATP-dependent Clp protease regulatory subunit (cluster IV)
sll0170	dnaK	DnaK protein
sll0414		hypothetical protein (cluster IV)
sll0430	htpG	heat shock protein
sll1614	pma1	cation-transporting ATPase (E1-E2 ATPase) (cluster IV)



k-means clustering using uncentered correlation: Cluster III



東京大学

sll0144	pyrH	uridine monophosphate kinase <mark>(cluster IV)</mark>
sll0587	pykF	pyruvate kinase <mark>(cluster IV)</mark>
sll0680	pstS or phoS	phosphate-binding periplasmic protein precursor (PBP) <mark>(cluster IV</mark>
sll1327	atpC	ATP synthase g subunit <mark>(cluster IV)</mark>
SII1096	rps12	30S ribosomal protein S12
SII1097	rps7	30S ribosomal protein S7
SII1260	rps2	30S ribosomal protein S2
SII1743	rpl11	50S ribosomal protein L11
SII1745	rpl10	50S ribosomal protein L10
SII1745	rpl3	50S ribosomal protein L3
SII1799	rpl23	50S ribosomal protein L23
SII1801	rpl23	50S ribosomal protein L2
SII1802	rpl2	30S ribosomal protein S3
SII1804	rps3	50S ribosomal protein S3
SII1807	rpl24	50S ribosomal protein S8
SII1809	rps8	50S ribosomal protein S8
SII1810	rpl6	50S ribosomal protein S8
SII1810	rps13	50S ribosomal protein S13
slr1237 slr1469 slr1557 ssr2554	codA rnpA	cytosine deaminase (cluster IV) protein subunit of ribonuclease P (RNase P) (cluster IV) unknown (cluster IV) unknown (cluster IV)



k-means clustering using uncentered correlation: Cluster IV

sll0185		unknown
sll1002	ycf22	hypothetical protein
slr0839	hemH	ferrochelatase



k-means clustering using uncentered correlation: Cluster V



sll0258	psbV	cytochrome c550
sll0901	purE or ade6	phosphoribosylaminoimidazole carboxylase
sll1091		391 aa (43 kD) bacteriochlorophyll synthase subunit
sll1214	PNIL34 or AT103	phytochrome-regulated gene
sll1626	lexA	SOS function regulatory protein
sll1712		DNA binding protein HU
slr0452	ilvD	dihydroxyacid dehydratase <mark>(cluster IV)</mark>
slr1459	apcF	phycobilisome core component
slr1545	rpoE	RNA polymerase sigma-E factor
slr1853		hypothetical protein
slr1855		hypothetical protein (cluster VI)
slr1856		anti-sigma B factor antagonist
slr1986	арсВ	allophycocyanin b chain (cluster VI)
slr2067	, apcA	allophycocyanin a chain (cluster VI)
_		



k-means clustering using uncentered correlation: Cluster VI



Hierarchical clustering using Euclidean distance: Cluster I



sll1028	сстК	carbon dioxide concentrating mechanism protein
sll1029	сстК	carbon dioxide concentrating mechanism protein
sll1031	сстМ	carbon dioxide concentrating mechanism protein
slr1963		hypothetical protein
slr2075	groES	10kD chaperonin
sll0416	groEL-2	60kD chaperonin 2



The University of Tokyo, Human Genome Center

東京大学

Hierarchical clustering using Euclidean distance: Cluster II/IV



sll0170	dnaK	DnaK protein
sll0430	htpG	heat shock protein
sll0521	ndhG	NADH dehydrogenase subunit 6
slr1280	ndhK	NADH dehydrogenase subunit NdhK
sir0011 sir1350 sir1604 sir1641 sil0814 sir0476 sir1687	rbcX desA ftsH clpB	unknown function fatty acid desaturase cell division protein ClpB protein hypothetical protein hypothetical protein hypothetical protein





Hierarchical clustering using Euclidean distance: Cluster III



sll1096	rps12	30S ribosomal protein S12
sll1097	rps7	30S ribosomal protein S7
sll1260	rps2	30S ribosomal protein S2
sll1743	rpl11	50S ribosomal protein L11
sll1745	rpl10	50S ribosomal protein L10
sll1799	rpl3	50S ribosomal protein L3
sll1801	rpl23	50S ribosomal protein L23
sll1802	rpl2	50S ribosomal protein L2
sll1804	rps3	30S ribosomal protein S3
sll1807	rpl24	50S ribosomal protein L24
sll1809	rps8	30S ribosomal protein S8
sll1810	rpl6	50S ribosomal protein L6
sll1816	rps13	30S ribosomal protein S13
sll0322 sll0927	hypF metX	transcriptional regulatory protein HypF (cluster I) S-adenosylmethionine synthetase (cluster I)



Hierarchical clustering using Euclidean distance: Cluster II/IV



sll0020	clpC	ATP-dependent Clp protease regulatory subunit
sll0144	pyrH	uridine monophosphate kinase
sll0185		unknown
sll0262	desD(des6)	delta-6 desaturase
sll0414		hypothetical protein
sll0519	ndhA	NADH dehydrogenase subunit 1
sll0587	pykF	pyruvate kinase
sll0680	pstS or phoS	phosphate-binding periplasmic protein precursor (PBP)
sll0854		unknown
sll1002	ycf22	hypothetical protein
sll1327	atpC	ATP synthase g subunit
sll1614	pma1	cation-transporting ATPase (E1-E2 ATPase)
slr0208		unknown
slr0642		integral membrane protein
slr0839	hemH	ferrochelatase
slr1237	codA	cytosine deaminase
slr1469	rnpA	protein subunit of ribonuclease P (RNase P)
slr1513		unknown
slr1557		unknown
slr1718		unknown
ssr2554		unknown



Hierarchical clustering using Euclidean distance: Cluster V



sll0258	psbV	cytochrome c550
sll0617	im30	chloroplast membrane-associated 30 kD protein
sll0851	psbC	photosystem II CP43 protein
sll0901	purE or ade6	phosphoribosylaminoimidazole carboxylase
sll1306		hypothetical protein
sll1471	cpcG	phycobilisome rod-core linker polypeptide CpcG
sll1626	lexA	SOS function regulatory protein
sll1694	hofG or hopG	general secretion pathway protein G
sll1712		DNA binding protein HU
slr0151		hypothetical protein
slr0329	xylR	xylose repressor
slr0374	-	cell division cycle protein
slr0452	ilvD	dihydroxyacid dehydratase (cluster IV)
slr1128		erthyrocyte band 7 integral embrane protein, protein 7.2B, stomatin
slr1348	cysE	serine acetyltransferase
slr1545	rpoE	RNA polymerase sigma-E factor
slr1793	talB	transaldolase
slr1853		hypothetical protein
slr1856		anti-sigma B factor antagonist
ssl1533		hypothetical protein



Hierarchical clustering using Euclidean distance: Cluster V/VI



sll1091 sll1214 ssl3093 slr1459 slr1855	PNIL34 or AT103 cpcD apcF	391 aa (43 kD) bacteriochlorophyll synthase subunit (cluster V) phytochrome-regulated gene (cluster V) phycocyanin associated linker protein (cluster V) phycobilisome core component (cluster V) hypothetical protein (cluster VI)
slr1986	арсВ	allophycocyanin b chain <mark>(cluster VI)</mark>
slr2067	арсА	allophycocyanin a chain <mark>(cluster VI)</mark>



Hierarchical clustering using Euclidean distance: Cluster VI



sll1577	срсВ	phycocyanin b subunit
sll1578	срсА	phycocyanin a subunit
sll1579	срсС	phycocyanin associated linker protein
sll1580	срсС	phycocyanin associated linker protein
slr2051	срсG	phycobilisome rod-core linker polypeptide
sll0819	psaF	photosystem I subunit III
slr0737	psaD	photosystem I subunit II
slr1655	psaL	photosystem I subunit XI
slr1834	psaA	P700 apoprotein subunit la
slr1835	psaB	P700 apoprotein subunit lb



Conclusions



- Different clustering methods yield different clustering solutions.
- *k*-means clustering is more reliable than hierarchical clustering.
- The *k*-means clustering algorithm may take a huge amount of time.
- We therefore recommend to use hierarchical clustering for preliminary analyses, and to check the results with k-means.
- Among the hierarchical clustering routines, complete linkage clustering was found to perform best.
- We found that the uncentered correlation performs better than the Euclidean distance.
- As we show only two examples, we cannot provide a final answer to the question which clustering method is best. We therefore suggest that in general, clustering results should be verified by using different methods. Differences in the clustering solutions should be included in publications.



Reprint requests



In case you would like to receive a reprint of this poster, please leave us your contact information below.



