

ΠΟΛΥΤΕΧΝΕΙΟ ΚΡΗΤΗΣ

Σχολή Ηλεκτρολόγων Μηχανικών και Μηχανικών Υπολογιστών



Διπλωματική Εργασία

**ΜΕΘΟΔΟΙ ΔΙΠΛΗΣ ΚΑΤΗΓΟΡΙΟΠΟΙΗΣΗΣ ΣΤΗΝ
ΑΝΑΛΥΣΗ ΓΟΝΙΔΙΩΝ**

“Biclustering Methods in Gene Analysis”

ΑΝΔΡΙΤΣΟΥ ΔΑΦΝΗ

Επιβλέπων Καθηγητής: Καθηγητής Ζερβάκης Μιχάλης

Εξεταστική Επιτροπή: Καθηγητής Ζερβάκης Μιχάλης

Καθηγητής Μπάλας Κωνσταντίνος

Αναπληρωτής Καθηγητής Λαγουδάκης Μιχαήλ

Χανιά, Ιούνιος 2017

Ευχαριστίες

Θα ήθελα να ευχαριστήσω

Πρωτίστως, τον καθηγητή μου κύριο Ζερβάκη Μιχάλη, για την πολύτιμη βοήθεια και καθοδήγησή του καθ' όλη τη διάρκεια της εκπόνησης της διπλωματικής μου εργασίας. Για την εμπιστοσύνη που μου έδειξε δίνοντας μου τη δυνατότητα να εκπονήσω την πτυχιακή μου εργασία στον επιστημονικό τομέα που επιθυμούσα. Επίσης, θα ήθελα να τον ευχαριστήσω για τη διάθεση του να με βοηθήσει και να μου λύσει οποιαδήποτε απορία οποιαδήποτε στιγμή το χρειάζομαι.

Τους καθηγητές κύριο Μπάλα Κωνσταντίνο και κύριο Λαγουδάκη Μιχαήλ για τη συνεισφορά τους ως μέλη της εξεταστικής επιτροπής.

Την Δρ. Μπέη Αικατερίνη για τον πολύτιμο χρόνο που αφιέρωσε από την αρχή μέχρι το τελικό στάδιο της διπλωματικής μου εργασίας καθώς και για την απλόχερη στήριξη και ενθάρρυνσή της προς το άτομό μου.

Την Δρ. Obermayr για την παροχή του dataset και τον διδάκτορα Σφακιανάκη Στέλιο για την βοήθειά του στην επεξεργασία του dataset.

Την Αλεβυζάκη Ανδρονίκη για την βοήθειά της στο αρχικό στάδιο της εργασίας μου.

Θα ήθελα να ευχαριστήσω θερμά όλους τους φίλους μου, κυρίως την Έλενα, την Άννα, τη Ρένια και τη Σέβη που πίστεψαν σε εμένα και με ενθάρρυναν σε κάθε στάδιο των σπουδών μου.

Ιδιαίτερες ευχαριστίες θέλω να εκφράσω προς τους γονείς μου, Ευαγγελία και Κωσταντίνο για την διαχρονική συμπαράστασή τους και την υλική και ηθική στήριξη των επιλογών μου.

Περίληψη

Οι μικροσυστοιχίες DNA αποτελούν μία από τις πιο διαδεδομένες πειραματικές μεθόδους ταυτόχρονης ανάλυσης του τρόπου έκφρασης χιλιάδων γονιδίων σε διαφορετικά δείγματα ή σε διαφορετικές συνθήκες. Το γεγονός αυτό τις καθιστά ιδανικό εργαλείο για την ανάλυση και μελέτη καρκινικών ιστών, με στόχο την εύρεση των μοριακών μηχανισμών που διέπουν την ογκογένεση σε διάφορους τύπους καρκίνου και την περαιτέρω κατανόηση της νόσου. Τα δεδομένα γονιδιακής έκφρασης αναπαρίστανται με τη μορφή πινάκων όπου οι γραμμές αντιστοιχούν σε γονίδια και οι στήλες σε διάφορες πειραματικές συνθήκες. Στόχος των διάφορων τεχνικών ομαδοποίησης αποτελεί η εξαγωγή σημαντικών βιολογικών πληροφοριών που αφορούν ομάδες γονιδίων κάτω από συγκεκριμένες συνθήκες. Στην παρούσα εργασία επιλέγεται η εφαρμογή της μεθόδου διπλής κατηγοριοποίησης δεδομένων γονιδιακής έκφρασης που προέρχονται από δείγματα κυτταρικών σειρών τεσσάρων τύπων καρκίνου - του μαστού, των ωοθηκών, του ενδομητρίου και του τραχήλου της μήτρας - συγκεντρωμένα σ' έναν πίνακα δεδομένων. Η μεθοδολογία που προτάθηκε, περιλάμβανε την εφαρμογή του αλγορίθμου διπλής κατηγοριοποίησης των Cheng και Church, εφαρμόζοντας σε αυτόν μια σειρά βελτιώσεων ούτως ώστε να εξάγουμε ταυτόχρονα ομάδες γονιδίων με ομοιόμορφη συμπεριφορά αλλά και συγκεκριμένο εύρος τιμών. Ο αριθμός γονιδίων (ανιχνευτών) που περιείχε ο πίνακας, ανερχόταν σε 33096, γεγονός που έκανε την μελέτη και την ανάλυση της συμπεριφοράς των γονιδίων δύσκολη. Γι' αυτό το λόγο, μειώσαμε το πλήθος των γονιδίων προς ανάλυση στα 1000 γονίδια ανά καρκινικό τύπο, καθώς επίσης, εστίασαμε στην εξέταση των τριών πρώτων ομάδων διπλής κατηγοριοποίησης (Biclusters) που προέκυπταν κάθε φορά. Με αυτή τη διπλή κατηγοριοποίηση **υποχώρου δεδομένων**, καταφέραμε να δημιουργήσουμε ομάδες γονιδίων οι οποίες εμφάνιζαν γονίδια με όμοια συμπεριφορά κατά μήκος των κυτταρικών σειρών και παράλληλα το εύρος των τιμών τους ήταν μικρό, κάνοντας έτσι τις ομάδες πιο συνεκτικές και ικανές να εξάγουν σημαντικές βιολογικές πληροφορίες, όπως η συμμετοχή των γονιδίων που τις απαρτίζουν σε συγκεκριμένες βιολογικές διεργασίες (π.χ. μεταβολισμός, βιοσύνθεση μακρομορίων, αναπαραγωγή) και εξειδικευμένα μοριακά μονοπάτια (π.χ. οξειδωτική φωσφορυλίωση, μεταφορά RNA, ενδοκύττωση).

Abstract

DNA microarrays comprise one of the most widespread experimental methods. This fact makes them a perfect tool for analysis and evaluation of cancer tissues, with the aim to find those molecular mechanisms governing the formation of tumors in different types of cancer and a further understanding of this disease. Molecular expression data are depicted in tables where the lines represent the genes and the columns represent the different experimental conditions.

The aim of the different grouping techniques is to elicit crucial biological information regarding gene groups under specific circumstances. In this essay the method selected was the double categorization of gene expression data coming from samples of cell arrays of four different types of cancer (breast, ovarian, endometrial, cervical) gathered on a data table. The methodology proposed included the application of the Cheng & Church biclustering algorithm, by applying a series of revisions to it so as to elicit gene groups with similar patterns and a specific value range simultaneously.

The number of genes (probes) included in the table rose to 33096, a fact that posed a serious burden to the study and analysis of the genes. For that reason, we brought the number of genes analyzed down to 1000 per type of cancer as well as focusing on studying the first three Biclusters each time. With this double categorization of data subspace we managed to create gene groups which showed genes with similar patterns along the cell arrays and, at the same time, their value range was small, making these groups more coherent and offering us the opportunity to elicit crucial biological information, such as the involvement of their genes in concrete biological processes (e.g., metabolism, macromolecular biosynthesis, reproduction) and specific molecular pathways (e.g., oxidative phosphorylation, RNA transfer, endocytosis).

ΠΕΡΙΕΧΟΜΕΝΑ

1 Εισαγωγή	11
1.1. Ανάλυση του Προβλήματος	11
Κίνητρο και Στόχοι	12
1.2. Υφιστάμενη Γνώση	13
1.3. Δομή της Εργασίας	15
Καινοτομία της Εργασίας	15
2 Θεωρητικό Υπόβαθρο	17
2.1. Βιολογικό Υπόβαθρο	17
2.1.1. Γενετικό υλικό και Χρωμοσώματα	18
2.1.2. DNA	19
2.1.3. Το μοντέλο της διπλής έλικας	20
2.1.4. RNA	21
2.1.5. Πρωτεΐνες	22
2.1.6. Γονιδιακή Έκφραση	23
2.2. Εισαγωγή στην τεχνολογία των Μικροσυστοιχιών (microarrays)	25
2.3. Ανάλυση γονιδιακής έκφρασης από μικροσυστοιχίες	28
2.4. ΒιοΠληροφορική	29
2.5. Αλγόριθμοι Κατηγοριοποίησης	31
Κατηγοριοποίηση Βάσει Γονιδίων	32
Κατηγοριοποίηση Βάσει Δειγμάτων	32
Κατηγοριοποίηση Υποδιαστημάτων (Subspace Clustering)	33
2.6. Τεχνικές Κατηγοριοποίησης	34
Ιεραρχική Ομαδοποίηση (Hierarchical Clustering)	34
Συσταδοποίηση δυο-τρόπων (Biclustering)	36
3 Αλγόριθμοι και Τεχνικές Διπλής Κατηγοριοποίησης (Biclustering)	38
3.1. Τύποι Διπλής Κατηγοριοποίησης	38
3.2. Δομή Διπλής Κατηγοριοποίησης	42
3.3. Ποιοτικά Μέτρα Διπλής Κατηγοριοποίησης	42
4 Δεδομένα και Προτεινόμενη Μεθοδολογία	52
4.1. Ανάπτυξη Μεθοδολογίας	53
4.2. Δεδομένα	53
4.2.1. Επιλογή αλγορίθμου Cheng και Church	54
4.2.2. Εισαγωγή στη χρήση του αλγορίθμου Cheng και Church	54
4.3. Προτεινόμενη Μεθοδολογία	61

5 Αποτελέσματα	67
5.1. Αποτελέσματα biclustering από κυτταρικές σειρές καρκίνου του μαστού	69
5.2. Αποτελέσματα biclustering από κυτταρικές σειρές του καρκίνου του τραχήλου της μήτρας.....	69
5.3. Αποτελέσματα biclustering από κυτταρικές σειρές καρκίνου των Ωοθηκών	71
5.4. Αποτελέσματα biclustering από κυτταρικές σειρές καρκίνου του Ενδομητρίου	73
5.5.Αποτελέσματα Biclustering από Κυτταρικές Σειρές Καρκίνου του Μαστού για 33096 Γονίδια	75
5.6. Αποτελέσματα Biclustering από Κυτταρικές Σειρές Καρκίνου του Τραχήλου της Μήτρας για 33096 Γονίδια.....	78
6 Αξιολόγηση Των Αποτελεσμάτων	81
Στατιστική Επικύρωση Αποτελεσμάτων	87
Οπτικοποίηση Αποτελεσμάτων Διπλής Κατηγοριοποίησης.....	88
7 Συμπεράσματα	97
7.1. Συμπεράσματα	97
7.2. Μελλοντικές Επεκτάσεις	97
Βιβλιογραφία	100
Παράρτημα	101

ΠΙΝΑΚΑΣ ΠΕΡΙΕΧΟΜΕΝΩΝ ΕΙΚΟΝΩΝ-ΔΙΑΓΡΑΜΜΑΤΩΝ

Κεφάλαιο 1

Εικόνα 1.1.	Πρόβλημα NP-Complete.....	13
-------------	---------------------------	----

Κεφάλαιο 2

Εικόνα 2.1.	Από το κύτταρο στο DNA.....	17
Εικόνα 2.2.	Γονίδιο και Ανθρώπινος καρυότυπος γυναίκας και άνδρα.....	18
Εικόνα 2.3.	Τα τέσσερα διαφορετικά δεοξυριβονουκλεοτίδια του DNA.....	19
Εικόνα 2.4.	Το μοντέλο της διπλής έλικας του DNA	20
Εικόνα 2.5.	Αμινοξύ – δομική μονάδα πρωτεϊνών.....	22
Εικόνα 2.6.	Μηχανισμός σύνθεσης πρωτεϊνών στα ευκαρυωτικά κύτταρα	23
Εικόνα 2.7.	Μικροσυστοιχίες DNA	26
Εικόνα 2.8.	Σάρωση και υβριδοποίηση μικροσυστοιχίας που περιέχει περισσότερα από 5000 γονίδια.....	28
Εικόνα 2.9.	Αύξηση βιολογικών δεδομένων από το 1971 έως το 2015.....	31
Εικόνα 2.10.	Υποδιάστημα συστάδων A & B.....	33
Εικόνα 2.11.	Ιεραρχική ομαδοποίηση (Hierarchical clustering)	35
Εικόνα 2.12.	Ενωτική (Agglomerative) και Διαχωριστική (Divisive) Ομαδοποίηση.....	36
Εικόνα 2.13.	Είδη αλγορίθμων διπλής κατηγοριοποίησης.....	37

Κεφάλαιο 3

Εικόνα 3.1.	Τύποι Biclusters.....	38
Εικόνα 3.2.	Biclusters με σταθερές τιμές.....	39
Εικόνα 3.3.	Biclusters με σταθερές τιμές σε γραμμές και στήλες.....	39
Εικόνα 3.4.	Biclusters με συνεκτικές τιμές (μοντέλο πρόσθεσης, μοντέλο πολλαπλασιασμού).....	40
Εικόνα 3.5.	Biclusters με συνεκτικές εξελίξεις.....	41
Εικόνα 3.6.	Δομές Bicluster.....	42

Κεφάλαιο 4

Εικόνα 4.1.	Σύνολο δεδομένων εκφρασμένο ως πίνακα πραγματικών τιμών.....	53
Εικόνα 4.1.α.	Πρώτη Εκτέλεση του Αλγορίθμου (<i>didj</i>) για 1000 Γονίδια του Καρκίνου του Μαστού	61
Εικόνα 4.1.β.	Επιλεκτική Εμφάνιση Γονιδίων μετά την Πρώτη εκτέλεση του αλγορίθμου (<i>didj</i>) για 1000 Γονίδια του Καρκίνου του Μαστού.....	62
Εικόνα 4.1.γ.	Εστιασμένη Εμφάνιση Γονιδίων μετά την Πρώτη εκτέλεση του αλγορίθμου (<i>didj</i>) για 1000 Γονίδια του Καρκίνου του Μαστού.....	62
Εικόνα 4.2.	Δεύτερη Εκτέλεση του Αλγορίθμου (<i>di_new</i>) για 1000 Γονίδια του Καρκίνου του Μαστού.....	66

Εικόνα 4.3.	Τρίτη Εκτέλεση του Αλγορίθμου (<i>di_new,dj_new</i>) για 1000 Γονίδια του Καρκίνου του Μαστού.....	66
Κεφάλαιο 5		
Εικόνα 5.1.	Συγκεντρωτική Παρουσίαση των Τριών εκτελέσεων του Αλγορίθμου Cheng και Church για 1000 Γονίδια από τις Κυτταρικές Σειρές Καρκίνου του Μαστού.....	69
Εικόνα 5.2.α.	Πρώτη Εκτέλεση του Αλγορίθμου (<i>didj</i>) για 1000 Γονίδια του Καρκίνου του Τραχήλου της Μήτρας.....	70
Εικόνα 5.2.β.	Δεύτερη Εκτέλεση του Αλγορίθμου (<i>di_new</i>) για 1000 Γονίδια του Καρκίνου του Τραχήλου της Μήτρας.....	70
Εικόνα 5.2.γ.	Τρίτη Εκτέλεση του Αλγορίθμου (<i>di_new,dj_new</i>) για 1000 Γονίδια του Καρκίνου του Τραχήλου της Μήτρας.....	70
Εικόνα 5.3.	Συγκεντρωτική Παρουσίαση των Τριών εκτελέσεων του Αλγορίθμου Cheng και Church για 1000 Γονίδια από τις Κυτταρικές Σειρές Καρκίνου του Τραχήλου της Μήτρας	71
Εικόνα 5.3.α.	Πρώτη Εκτέλεση του Αλγορίθμου (<i>di,dj</i>) για 1000 Γονίδια του Καρκίνου των Ωοθηκών..	72
Εικόνα 5.3.β.	Δεύτερη Εκτέλεση του Αλγορίθμου (<i>di_new,dj</i>) για 1000 Γονίδια του Καρκίνου των Ωοθηκών	72
Εικόνα 5.3.γ.	Τρίτη Εκτέλεση του Αλγορίθμου (<i>di_new,dj_new</i>) για 1000 Γονίδια του Καρκίνου των Ωοθηκών	72
Εικόνα 5.4.	Συγκεντρωτική Παρουσίαση των Τριών εκτελέσεων του Αλγορίθμου Cheng και Church για 1000 Γονίδια από τις Κυτταρικές Σειρές Καρκίνου των Ωοθηκών..	73
Εικόνα 5.4.α.	Πρώτη Εκτέλεση του Αλγορίθμου (<i>didj</i>) για 1000 Γονίδια του Καρκίνου του Ενδομητρίου	74
Εικόνα 5.4.β.	Δεύτερη Εκτέλεση του Αλγορίθμου (<i>di_new</i>) για 1000 Γονίδια του Καρκίνου του Ενδομητρίου	74
Εικόνα 5.4.γ.	Τρίτη Εκτέλεση του Αλγορίθμου (<i>di_new,dj_new</i>) για 1000 Γονίδια του Καρκίνου του Ενδομητρίου..	74
Εικόνα 5.5.	Συγκεντρωτική Παρουσίαση των Τριών εκτελέσεων του Αλγορίθμου Cheng και Church για 1000 Γονίδια από τις Κυτταρικές Σειρές Καρκίνου του Ενδομητρίου..	75
Εικόνα 5.5.α.	Πρώτη Εκτέλεση του Αλγορίθμου (<i>didj</i>) για 33096 Γονίδια του Καρκίνου του Μαστού ..	76
Εικόνα 5.5.β.	Δεύτερη Εκτέλεση του Αλγορίθμου (<i>di_new</i>) για 33096 Γονίδια του Καρκίνου του Μαστού.....	76
Εικόνα 5.5.γ.	Τρίτη Εκτέλεση του Αλγορίθμου (<i>di_new,dj_new</i>) για 33096 Γονίδια του Καρκίνου του Μαστού.....	77
Εικόνα 5.6.	Συγκεντρωτική Παρουσίαση των Τριών εκτελέσεων του Αλγορίθμου Cheng και Church για 33096 Γονίδια από τις Κυτταρικές Σειρές Καρκίνου του Μαστού	77

Εικόνα 5.6.α. Πρώτη Εκτέλεση του Αλγορίθμου (di,dj) για 33096 Γονίδια του Καρκίνου του Τραχήλου της Μήτρας.....	78
Εικόνα 5.6.β. Δεύτερη Εκτέλεση του Αλγορίθμου (di_new,dj) για 33096 Γονίδια του Καρκίνου του Τραχήλου της Μήτρας.....	79
Εικόνα 5.6.γ. Τρίτη Εκτέλεση του Αλγορίθμου (di_new,dj_new) για 33096 Γονίδια του Καρκίνου του Τραχήλου της Μήτρας.....	79
Εικόνα 5.7. Συγκεντρωτική Παρουσίαση των Τριών εκτελέσεων του Αλγορίθμου Cheng and Church για 33096 Γονίδια από τις Κυτταρικές Σειρές Καρκίνου του Τραχήλου της Μήτρας.....	80

Κεφάλαιο 6

Εικόνα 6.1. Συγκεντρωτική Παρουσίαση των Τριών εκτελέσεων του Αλγορίθμου Cheng και Church για 1000 Γονίδια από τις Κυτταρικές Σειρές Καρκίνου του Μαστού (Bicluster1)	82
Εικόνα 6.1.α. Πρώτη Εκτέλεση του Αλγορίθμου (di,dj) για 1000 Γονίδια του Καρκίνου του Μαστού- Ανάδειξη των μεγαλύτερων και μικρότερων τιμών των Καρκινικών Κυτταρικών Σειρών που Συμμετέχουν	83
Εικόνα 6.1.β. Πρώτη Εκτέλεση του Αλγορίθμου (di,dj) για 1000 Γονίδια του Καρκίνου του Μαστού Αναδεικνύοντας την Ανομοιομορφία μεταξύ των γονιδίων.....	84
Εικόνα 6.1.γ. Πρώτη Εκτέλεση του Αλγορίθμου (di,dj) για 1000 Γονίδια του Καρκίνου του Μαστού Αναδεικνύοντας το Εύρος της Απόστασης μεταξύ των γονιδίων	84
Εικόνα 6.2. Δεύτερη Εκτέλεση του Αλγορίθμου (di_new,dj) για 1000 Γονίδια του Καρκίνου του Μαστού - Ανάδειξη των Καρκινικών Κυτταρικών Σειρών που Συμμετέχουν	84
Εικόνα 6.3.α. Τρίτη Εκτέλεση του Αλγορίθμου (di_new,dj_new) για 1000 Γονίδια του Καρκίνου του Μαστού - Ανάδειξη των Καρκινικών Κυτταρικών Σειρών που Συμμετέχουν	85
Εικόνα 6.3.β. Πρώτη Εκτέλεση του Αλγορίθμου (di,dj) για 1000 Γονίδια του Καρκίνου του Μαστού .	86
Εικόνα 6.3.γ. Δεύτερη Εκτέλεση του Αλγορίθμου (di_new,dj) για 1000 Γονίδια του Καρκίνου του Μαστού	86
Εικόνα 6.3.δ. Τρίτη Εκτέλεση του Αλγορίθμου (di_new,dj_new) για 1000 Γονίδια του Καρκίνου του Μαστού	86
Εικόνα 6.4.α. Πρώτη Εκτέλεση του Αλγορίθμου ($didj$) για 1000 Γονίδια του Καρκίνου του Μαστού- Ανάδειξη των Καρκινικών Κυτταρικών Σειρών που Συμμετέχουν.....	86
Εικόνα 6.4.β. Πρώτη Εκτέλεση του Αλγορίθμου ($didj$) για 1000 Γονίδια του Καρκίνου του Μαστού	87
Εικόνα 6.4.γ. Δεύτερη Εκτέλεση του Αλγορίθμου (di_new,dj) για 1000 Γονίδια του Καρκίνου του Μαστού	87
Εικόνα 6.4.δ. Τρίτη Εκτέλεση του Αλγορίθμου (di_new,dj_new) για 1000 Γονίδια του Καρκίνου του Μαστού	87
Εικόνα 6.5. Στατιστική Επικύρωση Αποτελεσμάτων	

ΠΙΝΑΚΑΣ ΠΕΡΙΕΧΟΜΕΝΩΝ ΠΙΝΑΚΩΝ

Κεφάλαιο 1

Πίνακας 1.1. Εφαρμογές διπλής κατηγοριοποίησης τα τελευταία 14 χρόνια.....	14
--	----

Κεφάλαιο 3

Πίνακας 3.1. Μοντέλο Πρόσθεσης.....	40
Πίνακας 3.2. Μοντέλο Πολλαπλασιασμού.....	41

Κεφάλαιο 4

Πίνακας 4.1. Συγκεντρωτικά Αποτελέσματα Αριθμού Γονιδίων για Κάθε Ένα Από τα Τρία Biclusters που μελετήθηκαν.....	60
---	----

Κεφάλαιο 5

Πίνακας 5.1. Συγκεντρωτικός πίνακας αποτύπωσης προτεινόμενης μεθοδολογίας.....	68
--	----

Κεφάλαιο 6

Πίνακας 6.1. Παρουσίαση τριών πρώτων Biclusters του καρκίνου του μαστού	82
Πίνακας 6.2. Παρουσίαση των εμπλουτισμένων βιολογικών μονοπατιών ($p \leq 0.05$) για τα τρία Biclusters κατά την διπλή κατηγοριοποίηση του συνόλου των δεδομένων για τον καρκίνο του μαστού	94
Πίνακας 6.3. Παρουσίαση των εμπλουτισμένων βιολογικών διεργασιών ($p \leq 0.05$) για το δεύτερο Bicluster κατά την διπλή κατηγοριοποίηση του συνόλου των δεδομένων για τον καρκίνο του μαστού	95

ΠΙΝΑΚΑΣ ΠΕΡΙΕΧΟΜΕΝΩΝ ΣΧΗΜΑΤΩΝ

Σχήμα 4.1. Σχεδιάγραμμα ροής	59
Σχήμα 4.2. Σχηματική Αναπαράσταση της Συμπεριφοράς των Γονιδίων στις Τρεις Εκτελέσεις του Αλγορίθμου Cheng και Church.....	64
Σχήμα 6.5. Οπτικοποίηση Αποτελεσμάτων Διπλής Κατηγοριοποίησης Υποχώρου Δεδομένων – Καρκίνος του Μαστού.....	89
Σχήμα 6.6. Οπτικοποίηση Αποτελεσμάτων Διπλής Κατηγοριοποίησης Υποχώρου Δεδομένων – Καρκίνος του Τραχήλου της Μήτρας.....	90
Σχήμα 6.7. Οπτικοποίηση Αποτελεσμάτων Διπλής Κατηγοριοποίησης Υποχώρου Δεδομένων – Καρκίνος του Ενδομητρίου	91
Σχήμα 6.8. Οπτικοποίηση Αποτελεσμάτων Διπλής Κατηγοριοποίησης Υποχώρου Δεδομένων – Καρκίνος του Ενδομητρίου	92
Σχήμα 6.9. Οπτικοποίηση Αποτελεσμάτων Διπλής Κατηγοριοποίησης Συνόλου Δεδομένων – Καρκίνος του Μαστού.....	93

Εισαγωγή

1.1. Ανάλυση του Προβλήματος

Η ανάπτυξη της τεχνολογίας τα τελευταία χρόνια έχει οδηγήσει σε πολύ σημαντικές εξελίξεις, οι οποίες έχουν επηρεάσει σε σημαντικό βαθμό όλες τις πτυχές της ζωής μας. Μία από αυτές τις πτυχές, είναι και αυτή της βιοϊατρικής, δηλαδή του κλάδου που προσπαθεί να αναπτύξει της δυνατότητες της ιατρικής μέσω της εξέλιξης της τεχνολογίας.

Στα πλαίσια της ανάπτυξης νέων τεχνολογιών έχουν παραχθεί εργαλεία γονιδιωματικής απαραίτητα για την κατανόηση και τη διαλεύκανση των κυτταρικών λειτουργιών, των βιοχημικών οδών και των ρυθμιστικών μηχανισμών με βάση το πρότυπο της γονιδιακής έκφρασης (gene expression profile) ενός κυττάρου. Μία από αυτές τις τεχνολογίες είναι οι μικροσυστοιχίες DNA, οι οποίες αποτελούν ένα σημαντικό εργαλείο τα τελευταία χρόνια αλλάζοντας τον τρόπο με τον οποίο η βιοϊατρική και άλλες επιστήμες αντιμετωπίζουν διάφορες βιολογικές ερωτήσεις και δίνοντας ώθηση στη μετάφραση της γονιδιακής έρευνας στην κλινική. Από την άλλη μεριά, η τεχνολογία αυτή βασίζεται σε μεγάλο βαθμό σε συγκέντρωση μεγάλου όγκου δεδομένων, τα οποία απαιτούν επεξεργασία. Τις τελευταίες δεκαετίες παρουσιάζεται ραγδαία αύξηση του όγκου των βιολογικών δεδομένων προς επεξεργασία και η ανάγκη για τον μετασχηματισμό της αδόμητης αυτής πληροφορίας σε γνώση. Αυτός ήταν ο κύριος λόγος δημιουργίας του τομέα της βιοπληροφορικής (bioinformatics). Έτσι, καθώς αναπτύσσεται αυτός ο κλάδος, κρίνεται απαραίτητο να χρησιμοποιηθούν διάφοροι σχετικοί αλγόριθμοι, ώστε να είναι εφικτά τα καλύτερα αποτελέσματα.

Όπως αναφέρθηκε, η τεχνολογία των μικροσυστοιχιών DNA επιτρέπει να παρατηρήσουμε ταυτοχρόνως χιλιάδες γονίδια και να καθορίσουμε ποια από αυτά είναι διαφορετικά εκφρασμένα σε ένα συγκεκριμένο τύπο κυττάρου ή υπό συγκεκριμένες συνθήκες. Επιπλέον, τα προφίλ γονιδιακής έκφρασης ενός νοσούντος κυττάρου/ιστού, σε σύγκριση με τους φυσιολογικούς μάρτυρες, μπορεί να προάγουν την κατανόηση της παθολογίας της νόσου και τον εντοπισμό νέων θεραπευτικών σημείων παρέμβασης, τη βελτίωση της διάγνωσης και τη διασαφήνιση της πρόγνωσης. Στην παρούσα εργασία επιλέχθηκε να εφαρμοστεί η τεχνική διπλής κατηγοριοποίησης σε δεδομένα γονιδιακής έκφρασης που προέρχονται από δείγματα καρκινικών κυτταρικών σειρών -καρκίνου του μαστού, των ωοθηκών, του ενδομητρίου και του τραχήλου της μήτρας- με στόχο την βέλτιστη ανάλυση των γονιδίων που είναι συγκεντρωμένα σε ένα πίνακα δεδομένων. Πραγματώθηκε, λοιπόν, η μελέτη αλγορίθμων ευαίσθητων στην πληροφορία και όχι στο θόρυβο και με ιδιαίτερη λεπτομέρεια και ακρίβεια στην ομαδοποίηση.

Κίνητρο και Στόχοι

Οι πίνακες γονιδιακής έκφρασης έχουν εκτενώς αναλυθεί χρησιμοποιώντας ομαδοποίηση (clustering) σε μία από τις δύο διαστάσεις:

- Τη γονιδιακή διάσταση
- Τη διάσταση της κατάστασης

Πιο αναλυτικά :

- α) ανάλυση της έκφρασης προτύπων γονιδίων, συγκρίνοντας τις γραμμές του πίνακα και,
- β) ανάλυση της έκφρασης προτύπων γονιδίων από δείγματα, συγκρίνοντας τις στήλες του πίνακα.

Οι αλγόριθμοι ομαδοποίησης

- Μπορούν να εφαρμοσούν ομαδοποίηση είτε σε γονίδια ή σε δείγματα, έτσι η ανάλυση κατευθύνεται σε μία συγκεκριμένη πτυχή του υπό μελέτη συστήματος.
- Συνήθως, επιδιώκουν μία ασύνδετη κάλυψη του συνόλου των στοιχείων, «απαιτώντας» ότι κανένα γονίδιο ή δείγμα δεν θα ανήκει σε περισσότερες από μια ομάδες (clusters).

Πολλά ενεργά πρότυπα είναι κοινά σε μία ομάδα γονιδίων μόνο κάτω από συγκεκριμένες πειραματικές καταστάσεις. Ανακαλύπτοντας τέτοια τοπική έκφραση προτύπων μπορεί να είναι το κλειδί για την αποκάλυψη πολλών γενετικών μονοπατιών που δεν είναι εμφανή με άλλο τρόπο. Είναι ως εκ τούτου επιθυμητό να προχωρήσουμε πέρα από το παράδειγμα ομαδοποίησης και να αναπτύξουμε προσεγγίσεις που θα είναι σε θέση να ανακαλύπτουν τοπικά πρότυπα σε δεδομένα μικροσυστοιχιών.

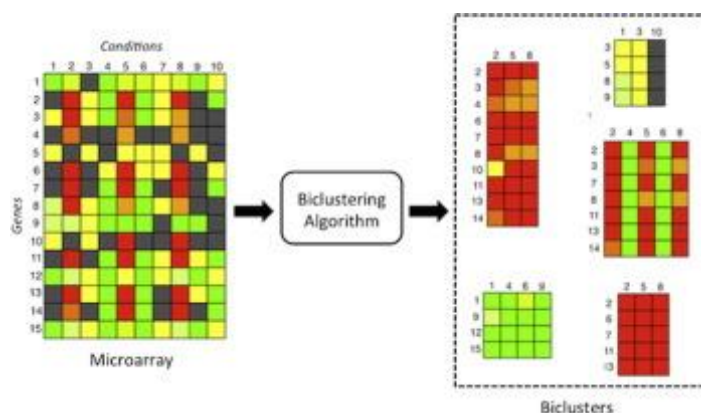
Η διπλή κατηγοριοποίηση αποτελεί μια πιο σύγχρονη μέθοδο ταυτόχρονης ομαδοποίησης γραμμών και στηλών, που αφορά στην ομαδοποίηση δεδομένων με βάση την ένταση της γονιδιακής έκφρασης (υψηλές ή χαμηλές αριθμητικές τιμές) των γονιδίων ή το εύρος της γονιδιακής έκφρασης (διακύμανση των αριθμητικών τιμών) των γονιδίων υπό συγκεκριμένες συνθήκες. Καινοτομία της συγκεκριμένης διπλωματικής αποτέλεσε η στόχευση και στις δύο αυτές παραμέτρους που ο συνδυασμός τους παρέχει αποτελέσματα με μεγάλη στατιστική σημαντικότητα.

Παρατηρούμε επίσης ότι ένας φυσικός τρόπος απεικόνισης μιας ομάδας Biclusters συνίσταται στην ανάθεση ενός διαφορετικού χρώματος για κάθε Biclusters και στην αναδιάταξη των σειρών και των στηλών του πίνακα δεδομένων, έτσι ώστε να αποκτήσουμε έναν πίνακα δεδομένων με έγχρωμα μπλοκς, όπου κάθε μπλοκ αντιπροσωπεύει ένα Biclusters (Εικόνα 1.1). Το Biclusters πρόβλημα μπορεί να διατυπωθεί ως εξής: Δοθέντος ενός πίνακα δεδομένων M , κατασκευάσε μια ομάδα Biclusters B_{opt} που σχετίζονται με το M έτσι ώστε:

$$f(B_{opt}) = \max_{B \in BC(M)} f(B)$$

Όπου f είναι μια αντικειμενική συνάρτηση που μετρά την ποιότητα, δηλ. το βαθμό συνοχής μιας ομάδας και τα $BC(M)$ είναι το σύνολο όλων των δυνατών ομάδων Biclusters που σχετίζονται με το M . Αυτό το πρόβλημα είναι NP-hard.[1],[2]

Ειδικότερα, η διπλή κατηγοριοποίηση (biclustering) και οι παραλλαγές της ανήκουν στην κλάση μη-ντετερμινιστικής πολωνυμικής (Non Deterministic Polynomial, NP) πολυπλοκότητας, δηλαδή είναι NP-πλήρης. Αυτό σημαίνει ότι τα προβλήματα που περιλαμβάνει μπορούν να επαληθευτούν εντός πολωνυμικού χρόνου (Εικόνα 1.1). Η πολυπλοκότητα του προβλήματος διπλής κατηγοριοποίησης είναι άρρηκτα συνδεδεμένη με το εκάστοτε πρόβλημα, αλλά και την εκτιμήτρια συνάρτηση που χρησιμοποιείται για την αξιολόγηση μιας δοσμένης ομάδας διπλής κατηγοριοποίησης (Biclustor).[1],[2],[3]



Εικόνα 1.1. – Πρόβλημα NP-hard [3]

Στην παρούσα διπλωματική εργασία στόχος είναι η επεξεργασία δεδομένων με τον αλγόριθμο διπλής κατηγοριοποίησης ο οποίος παρουσιάστηκε το 2000 από τους **Cheng και Church**. [4]

Σκοπός είναι η μελέτη του αλγορίθμου στην αρχική του μορφή και η εκτέλεση πειραμάτων για την εύρεση Biclusters των οποίων τα γονίδια θα παρουσιάζουν λιγότερο θόρυβο (πιο μικρή απόσταση μεταξύ τους), καθώς και ομοιόμορφη δομή. Πιο συγκεκριμένα για τον υπολογισμό της διακύμανσης της ομάδας όλων των στοιχείων του Biclustor, δηλαδή της βαθμολογίας μέσου τετραγωνικού υπολείμματος (mean squared residue score) θα χρησιμοποιηθεί η απόλυτη τιμή και θα αλλάξει ο τύπος υπολογισμού ανά γονίδιο και ανά κατάσταση. Η πειραματική διαδικασία εφαρμόστηκε σε πίνακες γονιδίων-κυτταρικών σειρών προερχόμενα από μικροσυστοιχίες DNA από δείγματα καρκίνου του μαστού, του τραχήλου της μήτρας, των ωοθηκών και του ενδομητρίου.

1.2. Υφιστάμενη Γνώση

Οι πίνακες γονιδιακής έκφρασης, όπως αναφέρθηκε, αναλύονται σε δύο διαστάσεις: τη διάσταση των γονιδίων και τη διάσταση των καταστάσεων. Η πρώτη διάσταση αναλύεται συγκρίνοντας τις γραμμές και η δεύτερη τις στήλες του πίνακα. Κατά την ανάλυση τέτοιων πινάκων τα πιο σημαντικά σημεία μελέτης περιλαμβάνουν: [5]

- Ομαδοποίηση γονιδίων με βάση την έκφρασή τους κάτω από διαφορετικές καταστάσεις.
- Πρόβλεψη νέων γονιδίων, με βάση την έκφραση άλλων γονιδίων με ήδη γνωστή πρόβλεψη.

- Ομαδοποίηση καταστάσεων με βάση την έκφραση ενός αριθμού γονιδίων.
- Πρόβλεψη ενός νέου συνολικού δείγματος, έχοντας ως γνωστή πληροφορία την έκφραση γονιδίων κάτω από συγκεκριμένες πειραματικές καταστάσεις.

Οι τεχνικές απλής ομαδοποίησης χρησιμοποιούνται με στόχο είτε να ομαδοποιήσουν γονίδια είτε καταστάσεις, δηλαδή είναι πιο άμεσες για τα σημεία 1 και 3. Σε αυτές όμως τις τεχνικές συναντάται η εξής δυσκολία. Ομαδοποιούν με ένα τρόπο καθολικό με αποτέλεσμα κάποιες ομάδες που παρουσιάζουν βιολογικό ενδιαφέρον σε μία μικρή ομάδα καταστάσεων να μην μπορούν να αξιοποιηθούν. Τέτοια τοπικά μοντέλα ομαδοποίησης μπορεί να αποτελούν κλειδιά για γενετικά μονοπάτια και αυτά στοχεύουν να ανακαλύψουν οι διάφορες τεχνικές διπλής κατηγοριοποίησης που πραγματοποιούν ταυτόχρονη ομαδοποίηση γραμμών και στηλών σε ένα πίνακα γονιδιακής έκφρασης.

- Η μέθοδος διπλής κατηγοριοποίησης παρουσιάστηκε αρχικά από τον J.A. Hartigan το 1972 [6] - αναδρομικός διαχωρισμός πίνακα μέσα σε blocks (Block clustering) - , ενώ ο όρος καθιερώθηκε το 1996 από τον Mirkin.[7]
- Όμως ο αλγόριθμος που περιγράφει τη γενική μέθοδο δόθηκε το 2000 από τους Y. Cheng και G. M. Church,[4] οι οποίοι και πρότειναν έναν αλγόριθμο διπλής κατηγοριοποίησης που βασίζεται στον υπολογισμό της διακύμανσης (variance) και τον εφάρμοσαν σε δεδομένα γονιδιακής έκφρασης και ακόμα αποτελεί τη βάση της διπλής κατηγοριοποίησης της γονιδιακής έκφρασης.

Από το 2000, ένας μεγάλος αριθμός από επιστημονικά άρθρα έχουν δημοσιευθεί σχετικά με το biclustering. Οι ενδιαφερόμενοι αναγνώστες μπορούν να ανατρέξουν στη βιβλιογραφία, σε άρθρα, όπως των [8],[9],[10],[11].

Ο πίνακας που ακολουθεί παραθέτει κάποιες ενδεικτικές εφαρμογές διπλής κατηγοριοποίησης σε δεδομένα μικροσυστοιχιών τα τελευταία 14 χρόνια:

Πίνακας 1.1. – Εφαρμογές διπλής κατηγοριοποίησης τα τελευταία 14 χρόνια [12]

Σύνολο δεδομένων	Αναφορές
Πολλαπλή σκλήρυνση	C. Tang et al 2001 [13]
Καρκίνος του μαστού	Ben-Dor et al 2002 [14]
Ζυμομύκτης-στρες	E. Segal et al 2003 [15]
Καρκίνος του παχέος εντέρου	T.M. Murali et al 2003 [16]
Λευχαιμία	Q. Sheng et al 2003 [17]
Ζυμομύκτης/Λέμφωμα	C. Cano et al 2007 [18]
Οξεία λεμφοβλαστική λευχαιμία/λέμφωμα/ ζυμομύκτης	U. Maulik & S. Bandyopadhyay 2009 [19]
5 τύποι καρκίνου του πνεύμονα	C.-P. Chen et al 2014 [20]

Η πλειοψηφία των βιολογικών εφαρμογών που εφαρμόζει την μέθοδο της διπλής κατηγοριοποίησης διεξάγεται με τη χρήση της τεχνολογίας των μικροσυστοιχιών η οποία επιτρέπει την μέτρηση της γονιδιακής έκφρασης από χιλιάδες γονίδια υπό σαφείς πειραματικές συνθήκες. Το μεγαλύτερο μέρος των μελετητών που έχουν εφαρμόσει μεθόδους διπλής κατηγοριοποίησης χρησιμοποιούν πίνακες οι οποίοι προκύπτουν από καρκινικά κύτταρα μέσω διαφόρων σταδίων της ασθένειας. Εκμεταλλεύονται επίσης δείγματα από διαφορετικούς ασθενείς και υγιή άτομα όπως και από πρότυπους οργανισμούς (π.χ. ζυμομύκητες) (Πίνακας 1.1.).

1.3. Δομή της Εργασίας

Η οργάνωση των κεφαλαίων που ακολουθούν, βασισμένη στον τρόπο ανάπτυξης της παρούσας εργασίας, έχει ως εξής:

Το πρώτο κεφάλαιο περιλαμβάνει το κίνητρο και τον στόχο της παρούσας διπλωματικής εργασίας καθώς και το σύνολο του αλγοριθμικού υπόβαθρου που μελετήθηκε.

Στο δεύτερο κεφάλαιο αναλύονται βασικές έννοιες όπως το κύτταρο, το γονίδιο και η γονιδιακή έκφραση καθώς και τα διάφορα επίπεδά της, και δίνονται συνοπτικά οι έννοιες των μεταγραφικών προτύπων και των κυτταρικών σειρών. Επιπλέον, περιγράφεται το RNA, το DNA, το μοντέλο διπλής έλικας του DNA καθώς και η επιστήμη της βιοπληροφορικής.

Στο τρίτο κεφάλαιο, αρχικά διαχωρίζονται οι αλγόριθμοι κατηγοριοποίησης βάσει γονιδίων, δειγμάτων και υποδιαστημάτων. Στη συνέχεια περιγράφεται η τεχνική διπλής κατηγοριοποίησης, η δομή και τα είδη των αλγορίθμων της.

Στο τέταρτο κεφάλαιο περιγράφεται ο αλγόριθμος Cheng και Church που εφαρμόστηκε αναπτύσσοντας τα επιμέρους βήματά του και το θεωρητικό υπόβαθρο στο οποίο στηρίζεται. Περιλαμβάνεται ακόμα, η προτεινόμενη μεθοδολογία μας με την εφαρμογή του αλγορίθμου Cheng και Church στον οποίο έχουμε επέμβει προτείνοντας βελτιώσεις που στοχεύουν σε πιο λεπτομερή αποτελέσματα.

Τα αποτελέσματα παρατίθενται στο πέμπτο κεφάλαιο για τις τρεις πρώτες ομάδες γονιδίων διπλής κατηγοριοποίησης: α) για τα 1000 πρώτα από τα 33096 γονίδια στον καρκίνο του μαστού, β) για τα 1000 πρώτα από τα 33096 γονίδια στον καρκίνο του τραχήλου της μήτρας, γ) για τα 1000 πρώτα από τα 33096 γονίδια στον καρκίνο των ωοθηκών, δ) για τα 1000 πρώτα από τα 33096 γονίδια στον καρκίνο του ενδομητρίου, και ε) για το σύνολο των 33096 γονιδίων στον καρκίνο του μαστού και του τραχήλου της μήτρας.

Στο έκτο κεφάλαιο γίνεται η αξιολόγηση των αποτελεσμάτων και τέλος, στο έβδομο κεφάλαιο περιγράφονται κάποια τελικά συμπεράσματα και μελλοντικές επεκτάσεις της δουλειάς μας.

Καινοτομία της Εργασίας

Στην παρούσα εργασία πραγματοποιείται η εφαρμογή της μεθόδου διπλής κατηγοριοποίησης δεδομένων γονιδιακής έκφρασης που προέρχονται από δείγματα παθολογίας του καρκίνου του μαστού, των

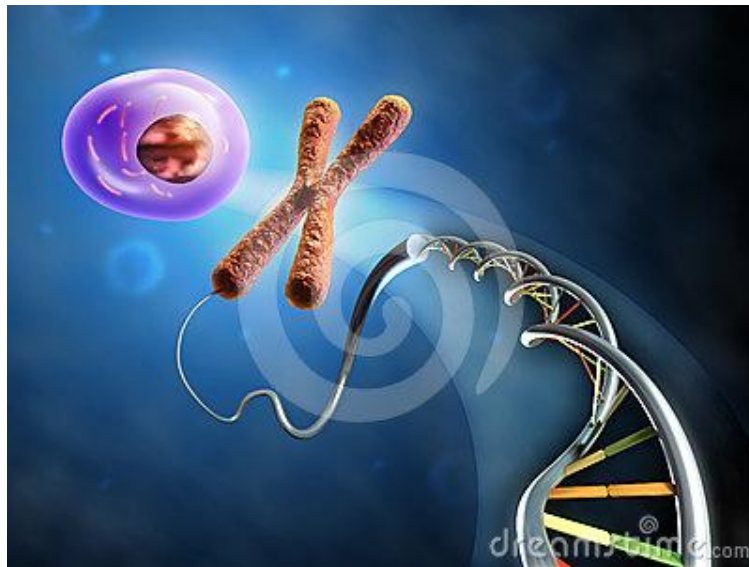
ωοθηκών, του ενδομητρίου και του τραχήλου της μήτρας. Η μεθοδολογία που προτείνεται περιλαμβάνει την εφαρμογή του αλγορίθμου διπλής κατηγοριοποίησης Cheng και Church εφαρμόζοντας σε αυτόν μια σειρά βελτιώσεων ούτως ώστε να εξάγουμε ταυτόχρονα ομάδες γονιδίων με ομοιόμορφη συμπεριφορά αλλά και συγκεκριμένο εύρος τιμών. Για να μελετήσουμε τη συμπεριφορά των γονιδίων, ώστε να διακρίνουμε τις βελτιώσεις που έπρεπε να κάνουμε στον αλγόριθμο, μειώσαμε τον αριθμό των δεδομένων από 33096 στα 1000 γονίδια και εστίασαμε στα τρία πρώτα εξαγόμενα Biclusters. Συγκεκριμένα, οι βελτιώσεις που κάναμε στον αλγόριθμο αφορούν το $d(i)$ και $d(j)$, τα οποία αντιπροσωπεύουν το «mean residue score» για κάθε γραμμή και στήλη, αντίστοιχα. Ύστερα από αρκετές δοκιμές στους τύπους, καταλήξαμε και εφαρμόσαμε τις βελτιωμένες εξισώσεις στον αλγόριθμο διπλής κατηγοριοποίησης Cheng και Church. Εκτελέσαμε τον αλγόριθμο τρεις φορές για κάθε τύπο καρκίνου. Κατά την πρώτη εκτέλεση, εξαγάγαμε τα αποτελέσματα χωρίς να έχουμε κάνει καμία αλλαγή στις εξισώσεις $d(i)$ και $d(j)$. Στην δεύτερη εκτέλεση, αλλάξαμε μόνο τον τύπο της εξίσωσης του $d(i)$, και τέλος, κατά την τρίτη εκτέλεση, έχοντας κάνει τις επιθυμητές αλλαγές και στις δύο εξισώσεις ($d(i)$ και $d(j)$), εξαγάγαμε τα αποτελέσματα. Μελετώντας στη συνέχεια τα αποτελέσματα αυτά, παρατηρήσαμε ότι είχαμε πετύχει τον αρχικό μας στόχο σχετικά με την ομοιόμορφη ομαδοποίηση των γονιδίων και τη μείωση του εύρους τιμών σε κάθε Bicluster.

2

Θεωρητικό Υπόβαθρο

2.1. Βιολογικό Υπόβαθρο

Το κύτταρο αποτελεί τη δομική και λειτουργική μονάδα της ζωής. Στο κύτταρο η λειτουργία και η μορφολογία είναι αλληλένδετες έννοιες, δεδομένου ότι η λειτουργία του κυττάρου καθορίζει τη μορφολογία του κυττάρου και η μορφολογία του κυττάρου εξυπηρετεί τη λειτουργία του. Όλα τα κύτταρα δομούνται από τις ίδιες χημικές ενώσεις και εκδηλώνουν παρόμοιες μεταβολικές διεργασίες. Τα κύτταρα περιέχουν επίσης το γενετικό υλικό του οργανισμού, το DNA (Εικόνα 2.1.).[5]



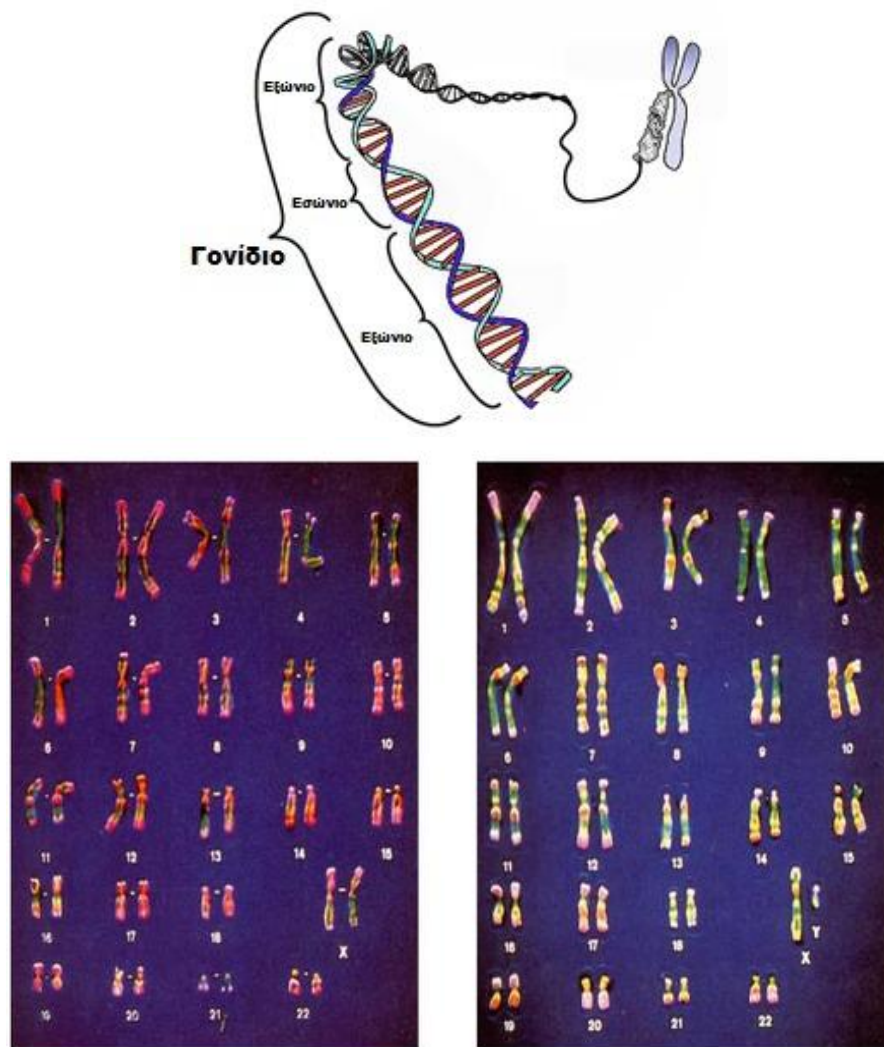
Εικόνα 2.1. – Από το κύτταρο στο DNA [21]

Το DNA ενός οργανισμού είναι το μοριακό πρόγραμμα που περιέχει ακριβείς οδηγίες, οι οποίες καθορίζουν τη δομή και τη λειτουργία του οργανισμού. Επιπλέον, περιέχει τις οδηγίες για τον αυτοδιπλασιασμό του. Με τον αυτοδιπλασιασμό επιτυγχάνεται η μεταβίβαση των γενετικών πληροφοριών από ένα κύτταρο (μητρικό) στα θυγατρικά του κύτταρα και από έναν οργανισμό στους απογόνους του μέσω της διαδικασίας της αντιγραφής (replication). Οι εντολές που δέχεται το κύτταρο για να επιτελέσει ορισμένες διαδικασίες μεταφέρονται μέσω των παραγομένων πρωτεϊνών.[5]

2.1.1. Γενετικό υλικό και Χρωμοσώματα

Η γενετική πληροφορία είναι η καθορισμένη σειρά των αζωτούχων βάσεων των νουκλεοτιδίων. Η πληροφορία υπάρχει σε τμήματα του DNA με συγκεκριμένη ακολουθία, τα γονίδια.[5] Με άλλα λόγια, γονίδια ονομάζονται οι αλληλουχίες νουκλεοτιδίων τμήματος του DNA, οι οποίες ελέγχουν τα κληρονομικά γνωρίσματα ενός οργανισμού. Οι αλληλουχίες που μεταφράζονται σε αμινοξέα αποτελούν τα εξώνια, ενώ οι ενδιάμεσες αλληλουχίες ονομάζονται εσώνια και δε μεταφράζονται σε αμινοξέα (Εικόνα 2.2 - άνω μέρος). Όλα τα κύτταρα ενός πολυκύτταρου οργανισμού έχουν το ίδιο DNA. Σε κάθε ομάδα κυττάρων όμως εκφράζονται διαφορετικά γονίδια.

Το γενετικό υλικό ενός κυττάρου αποτελεί το γονιδίωμά του. Κύτταρα στα οποία το γονιδίωμα υπάρχει σε ένα μόνο αντίγραφο, ονομάζονται απλοειδή (προκαρυωτικά κύτταρα, γαμέτες διπλοειδών οργανισμών). Τα κύτταρα στα οποία το γονιδίωμα υπάρχει σε δυο αντίγραφα, ονομάζονται διπλοειδή (σωματικά κύτταρα ανώτερων ευκαρυωτικών οργανισμών). Το γονιδίωμα καθορίζει τη γενετική κατασκευή ενός οργανισμού ή ενός κυττάρου, ή τον γονότυπό του. Ο φαινότυπος, είναι το σύνολο των χαρακτηριστικών που εμφανίζει ένας οργανισμός υπό την επιρροή ενός συνόλου περιβαλλοντικών παραγόντων.[22]



Εικόνα 2.2. – Γονίδιο (επάνω) και Ανθρώπινος καρυότυπος γυναίκας (κάτω-αριστερά) και άνδρα (κάτω-δεξιά) [23]

Στα ευκαρυωτικά κύτταρα το DNA κατανέμεται στον πυρήνα, στα μιτοχόνδρια και στους χλωροπλάστες (φυτικά κύτταρα). Το συνολικό μήκος του DNA στα διπλοειδή κύτταρα του ανθρώπου είναι περίπου 2 m (6×10^9 ζεύγη βάσεων) και συσπειρώνεται σε τέτοιο βαθμό ώστε να χωράει στον πυρήνα του κυττάρου που έχει διάμετρο 10 μm , δηλαδή το DNA είναι ισχυρά «πακεταρισμένο» στο χρωμόσωμα. Από τα παραπάνω προκύπτει ότι η σημαντικότερη λειτουργία των χρωμοσωμάτων είναι να μεταφέρουν τα γονίδια-τις λειτουργικές μονάδες της κληρονομικότητας. Το ένα χρωμόσωμα κάθε ζεύγους είναι πατρικής και το άλλο μητρικής προέλευσης και ελέγχουν τις ίδιες ιδιότητες. Από τα 23 ζεύγη τα 22 είναι μορφολογικά ίδια στα αρσενικά και στα θηλυκά άτομα και ονομάζονται αυτοσωμικά χρωμοσώματα. Το 23ο ζεύγος στα θηλυκά άτομα αποτελείται από δύο X χρωμοσώματα, ενώ στα αρσενικά από ένα X και ένα Y χρωμόσωμα. Τα χρωμοσώματα αυτά ονομάζονται φυλετικά και σε πολλούς οργανισμούς, συμπεριλαμβανομένου και του ανθρώπου, καθορίζουν το φύλο. Στον άνθρωπο η παρουσία του Y χρωμοσώματος καθορίζει το αρσενικό άτομο, ενώ η απουσία του το θηλυκό άτομο.[23] Η παρουσίαση του συνόλου των 46 χρωμοσωμάτων του ανθρώπου αποκαλείται καρυότυπος του ανθρώπου (Εικόνα 2.2 - κάτω μέρος).

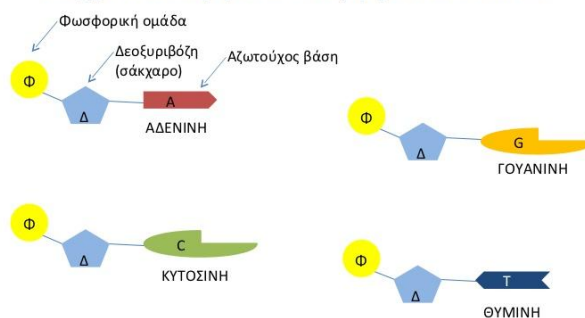
2.1.2. DNA

Το DNA, όπως και το RNA, είναι ένα μακρομόριο, που αποτελείται από νουκλεοτίδια. Κάθε νουκλεοτίδιο του DNA αποτελείται από μία πεντόζη, τη δεοξυριβόζη, ενωμένη με μία φωσφορική ομάδα και μία αζωτούχος βάση. Στα νουκλεοτίδια του DNA η αζωτούχος βάση μπορεί να είναι μια από τις: αδενίνη (A), γουανίνη (G), κυτοσίνη (C) και θυμίνη (T) (Εικόνα 2.3).

Κάθε νουκλεϊκό οξύ σχηματίζεται από πολλά νουκλεοτίδια.

Το DNA σχηματίζεται από πολλά δεοξυριβονουκλεοτίδια.

Υπάρχουν 4 διαφορετικά δεοξυριβονουκλεοτίδια:



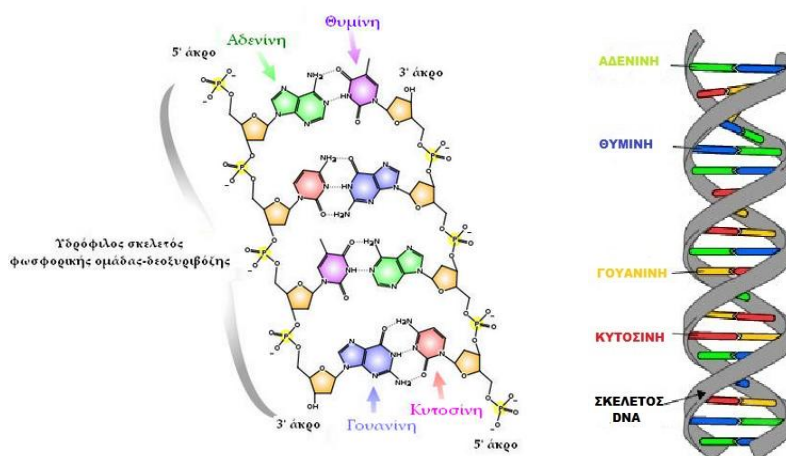
Εικόνα 2.3. – Τα τέσσερα διαφορετικά δεοξυριβονουκλεοτίδια του DNA [24]

Σε κάθε νουκλεοτίδιο η αζωτούχος βάση συνδέεται με τον 1' άνθρακα της δεοξυριβόζης και η φωσφορική ομάδα με τον 5' άνθρακα. Μια πολυνουκλεοτιδική αλυσίδα σχηματίζεται από την ένωση πολλών νουκλεοτιδίων με ομοιοπολικό δεσμό. Ο δεσμός αυτός δημιουργείται μεταξύ του υδροξυλίου του 3' άνθρακα της πεντόζης του πρώτου νουκλεοτιδίου και της φωσφορικής ομάδας που είναι συνδεδεμένη στον 5' άνθρακα

της πεντόζης του επόμενου νουκλεοτιδίου. Ο δεσμός αυτός ονομάζεται 3'-5' φωσφοδιεστερικός δεσμός. Με τον τρόπο αυτό η πολυνουκλεοτιδική αλυσίδα που δημιουργείται έχει ένα σκελετό, που αποτελείται από επανάληψη των μορίων φωσφορική ομάδα- πεντόζη-φωσφορική ομάδα-πεντόζη. Ανεξάρτητα από τον αριθμό των νουκλεοτιδίων από τα οποία αποτελείται η πολυνουκλεοτιδική αλυσίδα, το πρώτο της νουκλεοτίδιο έχει πάντα μία ελεύθερη φωσφορική ομάδα συνδεδεμένη στον 5' άνθρακα της πεντόζης του και το τελευταίο νουκλεοτίδιο της έχει ελεύθερο το υδροξύλιο του 3' άνθρακα της πεντόζης του. Για το λόγο αυτό αναφέρεται ότι ο προσανατολισμός της πολυνουκλεοτιδικής αλυσίδας είναι 5'→ 3' (Εικόνα 2.4).[22]

2.1.3. Το μοντέλο της διπλής έλικας

Το 1953 οι Watson και Crick διατύπωσαν το μοντέλο της διπλής έλικας του DNA, που αναφέρεται στη δομή του DNA στο χώρο. Σύμφωνα με το μοντέλο αυτό, το DNA αποτελείται από δύο πολυνουκλεοτιδικές αλυσίδες που σχηματίζουν στο χώρο μία δεξιόστροφη διπλή έλικα. Η διπλή έλικα έχει ένα σταθερό σκελετό, που αποτελείται από επαναλαμβανόμενα μόρια φωσφορικής ομάδας-δεοξυριβόζης ενωμένων με φωσφοδιεστερικό δεσμό. Ο σκελετός αυτός είναι υδρόφιλος και βρίσκεται στο εξωτερικό του μορίου. Προς το εσωτερικό του σταθερού αυτού σκελετού βρίσκονται οι αζωτούχες βάσεις που είναι υδρόφοβες. Οι αζωτούχες βάσεις της μιας αλυσίδας συνδέονται με δεσμούς υδρογόνου με τις αζωτούχες βάσεις της απέναντι αλυσίδας με βάση τον κανόνα της συμπληρωματικότητας. Η αδενίνη συνδέεται μόνο με θυμίνη και αντίστροφα, ενώ η κυτοσίνη μόνο με γουανίνη και αντίστροφα. Ανάμεσα στη θυμίνη και την αδενίνη αναπτύσσονται δύο δεσμοί υδρογόνου, ενώ ανάμεσα στη κυτοσίνη και τη γουανίνη αναπτύσσονται τρεις δεσμοί υδρογόνου (Εικόνα 2.4).



Εικόνα 2.4. – Το μοντέλο της διπλής έλικας του DNA [25]

Οι δεσμοί υδρογόνου που σχηματίζονται μεταξύ των βάσεων σταθεροποιούν τη δευτεροταγή δομή του μορίου. Οι δύο αλυσίδες ενός μορίου DNA είναι συμπληρωματικές, και αυτό υποδηλώνει ότι η αλληλουχία της μιας καθορίζει την αλληλουχία της άλλης. Η συμπληρωματικότητα έχει τεράστια σημασία

για τον αυτοδιπλασιασμό του DNA, μια ιδιότητα που το καθιστά το καταλληλότερο μόριο για τη διατήρηση και τη μεταβίβαση της γενετικής πληροφορίας.[22],[25]

2.1.4. RNA

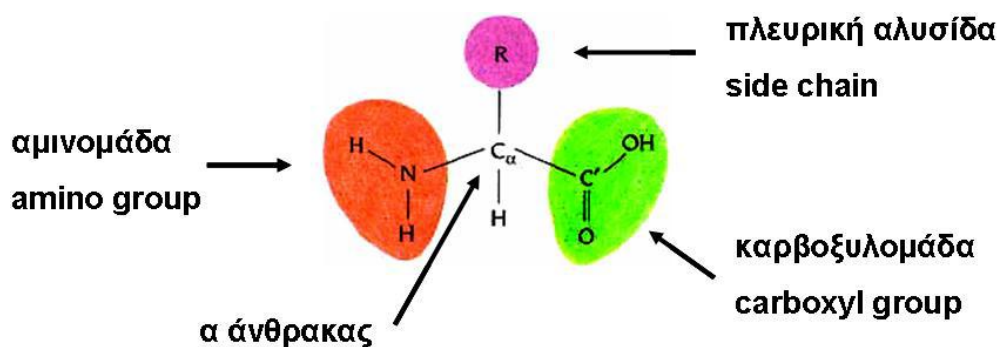
Το RNA ή διαφορετικά ριβονουκλεϊκό οξύ είναι ένα πολυμερές νουκλεϊκό οξύ το οποίο απαρτίζεται από μονομερή νουκλεοτίδια που επιτελούν βασικό ρόλο κατά τη μετάφραση του δεοξυριβονουκλεϊκού οξέος (DNA) σε πρωτεϊνικά προϊόντα. Το RNA λειτουργεί ως ο αγγελιοφόρος του DNA (αγγελιοφόρο RNA ή mRNA). Το RNA εκτός από βασικό στοιχείο των ριβοσωμάτων (rRNA) είναι μεταφορέας αμινοξέων (tRNA - διαθέτει ειδικούς υποδοχείς για την πρόσδεση αμινοξέων και μια περιοχή αντικωδικονίου για την αναγνώριση του κωδικονίου που δεσμεύεται με την ειδική ακολουθία του αγγελιοφόρου RNA με ισχυρούς δεσμούς υδρογόνου) που χρησιμοποιούνται στη διαδικασία της πρωτεϊνοσύνθεσης. Το RNA είναι σχεδόν πανομοιότυπο με το DNA. Η διαφορά τους είναι στο ότι τα μόρια RNA αποτελούνται από ριβόζη αντί για δεοξυριβόζη σαν το κύριο σάκχαρο και επίσης περιέχει τη βάση ουρακίλη αντί για τη θυμίνη που απαντάται στο DNA. Το RNA μεταγράφεται από το DNA με τη βοήθεια κυρίως ενός ενζύμου που ονομάζεται RNA πολυμεράση και στη συνέχεια επεξεργάζεται με έναν αριθμό άλλων δευτερευόντων ενζύμων. Στη συνέχεια χρησιμοποιείται σαν βάση για τη μετάφραση των γονιδίων σε πρωτεΐνες, μεταφέροντας αμινοξέα στα ριβοσώματα για να δημιουργηθούν πρωτεΐνες (πάντα βάσει της αρχής της συμπληρωματικότητας).

Στους ευκαρυωτικούς οργανισμούς το mRNA που παράγεται με τη μεταγραφή υφίσταται μια διαδικασία ωρίμανσης προτού να είναι έτοιμο να προχωρήσει στα ριβοσώματα για τη μετάφραση. Αυτό γίνεται γιατί τα περισσότερα γονίδια είναι ασυνεχή και εκτός των αλληλουχιών που μεταφράζονται υπάρχουν αλληλουχίες οι οποίες δε μεταφράζονται σε αμινοξέα. Οι αλληλουχίες που μεταφράζονται είναι τα εξώνια και εκείνες που δε μεταφράζονται είναι τα εσώνια. Έτσι, το mRNA που μόλις έχει σχηματιστεί από τη μεταγραφή ενός γονιδίου ονομάζεται πρόδρομο mRNA και περιέχει εσώνια και εξώνια. Το πρόδρομο mRNA μετατρέπεται σε mRNA με τη διαδικασία της ωρίμανσης κατά την οποία τα εσώνια κόβονται από μικρά ριβοζονουκλεοπρωτεϊνικά σωματίδια (snRNPs), που λειτουργούν ως ένζυμα. Τα εξώνια που απομένουν συρράπτονται μεταξύ τους και με αυτό τον τρόπο σχηματίζεται το ώριμο mRNA που μεταφέρεται στα ριβοσώματα για την πρωτεϊνοσύνθεση. Δύο περιοχές του ώριμου mRNA δε μεταφράζονται σε αμινοξέα, η μια βρίσκεται στο άκρο 5' και η άλλη στο άκρο 3'. [22]

2.1.5. Πρωτεΐνες

Οι πρωτεΐνες είναι μακρομόρια, πολυμερή, με σαφώς καθορισμένη δομή και παρουσιάζουν μεγάλη ποικιλία ιδιοτήτων και σχήματος. Οι πρωτεΐνες είναι απαραίτητες σε όλες σχεδόν τις βιολογικές λειτουργίες και η σημασία τους συνοψίζεται στα εξής: ενζυμική κατάλυση, μεταφορά, αποθήκευση, κίνηση, μηχανική στήριξη, ανοσοπροστασία, δημιουργία και μετάδοση νευρικών παλμών και έλεγχος της ανάπτυξης, της διαφοροποίησης.

Οι δομικές μονάδες των πρωτεϊνών είναι 20 διαφορετικά αμινοξέα. Τα αμινοξέα αποτελούνται από μία αμινομάδα ($-NH_2$) και μία καρβοξυλομάδα ($-COOH$) συνδεδεμένες σε ένα άτομο άνθρακα C_α . Στο ίδιο άτομο ενώνεται και μία πλευρική αλυσίδα, η οποία διαφέρει μεταξύ των αμινοξέων και καθορίζει τις ιδιότητές τους. Τα αμινοξέα ενώνονται μεταξύ τους με πεπτιδικούς δεσμούς για τη δημιουργία πολυπεπτιδικών αλυσίδων - πρωτεϊνών. Τα πρωτεϊνικά μόρια αποτελούνται από μια ή περισσότερες πολυπεπτιδικές αλυσίδες, που η κάθε μία τους μπορεί να αποτελείται από μερικές εκατοντάδες αμινοξέα με μια συγκεκριμένη σειρά γνωστή σαν αμινοξική ακολουθία.



Εικόνα 2.5. – Αμινοξύ – δομική μονάδα πρωτεϊνών

Στις πρωτεΐνες μπορούν να αναγνωριστούν 4 επίπεδα οργάνωσης:

α. Πρωτοταγής δομή: αντιστοιχεί στην αμινοξική ακολουθία, δηλαδή στην διάταξη των αμινοξέων σε μια πολυπεπτιδική αλυσίδα.

β. Δευτεροταγής δομή: αναφέρεται στην κανονική στερεοδιάταξη τμημάτων της πολυπεπτιδικής αλυσίδας.

γ. Τριτοταγής δομή: αναφέρεται στην τρισδιάστατη δομή.

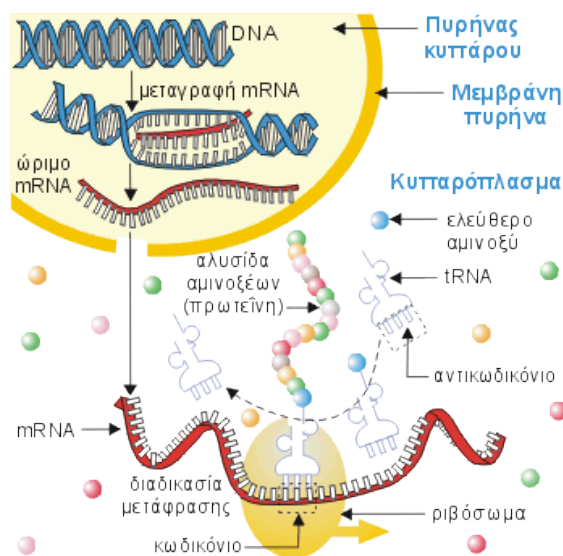
δ. Τεταρτογενής Δομή: η διάταξη στο χώρο των πολυπεπτιδικών αλυσίδων από τις οποίες αποτελείται.

Η ακολουθία μιας πρωτεΐνης καθορίζει την τρισδιάστατη δομή της, η οποία με τη σειρά της καθορίζει τις αλληλεπιδράσεις της με άλλες πρωτεΐνες, με νουκλεϊνικά οξέα και μικρά μόρια. Επομένως, η δομή μιας πρωτεΐνης καθορίζει τη λειτουργία της και οι ερευνητικές προσπάθειες κατευθύνονται τόσο στον καθορισμό των αμινοξικών ακολουθιών της κάθε πρωτεΐνης όσο και την εύρεση της τρισδιάστατης της δομής. [22], [26]

2.1.6. Γονιδιακή Έκφραση

Γονιδιακή έκφραση καλείται η διαδικασία κατά την οποία η πληροφορία από ένα γονίδιο χρησιμοποιείται για την σύνθεση ενός λειτουργικού γονιδιακού προϊόντος. Τα προϊόντα αυτά είναι συχνά πρωτεΐνες. Εκτός από τις πρωτεΐνες, ως λειτουργικό γονιδιακό προϊόν θεωρείται επίσης το RNA (π.χ. rRNA και tRNA). Ένα γονίδιο εκφράζεται μέσω των μηχανισμών της μεταγραφής και της μετάφρασης. Το πρώτο βήμα για την έκφραση της πληροφορίας που υπάρχει στο DNA είναι η μεταβίβασή της στο RNA μέσω της διαδικασίας της μεταγραφής (transcription). Στη συνέχεια, το RNA, μέσω της διαδικασίας της μετάφρασης (translation), μεταφέρει την πληροφορία στις πρωτεΐνες που είναι υπεύθυνες για τη δομή και λειτουργία των κυττάρων και κατ' επέκταση και των οργανισμών. Έτσι, οι πορείες της μεταγραφής και της μετάφρασης των γονιδίων αποτελούν τη γονιδιακή έκφραση (Εικόνα 2.6).[27] Η διαδικασία της γονιδιακής έκφρασης πραγματοποιείται σε όλους τους έμβιους οργανισμούς, ευκαρυωτικούς, προκαρυωτικούς και στους ιούς, προκειμένου να παραχθεί ο μακρομοριακός μηχανισμός της ζωής. Στην γενετική η γονιδιακή έκφραση αποτελεί το θεμελιώδες επίπεδο μέσω του οποίου ο γονότυπος δίνει γέννηση στο φαινότυπο. Ο γενετικός κώδικας ερμηνεύεται μέσω της γονιδιακής έκφρασης και οι ιδιότητες των εκφραζόμενων προϊόντων μας δίνουν τον φαινότυπο του οργανισμού.

Η μέτρηση της γονιδιακής έκφρασης, δηλαδή η ποσοτικοποίηση του επιπέδου έκφρασης ενός γονιδίου μέσα σε ένα κύτταρο, αποτελεί μία τεράστια ποσότητα πληροφορίας.



Εικόνα 2.6. - Μηχανισμός σύνθεσης πρωτεϊνών στα ευκαρυωτικά κύτταρα [27]

Η πληροφορία της γονιδιακής έκφρασης εστιάζει είτε:

(α) στον καθορισμό των ιστών οι οποίοι συνδέονται με τη φυσιολογική λειτουργία ενός συγκεκριμένου γονιδίου, είτε

(β) στον καθορισμό των παραγόντων εκείνων που ρυθμίζουν την έκφραση ενός συγκεκριμένου γονιδίου, οι οποίοι συνήθως χωρίζονται σε διατροφικούς, ορμονικούς, ή περιβαλλοντικούς.[22]

Οπότε έχοντας την παραπάνω πληροφορία δύνανται να πραγματοποιηθεί πρόβλεψη για ένα άτομο εάν έχει προδιάθεση για κάποια μορφή καρκίνου ή να βρεθεί εάν κάποια κύτταρα παρουσιάζουν ανθεκτικότητα σε κάποιους φαρμακευτικούς παράγοντες κ.α..

Όπως αναφέρθηκε, ο όρος γονιδιακή έκφραση αναφέρεται συνήθως σε όλη τη διαδικασία με την οποία ένα γονίδιο ενεργοποιείται για να παραγάγει μια πρωτεΐνη. Όμως σε κάθε κύτταρο δεν παράγονται όλες οι πρωτεΐνες σε κάθε χρονική στιγμή. Επιπλέον, επειδή το κύτταρο χρειάζεται κάθε πρωτεΐνη σε συγκεκριμένη ποσότητα, οι πρωτεΐνες ενός κυττάρου δεν παράγονται σε ίσες ποσότητες. Έτσι, είναι απαραίτητη η ύπαρξη και η λειτουργία ενός προγράμματος ρύθμισης της γονιδιακής έκφρασης, που παρέχει τις οδηγίες για το είδος και την ποσότητα των πρωτεϊνών οι οποίες πρέπει να παραχθούν σε κάθε συγκεκριμένη χρονική στιγμή. Στα ευκαρυωτικά κύτταρα η γονιδιακή έκφραση ρυθμίζεται σε τέσσερα επίπεδα:

- **Στο επίπεδο της μεταγραφής:**

Ένας αριθμός μηχανισμών ελέγχουν ποια γονίδια θα μεταγραφούν ή και με ποια ταχύτητα θα γίνει η μεταγραφή. Το DNA των ευκαρυωτικών κυττάρων δεν οργανώνεται σε οπερόνια αλλά κάθε γονίδιο έχει τον δικό του υποκινητή και μεταγράφεται αυτόνομα. Η μεταγραφή περιλαμβάνει τη μεταφορά της γενετικής πληροφορίας από το DNA μέσω ενζυματικής σύνθεσης μιας συμπληρωματικής αλυσίδας RNA που καταλύεται από το ένζυμο της RNA πολυμεράσης. Η RNA πολυμεράση λειτουργεί (όπως και στους προκαρυωτικούς οργανισμούς) με τη βοήθεια πρωτεϊνών, που ονομάζονται μεταγραφικοί παράγοντες. Μόνο που στους ευκαρυωτικούς οργανισμούς οι μεταγραφικοί παράγοντες παρουσιάζουν τεράστια ποικιλία. Κάθε κυτταρικός τύπος περιέχει διαφορετικά είδη μεταγραφικών παραγόντων. Διαφορετικός συνδυασμός μεταγραφικών παραγόντων ρυθμίζει τη μεταγραφή κάθε γονιδίου. Μόνο όταν ο σωστός συνδυασμός των μεταγραφικών παραγόντων προσδεθεί στον υποκινητή ενός γονιδίου, αρχίζει η RNA πολυμεράση τη μεταγραφή ενός γονιδίου.

- **Στο επίπεδο μετά τη μεταγραφή:**

Περιλαμβάνονται οι μηχανισμοί με τους οποίους γίνεται η ωρίμανση του πρόδρομου mRNA και καθορίζεται η ταχύτητα με την οποία το ώριμο mRNA αφήνει τον πυρήνα και εισέρχεται στο κυτταρόπλασμα.

- **Στο επίπεδο της μετάφρασης:**

Ο χρόνος που "ζουν" τα μόρια mRNA στο κυτταρόπλασμα δεν είναι ο ίδιος για όλα τα είδη RNA, επειδή μετά από κάποιο χρονικό διάστημα αποικοδομούνται. Επίσης, ποικίλλει και η ικανότητα πρόσδεσης του mRNA στα ριβοσώματα.

- **Στο επίπεδο μετά τη μετάφραση:**

Ακόμη και όταν πραγματοποιηθεί η πρωτεϊνόςύνθεση και παραχθεί η κατάλληλη πρωτεΐνη, μπορεί να χρειαστεί να υποστεί τροποποιήσεις, για να γίνει βιολογικά λειτουργική.

Μερικά γονίδια που κωδικοποιούν πρωτεΐνες μεταγράφονται λιγότερο ή περισσότερο συχνά, και ονομάζονται «housekeeping» γονίδια και απαιτούνται πάντα για τις βασικές αντιδράσεις. Άλλα γονίδια δεν μεταγράφονται ή, μεταγράφονται για συγκεκριμένες λειτουργίες του οργανισμού, μόνο σε ιδιαίτερες στιγμές και κάτω από ιδιαίτερες εξωτερικές συνθήκες. Το σήμα που «καλύπτει» ή «αποκαλύπτει» ένα γονίδιο μπορεί

να προέλθει από το εξωτερικού του κυττάρου, όπως μια θρεπτική ουσία ή μια ορμόνη. Πρόσθετες ρυθμιστικές ακολουθίες στο DNA υπαγορεύουν εάν ένα γονίδιο θα ανταποκριθεί στα σήματα και στη συνέχεια επηρεάζουν την μεταγραφή του γονιδίου που κωδικοποιεί την πρωτεΐνη. Η αποκωδικοποίηση του ανθρώπινου γονιδιώματος (Human Genome Project) αλλά και άλλων σημαντικών οργανισμών αποτέλεσε σταθμό στη Βιολογία, μεταβάλλοντας ριζικά τις πρακτικές της και την οπτική της έρευνας.[22]

Προφίλ γονιδιακής έκφρασης και ανάλυση μεταγραφώματος

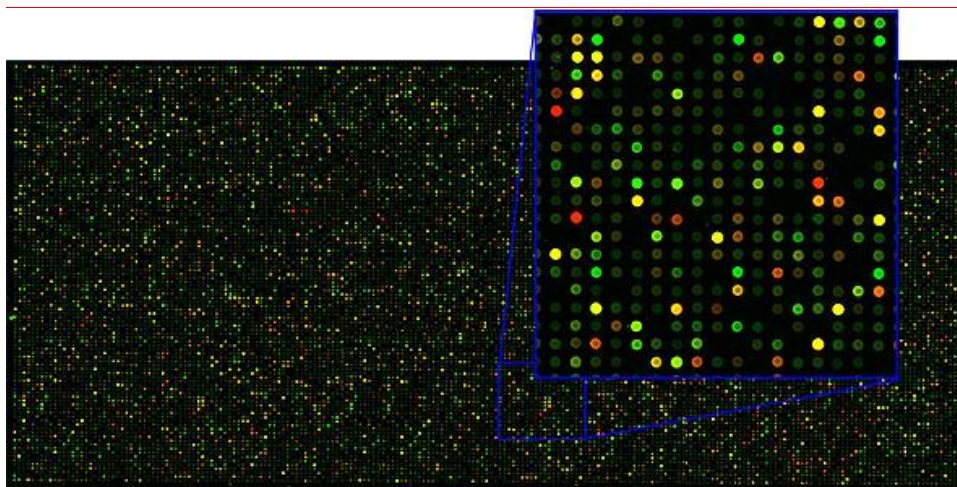
Το προφίλ γονιδιακής έκφρασης είναι ο προσδιορισμός του προτύπου γονιδίων, που εκφράζονται στο επίπεδο της μεταγραφής, κάτω από ειδικές περιστάσεις ή σε ένα ειδικό κύτταρο, το οποίο αποτυπώνει μια συνολική εικόνα της κυτταρικής λειτουργίας.

Το μεταγράφομα μετράει το αγγελιαφόρο RNA (mRNA), δηλαδή το μήνυμα που μεταφράζεται σε πρωτεΐνη. Σε αντίθεση με το γονιδίωμα, το μεταγράφομα δεν λέει ποιά γονίδια είναι παρόντα αλλά ποιά γονίδια μεταγράφονται και μεταφράζονται σε πρωτεΐνες. Δεν υπάρχει ακριβής συσχέτιση, επειδή υψηλά επίπεδα mRNA δεν μεταφράζονται αναγκαστικά σε αντίστοιχα υψηλά επίπεδα παραγωγής πρωτεϊνών και τα υψηλά επίπεδα παραγωγής πρωτεϊνών μπορούν να ακολουθήσουν χαμηλά επίπεδα mRNA. Επιπλέον, κάποιο mRNA δεν μεταφράζεται σε πρωτεΐνη και διαφορετικές μορφές της ίδιας πρωτεΐνης μπορεί να προκύψουν από ένα δεδομένο μετάγραφο.[28] Σήμερα, η μελέτη του μεταγραφώματος υπερτερεί της μελέτης του πρωτεώματος, και ενδείκνυται για την ποσοτικοποίηση της γονιδιακής έκφρασης σε πολλές περιπτώσεις. Η μεταγραφή των mRNA, όπως αναφέρεται παραπάνω, αποτελεί το πρώτο χρονικά βήμα, και είναι αυτό που συγκεντρώνει και τα περισσότερα πλεονεκτήματα για αξιόπιστες και αποτελεσματικές μετρήσεις γονιδιακής έκφρασης. Το mRNA απομονώνεται εύκολα και σε ποσότητες ικανές για αποτελεσματική ποσοτικοποίηση, οι μεθοδολογίες που υπάρχουν μπορούν να ενισχύσουν το δείγμα και να ταυτοποιήσουν έτσι ακόμα και μόρια mRNA που βρίσκονται σε πολύ μικρές συγκεντρώσεις, ενώ, επιπλέον, απ' όσο γνωρίζουμε η μεταγραφή είναι το στάδιο στη διαδικασία της έκφρασης γονιδίων, που υπόκειται στην πιο αυστηρή ρύθμιση. Με βάση τα παραπάνω, οι μετρήσεις σε επίπεδο mRNA, εξασφαλίζουν ομοιογένεια του δείγματος, είναι πιο εύκολες πειραματικά και μπορούν να είναι ποσοτικά ακριβείς. Επιπλέον, τα επίπεδα του mRNA αντανακλούν πιο άμεσα τις αλλαγές που επισυμβαίνουν στην έκφραση των γονιδίων, καθώς ο χρόνος που μεσολαβεί από τη στιγμή που ένα γονίδιο ενεργοποιείται ως την παραγωγή της πρωτεΐνης που κωδικοποιεί μπορεί συχνά να είναι απαγορευτικός για την μελέτη ταχείας απόκρισης σε ερεθίσματα. Τέλος, μια σειρά από μελέτες έχουν δείξει ότι σε σταθερές συνθήκες (steady state) τα επίπεδα των mRNA συσχετίζονται σε μεγάλο βαθμό με τα αντίστοιχα των πρωτεϊνών και παρά το γεγονός ότι ο βαθμός συσχέτισης είναι μειωμένος για τα mRNA που βρίσκονται σε χαμηλή συγκέντρωση, η ποσοτικοποίηση των mRNA θεωρείται ο πιο διαδεδομένος, άμεσος και συστηματικός τρόπος για τη μελέτη της γονιδιακής έκφρασης.[29] Από τις βασικές μεθοδολογίες για την πειραματική μελέτη των επιπέδων mRNA αποτελούν οι μικροσυστοιχίες DNA (DNA microarrays).

2.2. Εισαγωγή στην τεχνολογία των Μικροσυστοιχιών (microarrays)

Οι **μικροσυστοιχίες** (microarrays) είναι μία διάταξη μικροσκοπικών σημείων που αντιπροσωπεύουν μοναδικά γονίδια και ακινητοποιούνται με ομοιοπολικούς δεσμούς σε μία στερεή επιφάνεια (συνήθως

γυάλινη). Χρησιμοποιούνται για τη μέτρηση DNA ή χρησιμοποιούν DNA για το σύστημα ανίχνευσής τους. Ποσοτικές ή ποιοτικές μετρήσεις με μικροσυστοιχίες γονιδίων εκμεταλλεύονται την εκλεκτική φύση της αρχής της συμπληρωματικότητας μεταξύ νουκλεϊκών οξέων DNA-DNA ή DNA-RNA ή πρόσφατα και μεταξύ των αμινοξέων των πρωτεϊνών, υπό αυστηρά ελεγχόμενες συνθήκες θερμοκρασίας και με τη χρήση φθορίζουσων ουσιών. Οι μικροσυστοιχίες γονιδίων χρησιμοποιούνται σήμερα κατά κόρον για την εξέταση της γονιδιακής έκφρασης υπό ειδικές συνθήκες και για την ανίχνευση νουκλεϊκών οξέων παθογόνων οργανισμών π.χ. επιβλαβών ιών σε δείγματα ελέγχου. Οι μικροσυστοιχίες αποτελούν πολύτιμη μέθοδο της μοριακής βιολογίας, δεδομένου ότι δύνανται να εξετάζουν ταυτόχρονα την έκφραση χιλιάδων γονιδίων, και ενδείκνυνται για συγκριτικές μελέτες γονιδιωμάτων. Τα μειονεκτήματα της μεθόδου έγκεινται στο υψηλό κόστος και στη συχνή ανακρίβεια των αποτελεσμάτων λόγω τεχνικών προβλημάτων όπως η μη ειδική υβριδοποίηση φθορίζουσων χρωστικών σε λάθος γονίδια κτλ. Σήμερα, η χρήση μικροσυστοιχιών γονιδίων από την επιστημονική κοινότητα απαιτεί την παράλληλη χρήση και άλλων συμπληρωματικών μεθόδων για την επαλήθευση των αποτελεσμάτων του τσιπ όπως είναι το στύπωμα Northern (Northern Blot) και η ποσοτική αλυσιδωτή αντίδραση πολυμεράσης με τη χρήση αντίστροφης μεταγραφάσης σε πραγματικό χρόνο (Quantitative Real time (reverse transcriptase) PCR)).[30]



Εικόνα 2.7. - Μικροσυστοιχίες DNA [31]

Με τη χρήση της τεχνολογίας των μικροσυστοιχιών DNA κατέστη εφικτή η ταυτόχρονη παρακολούθηση εκατοντάδων γονιδίων κατά την διάρκεια σημαντικών βιολογικών διαδικασιών. Η καταγραφή και η εύρεση κρυμμένων προτύπων στα δεδομένα έκφρασης των γονιδίων βοήθησε πολύ στην κατανόηση των γενετικών λειτουργιών. Από την άλλη ο μεγάλος αριθμός γονιδίων και η πολυπλοκότητα των βιολογικών δικτύων αύξησε τις προκλήσεις για κατανόηση και μετάφραση των μαζικών αποτελεσμάτων, τα οποία συχνά προκύπτουν από χιλιάδες μετρήσεις. Για αυτό το λόγο με τη χρήση τεχνικών κατηγοριοποίησης (συσταδοποίησης), είναι εφικτή η εύρεση φυσικών δομών και προτύπων. Η ανάλυση συστάδων (cluster analysis) κατηγοριοποιεί τα δεδομένα σε ομάδες βάσει κοινών χαρακτηριστικών, ώστε τα μέλη κάθε ομάδας να μοιάζουν περισσότερο μεταξύ τους από ότι με μέλη άλλων ομάδων. Οι αλγόριθμοι κατηγοριοποίησης χρησιμοποιούνται ευρέως στην αναγνώριση βιολογικά όμοιων ομάδων γονιδίων και δειγμάτων.

Σε αντίθεση με τις κοινές τεχνικές στην γονιδιακή έρευνα η οποία επικεντρωνόταν στην τοπική μελέτη και συλλογή δεδομένων για μεμονωμένα γονίδια, με τη χρήση των μικροσυστοιχιών έγινε εφικτή η

ταυτόχρονη παρακολούθηση των επιπέδων έκφρασης χιλιάδων γονιδίων παράλληλα . Οι μικροσυστοιχίες διαφέρουν από την κλασσική δομή των πινάκων όπου αποτελούνται από σειρές με απλά δεδομένα, στο ότι υπολογίζουν την «ποσότητα του στόχου» έμμεσα μετρώντας μια άλλη φυσική ποσότητα – την ένταση του φθορισμού κάθε τοποθεσίας (spot) στον πίνακα αναπαράστασης φθορισμού. Έτσι στην πραγματικότητα οι σειρές δεδομένων που παράγονται από τις μικροσυστοιχίες είναι μονόχρωμες εικόνες. Μετασχηματίζοντας αυτές τις εικόνες μέσα στον «πίνακα γονιδιακής έκφρασης» (gene expression matrix), το οποίο αποτελεί μια δύσκολη διαδικασία, μπορούμε να αναλύσουμε αυτόν τον πίνακα και να προσπαθήσουμε να εξάγουμε από αυτόν κάποια συμπεράσματα και γνώσεις σχετικά με θεμελιώδεις βιολογικές διαδικασίες.[32],[33]

Δύο είναι οι βασικοί τρόποι πειραματισμού:

- cDNA μικροπίνακες
- ολιγονουκλεοτιδικοί πίνακες (oligo arrays).

Τα κοινά στάδια των δύο πρωτοκόλλων είναι τα παρακάτω:

Κατασκευή chip: ένα microarray είναι ένα μικρό chip (από nylon, γυαλί ή σιλικόνη) πάνω στο οποίο προσαρμόζονται μόρια DNA σε καθορισμένο πλέγμα. Κάθε κελί του πλέγματος συνδέεται με μία ακολουθία DNA.

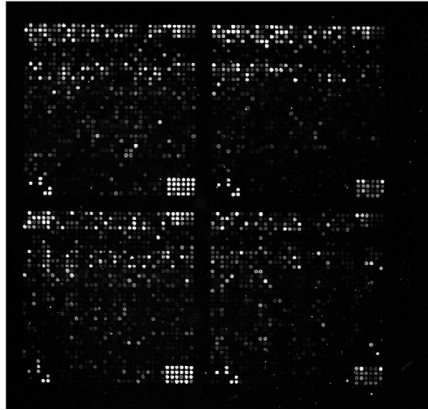
Προετοιμασία, χαρακτηρισμός και υβριδοποίηση στόχου: δύο δείγματα mRNA (ένα δείγμα ελέγχου και ένα πειραματικό δείγμα) μεταγράφονται σε cDNA (στόχοι), μαρκάρονται με ραδιοϊσότοπα ή φθοριούχες ενώσεις, και υβριδοποιούνται με τα μόρια στην επιφάνεια του πλέγματος.

Τα chips σαρώνονται για να διαβαστεί η ένταση του σήματος που εκπέμπουν από τους μαρκαρισμένους υβριδοποιημένους στόχους.

Και οι δύο μέθοδοι μετράνε την έκφραση των επιπέδων για κάθε ακολουθία DNA από τον λόγο της έντασης των σημάτων μεταξύ του δείγματος ελέγχου και του πειραματικού δείγματος.[32],[33]

Οι μικροσυστοιχίες μπορούν να χρησιμοποιηθούν για να μετρήσουν τη γονιδιακή έκφραση με πολλούς τρόπους αλλά μια από τις πιο διάσημες εφαρμογές τους είναι η σύγκριση των επιπέδων γονιδιακής έκφρασης μεταξύ δύο διαφορετικών δειγμάτων.(πχ ίδια κύτταρα κάτω από διαφορετικές συνθήκες). Οι μικροσυστοιχίες έχουν εφαρμοστεί επιτυχώς για τη λύση διαφόρων βιολογικών προβλημάτων και όσο οι συστοιχίες αυτές γίνονται περισσότερο διαθέσιμες στους ερευνητές τόσο η δημοτικότητα τους θα αυξηθεί. Για τη σωστή ανάλυση των πειραματικών δεδομένων είναι απαραίτητο να κατανοήσουμε τα πειράματα που αναπαράγουν αυτά τα δεδομένα. Τα πειράματα μικροσυστοιχιών αποτελούνται από πολλά βήματα, το καθένα από τα οποία εισάγει τον δικό του θόρυβο και τη δική του διακύμανση στο σύστημα. Το τελικό αποτέλεσμα μπορεί να επηρεαστεί από το τυχαίο σφάλμα κάθε βήματος.

Για αυτόν τον λόγο ένας σωστός πειραματικός σχεδιασμός και μια προσεκτική στατιστική ανάλυση απαιτούνται για τη σωστή ερμηνεία των δεδομένων των μικροσυστοιχιών.[34]



Εικόνα 2.8. – Σάρωση και υβριδοποίηση μικροσυστοιχίας που περιέχει περισσότερα από 5000 γονίδια κάθε σημείο χαρακτηρίζει ένα σύνολο από πανομοιότυπα μονοκλωνικά μόρια DNA που αναπαριστούν ένα μόνο γονίδιο. Η ένταση της φωτεινότητας δείχνει την ποσότητα φθορίζουσας ουσίας που περιέχει το υβριδοποιημένο mRNA

2.3. Ανάλυση γονιδιακής έκφρασης από μικροσυστοιχίες

Τα πειράματα μικροσυστοιχιών καθορίζουν την έκφραση μεταγραφής των γονιδίων ενός οργανισμού υπό διαφορετικές καταστάσεις. Η ανάλυση μικροσυστοιχιών προσπαθεί να εντοπίσει ομάδες γονιδίων που παρουσιάζουν παρόμοια συμπεριφορά υπό ορισμένες συνθήκες. Τα δεδομένα που προκύπτουν από τα πειράματα με μικροσυστοιχίες έχουν συνήθως πολύ θόρυβο. Γι' αυτό τον λόγο ο αλγόριθμος κατηγοριοποίησης πρέπει να είναι ικανός να ξεχωρίζει τη χρήσιμη πληροφορία από τον θόρυβο. Τα τελευταία χρόνια, μία από τις κύριες μεθόδους για την ανάλυση των δεδομένων μικροσυστοιχιών είναι η διπλή κατηγοριοποίηση, μια πολύ ανεπτυγμένη τεχνική. Η τεχνική biclustering υπερέχει της παραδοσιακής τεχνικής ομαδοποίησης λόγω των δύο κύριων χαρακτηριστικών της: την ταυτόχρονη ομαδοποίηση γονιδίων και συνθηκών και την επικάλυψη. Ταυτόχρονη ομαδοποίηση σημαίνει ότι οι ομάδες διπλής κατηγοριοποίησης (οι ομάδες που βρέθηκαν με τους αλγορίθμους biclustering) ομαδοποιούν γονίδια με παρόμοια συμπεριφορά κάτω από έναν ορισμένο αριθμό συνθηκών (έτσι, το Bicluster θα ομαδοποιεί γονίδια και συνθήκες), ενώ οι παραδοσιακές τεχνικές ομαδοποίησης συγκεντρώνουν μόνο γονίδια με παρόμοια συμπεριφορά σε όλες τις συνθήκες(ή αντιστρόφως). Αυτό το χαρακτηριστικό καθιστά τις ομάδες διπλής κατηγοριοποίησης κατάλληλες για την αποτύπωση της βιολογικής συμπεριφοράς σε διάφορες περιστάσεις, για παράδειγμα, όταν μια ενδιαφέρουσα κυτταρική διαδικασία είναι ενεργή μόνο σε ένα υποσύνολο συνθηκών. Αν και είναι ασυνήθιστο το γεγονός ότι τα υποσύνολα των γονιδίων που ομαδοποιούνται με δύο διαφορετικές ομάδες να διατέμνονται, η αλληλοεπικάλυψη είναι ένα εγγενές χαρακτηριστικό των Biclusters. Αν δύο ομάδες διπλής κατηγοριοποίησης $B1$ και $B2$ που ομαδοποιούν τα γονίδια $G1$ και $G2$ και τις συνθήκες $C1$ και $C2$, αντίστοιχα, έχουν $G1 \cap G2 \neq \emptyset$ και $C1 \cap C2 \neq \emptyset$ λέγεται ότι οι $B1$ και $B2$ αλληλεπικαλύπτονται. Η αλληλοεπικάλυψη δίνει στις ομάδες διπλής κατηγοριοποίησης την ευελιξία να αναπαριστούν βιολογικές καταστάσεις όπως είναι τα γονίδια που συμμετέχουν σε πολλαπλά μονοπάτια που λειτουργούν κάτω από ένα υποσύνολο συνθηκών.[35] Τις τελευταίες δεκαετίες παρουσιάζεται ραγδαία αύξηση του όγκου των βιολογικών δεδομένων και η ανάγκη για τον μετασχηματισμό της αδόμητης πληροφορίας που προκύπτει από αυτά τα δεδομένα σε γνώση. Αυτός ήταν ο κύριος λόγος δημιουργίας του νέου επιστημονικού πεδίου της Βιοπληροφορικής (Bioinformatics) που θα αναλυθεί στο επόμενο υποκεφάλαιο.

2.4. ΒιοΠληροφορική

Η Βιοπληροφορική αποτελεί στην ουσία σύμπραξη επιστημονικών κλάδων όπως η βιολογία, η βιοχημεία, τα μαθηματικά, η επιστήμη των υπολογιστών, η τεχνολογία πληροφοριών και άλλα. Το έργο της Βιοπληροφορικής συνίσταται στην έρευνα, ανάπτυξη, ή εφαρμογή υπολογιστικών εργαλείων και προσεγγίσεων για την επέκταση της χρήσης δεδομένων βιολογίας, ιατρικής, συμπεριφοράς ή υγείας, συμπεριλαμβανομένων εκείνων για την απόκτηση, αποθήκευση, οργάνωση, αρχειοθέτηση, ανάλυση, ή οπτικοποίηση αυτών των δεδομένων.[36] Για παράδειγμα, με τη βοήθεια της Βιοπληροφορικής μέσα από τη χρήση εξελιγμένων συστημάτων βάσεων δεδομένων, αλγορίθμων ομαδοποίησης και μια σειρά από στατιστικά εργαλεία αναγνωρίζεται η ταυτότητα των βιολογικών δεδομένων, κατανοούνται βιολογικοί μηχανισμοί και ασθένειες, και σχεδιάζονται φάρμακα.[26]

Ένα από τα πιο σημαντικά ίσως ερευνητικά πεδία της Βιοπληροφορικής είναι η διαχείριση και εξόρυξη γνώσης με τις βιολογικές βάσεις δεδομένων να αποτελούν θεμελιώδες τμήμα τους.[31],[37] Οι βιολογικές βάσεις δεδομένων είναι μεγάλα, οργανωμένα συστήματα δεδομένων που συνδέονται συνήθως με κατάλληλο λογισμικό για την ενημέρωση, αναζήτηση, και ανάκτηση στοιχείων των δεδομένων που έχουν αποθηκευθεί στο σύστημα. Από αυτά τα αποθετήρια γνώσης, ένας ερευνητής μπορεί να αντλήσει τις κατάλληλες πληροφορίες (MEDLINE, GenBank)[38],[39] ή τα δεδομένα (GEO)[40] πάνω στα οποία θα εφαρμόσει την ανάλυση και τις μεθοδολογικές προσεγγίσεις για την επίλυση ενός ή πολλαπλών ερωτημάτων. Γενικά, οι βιολογικές βάσεις δεδομένων μπορούν να διακριθούν σε δυο κύριες κατηγορίες, στις πρωτογενείς και τις δευτερογενείς βάσεις δεδομένων. Η πρώτη κατηγορία απαρτίζεται από βάσεις δεδομένων με πρωτογενή στοιχεία, τα οποία απορρέουν από διάφορες πειραματικές μεθοδολογίες, ενώ η δεύτερη από βάσεις δεδομένων με στοιχεία που δημιουργούνται από ταξινομήσεις των πρωτογενών δεδομένων και είναι χρήσιμα για αναλυτικούς σκοπούς. Ενδεικτικές βιολογικές βάσεις δεδομένων της πρώτης και δεύτερης κατηγορίας είναι οι ακόλουθες:

Πρωτογενείς βάσεις δεδομένων

Βάσεις δεδομένων ακολουθιών

Βάσεις νουκλεοτιδικών ακολουθιών ελεύθερα διαθέσιμες, οι οποίες συνεργάζονται μεταξύ τους ανταλλάσσοντας εγγραφές και δημιουργώντας κοινούς κανόνες για την ταξινόμηση και το σχολιασμό των δεδομένων:

- DNA Data Bank of Japan (DDBJ, <http://www.ddbj.nig.ac.jp/>) στο Center for Information Biology (CIB).
- GenBank (<http://www.ncbi.nlm.nih.gov/Genbank/>) στο National Center for Biotechnology Information (NCBI).
- EMBL_Bank (<http://www.ebi.ac.uk/embl/index.html>) στο European Bioinformatics Institute (EBI).

Εξειδικευμένες βάσεις δεδομένων που συνδυάζουν τα δεδομένα γονιδιωματικών ακολουθιών και το σχολιασμό τους με άλλα στοιχεία για τα συγκεκριμένα είδη.

- Ensembl (<http://www.ensembl.org/index.html>) αποτέλεσμα συνεργασίας του EBI και του Wellcome Trust Sanger Institute.

- Entrez Genomes (<http://www.ncbi.nlm.nih.gov/sites/entrez?db=genome>) στο National Center for Biotechnology Information (NCBI).

Βάσεις πρωτεϊνικών ακολουθιών που παρέχουν υψηλό επίπεδο σχολιασμού (όπως περιγραφή της λειτουργίας μιας πρωτεΐνης, μετα-μεταφραστικές τροποποιήσεις κ.λπ.) και διασυνδέσεις με άλλες βάσεις δεδομένων:

- Swiss-Prot (<http://www.expasy.ch/sprot/>).
- TrEMBL (<http://www.ebi.ac.uk/trembl/>).
- UniProt (<http://www.ebi.ac.uk/uniprot/index.html/>) που προέκυψε από τη συνεργασία των Swiss-Prot, TrEMBL και PIR.

Βάσεις δεδομένων γονιδιακής έκφρασης

Με την εξέλιξη της τεχνολογίας και τη δημιουργία νέων οικονομικότερων τσιπ μικροσυστοιχιών, αλλά και με την εμφάνιση των τεχνολογιών αλληλούχισης νέας γενιάς (next generation sequencing), τα πειράματα ανάλυσης γονιδιακής έκφρασης πραγματοποιούνται με μεγαλύτερο ρυθμό δίνοντας μεγάλο όγκο δεδομένων και δημιουργώντας την ανάγκη αποθήκευσης και ανάλυσης αυτών των δεδομένων. Οι βάσεις δεδομένων γονιδιακής έκφρασης επιτρέπουν την καταχώρηση αποτελεσμάτων από πειράματα μικροσυστοιχιών, ενώ κάποιες από αυτές προσφέρουν και επιπλέον εργαλεία ανάλυσης.¹¹

- GEO (Gene Expression Omnibus, <https://www.ncbi.nlm.nih.gov/geo/>). Η βάση δεδομένων GEO του NCBI παρέχει δεδομένα γονιδιακής έκφρασης, τόσο από μικροσυστοιχίες, όσο και από την αλληλούχιση νέας γενιάς (next generation sequencing), ενώ στον ίδιο ιστότοπο υπάρχουν διαθέσιμα και κάποια διαδικτυακά εργαλεία που επιτρέπουν απλές αναλύσεις των δεδομένων της βάσης. Τα δεδομένα υπάρχουν τόσο σε ακατέργαστη (raw) όσο και σε επεξεργασμένη μορφή (με κανονικοποιήσεις, κ.ά.).[36]

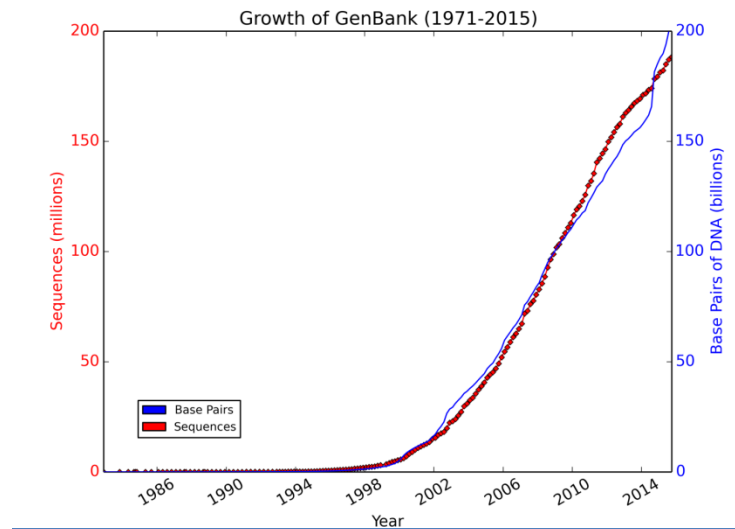
Βιβλιογραφικές βάσεις δεδομένων

- Η MEDLINE (US National Library of Medicine) είναι η βιβλιογραφική βάση δεδομένων της NLM (National Library of Medicine, USA) που καλύπτει τους τομείς της ιατρικής, της υγειονομικής περίθαλψης, των προκλινικών επιστημών, της βιολογίας καθώς και θεμάτων βιοϊατρικής τεχνολογίας. Περιέχει βιβλιογραφικές παραπομπές και περιλήψεις άρθρων από περισσότερα από 4800 βιοϊατρικά περιοδικά που δημοσιεύονται στις Ηνωμένες Πολιτείες και σε 70 άλλες χώρες. Η πρόσβαση στα περιεχόμενά της γίνεται από την ελεύθερη (δωρεάν) μηχανή αναζήτησης PubMed (<http://www.ncbi.nlm.nih.gov/PubMed/>), που είναι μέρος του συστήματος ανάκτησης πληροφοριών Entrez. Μια επιπρόσθετη βάση δεδομένων του Entrez είναι η MeSH (Medical Subject Headings) η οποία είναι ένα λεξικό όρων βιολογίας και ιατρικής. Μέσω της MeSH είναι δυνατή η άντληση ορισμών αλλά και η ανάκληση άρθρων από την PubMed.

Δευτερογενείς βάσεις δεδομένων

- SCOP (Structural Classification of Proteins, <http://scop.mrc-lmb.cam.ac.uk/scop/>). Η βάση SCOP έχει ως στόχο την ανάλυση των δομικών και εξελικτικών σχέσεων που παρατηρούνται μεταξύ όλων των πρωτεϊνών γνωστής δομής καταχωρημένων στην PDB (Protein Data Bank, <http://www.rcsb.org/pdb/home/home.do>). Η ταξινόμηση των πρωτεϊνών βασίζεται σ' αυτές τις δομικές και εξελικτικές σχέσεις. Τα βασικά επίπεδα ταξινόμησης είναι τέσσερα: η οικογένεια (family), η υπερ-οικογένεια (superfamily), η αναδίπλωση (fold), και η τάξη (class).

Με βάση τα παραπάνω, μπορεί να γίνει κατανοητή η παρακάτω Εικόνα, στην οποία παρουσιάζεται η αλματώδης αύξηση των βιολογικών δεδομένων από το 1971 έως το 2015.



Εικόνα 2.9. - Αύξηση βιολογικών δεδομένων από το 1971 έως το 2015 [41]

2.5. Αλγόριθμοι Κατηγοριοποίησης

Με τον όρο κατηγοριοποίηση αναφερόμαστε στη διαδικασία ομαδοποίησης αντικειμένων σε κλάσεις διαφορετικές κλάσεις, ξένες μεταξύ τους, με την προϋπόθεση ότι τα αντικείμενα τα οποία ανήκουν σε μία κλάση να παρουσιάζουν μεγαλύτερη ομοιότητα μεταξύ τους από ότι αντικείμενα άλλων κλάσεων. Σε δεδομένα γονιδιακής έκφρασης, ο διαχωρισμός των κατηγοριών βασίζεται σε μια σειρά από κριτήρια όπως είναι:

Η **κατηγοριοποίηση βάσει γονιδίων**, όπου τα γονίδια θεωρούνται ως αντικείμενα και τα δείγματα (συνθήκες) ως χαρακτηριστικά.

Η **κατηγοριοποίηση βάσει δειγμάτων**, όπου τα δείγματα θεωρούνται ως αντικείμενα, και τα χαρακτηριστικά πάνω στα οποία γίνεται η κατηγοριοποίηση είναι τα γονίδια.

Η διαφορά μεταξύ των δύο τύπων κατηγοριοποίησης είναι ότι βασίζονται σε διαφορετικά χαρακτηριστικά για να πετύχουν την κατηγοριοποίηση των δεδομένων. Παρ' όλα αυτά, αλγόριθμοι όπως ο K-means και οι ιεραρχικές προσεγγίσεις, μπορούν να εφαρμοστούν και στις δύο περιπτώσεις. Νεότερες όμως προσεγγίσεις στην μοριακή βιολογία θεωρούν ότι μόνο ένα μικρό μέρος των γονιδίων παίρνει μέρος στις κυτταρικές λειτουργίες. Έτσι έγινε απαραίτητη η εισαγωγή μίας νέας μεθόδου κατηγοριοποίησης:

Η **κατηγοριοποίηση υποδιαστημάτων (subspace clustering)**: Σ' αυτή την περίπτωση η κατηγοριοποίηση εφαρμόζεται σε ένα υποσύνολο των γονιδίων βάση ενός υποσυνόλου των δειγμάτων. Σε αυτή την προσέγγιση τα δείγματα και τα γονίδια αντιμετωπίζονται συμμετρικά, οπότε άλλοτε τα γονίδια και άλλοτε τα δείγματα θεωρούνται αντικείμενα ή χαρακτηριστικά.

Καθεμία από τις τρεις κατηγοριοποιήσεις βασίζεται σε διαφορετικό τρόπο ανάλυσης του πίνακα έκφρασης των γονιδίων και αντιμετωπίζει διαφορετικές προκλήσεις.[32]

Στην παρακάτω ενότητα θα παρουσιαστούν οι αλγόριθμοι κατηγοριοποίησης βάσει γονιδίων, δειγμάτων και υποδιαστημάτων.[42]

Κατηγοριοποίηση Βάσει Γονιδίων

Η ιδιαιτερότητα των δεδομένων έκφρασης γονιδίων και οι ιδιόμορφες απαιτήσεις του βιολογικού τομέα, δημιουργούν στην κατηγοριοποίηση γονιδίων προκλήσεις και προβλήματα, αρκετά από τα οποία εκκρεμούν.

Η ανάλυση συστάδων είναι το πρώτο βήμα για την εξόρυξη γνώσης (data mining). Ο κύριος στόχος της κατηγοριοποίησης δεδομένων έκφρασης γονιδίων είναι η εύρεση των δομών των δεδομένων και των κατανομών των δεδομένων. Χαρακτηριστικό ενός αποτελεσματικού αλγορίθμου κατηγοριοποίησης είναι η όσο το δυνατόν μικρότερη δυνατή εξάρτηση από προηγούμενη γνώση, η οποία συνήθως δεν είναι καν διαθέσιμη πριν την ανάλυση συστάδων.

Τα δεδομένα που προκύπτουν από τα πειράματα με μικροσυστοιχίες έχουν συνήθως πολύ θόρυβο. Γι' αυτό τον λόγο ο αλγόριθμος κατηγοριοποίησης πρέπει να είναι ικανός να ξεχωρίζει τη χρήσιμη πληροφορία από τον θόρυβο.

Μεταξύ των δεδομένων γονιδιακής έκφρασης υπάρχουν ισχυρές συνδέσεις, με τις διάφορες κατηγορίες να είναι επικαλυπτόμενες ή ακόμα να εξαπλώνεται η μία πάνω στην άλλη. Οπότε οι αλγόριθμοι gene-based κατηγοριοποίησης πρέπει να μπορούν να το αντιμετωπίσουν.

Τέλος οι χρήστες μπορεί να μην ενδιαφέρονται μόνο για τις συστάδες γονιδίων, αλλά ακόμα και για τις σχέσεις μεταξύ των συστάδων, σχέσεις γονιδίων μέσα στην ίδια συστάδα κτλ. Γι' αυτό το λόγο, ένας αλγόριθμος κατηγοριοποίησης θα ήταν καλό να μην σταματά στον διαμερισμό των δεδομένων, αλλά να παρέχει και γραφική αναπαράσταση της δομής των συστάδων.[43]

Οι βασικότεροι αλγόριθμοι κατηγοριοποίησης γονιδίων είναι οι παρακάτω:[32]

- K-means.
- Αυτό-οργανωμένοι χάρτες S.O.M. (Self-organized maps).
- Ιεραρχική κατηγοριοποίηση.
- Γραφοθεωρητικές μέθοδοι κατηγοριοποίησης.
- Κατηγοριοποίηση βάσει μοντέλου.
- Κατηγοριοποίηση βάσει πυκνότητας (ιεραρχική προσέγγιση DHC).

Κατηγοριοποίηση Βάσει Δειγμάτων

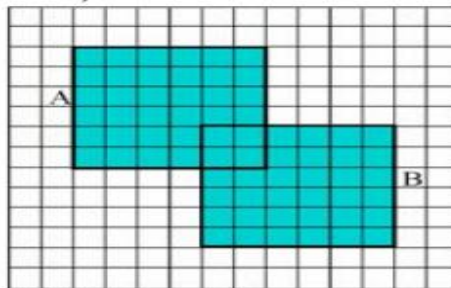
Ο πίνακας έκφρασης γονιδίων τις πιο πολλές φορές αποτελείται από μακροσκοπικούς φαινοτύπους δειγμάτων οι οποίοι συνδέονται με ασθένειες ή επίδραση φαρμάκων, όπως δείγματα ασθενειών, φυσιολογικά δείγματα και δείγματα θεραπείας. Ο κύριος στόχος της κατηγοριοποίησης βάσει δειγμάτων είναι η ανακάλυψη φαινοτυπικών δομών ή υπό-δομών των δειγμάτων. Προηγούμενες μελέτες [44] έχουν δείξει ότι ο φαινότυπος των δειγμάτων μπορεί να διακριθεί μόνο μέσα από μία μικρή υπομονάδα γονιδίων των οποίων τα

επίπεδα έκφρασης σχετίζονται στενά με την κλάση διάκρισης. Τα γονίδια αυτά αποκαλούνται και πληροφοριακά γονίδια. Τα υπόλοιπα γονίδια που απομένουν στον πίνακα έκφρασης γονιδίων θεωρούνται απλά θόρυβος. Οι κλασικές μέθοδοι κατηγοριοποίησης, όπως αυτές που αναφέρθηκαν πριν (π.χ. αυτοοργανωμένοι χάρτες, K-means κτλ.), μπορούν να εφαρμοστούν θεωρώντας όλα τα γονίδια σαν χαρακτηριστικά, όμως ο υψηλός λόγος σήματος-θορύβου (που είναι της τάξης του 1/10) μειώνει δραματικά την αξιοπιστία των αποτελεσμάτων. Για τον λόγο αυτό έχουν προταθεί καινούριοι αλγόριθμοι, οι οποίοι χωρίζονται σε δύο κύριες κατηγορίες:[44]

- Την εποπτευόμενη ανάλυση και
- Την μη-εποπτευόμενη ανάλυση.

Κατηγοριοποίηση Υποδιαστημάτων (Subspace Clustering)

Οι παραπάνω τύποι αλγορίθμων κατηγοριοποίησης θεωρούνται όλοι δείγματα “ολικής κατηγοριοποίησης” δεδομένου ότι για ένα συγκεκριμένο σύνολο δεδομένων προς κατηγοριοποίηση, ο χώρος των χαρακτηριστικών είναι καθορισμένος και μοιράζεται από όλες τις εξαγόμενες συστάδες, με τις συστάδες να είναι αποκλειστικές και εξαντλητικές. Παρ’ όλα αυτά, στην μοντέρνα μοριακή βιολογία είναι γνωστό ότι μόνο ένα μικρό μέρος των γονιδίων παίρνει μέρος στις κυτταρικές λειτουργίες. Επίσης, είναι σύνηθες ένα γονίδιο να έχει ενεργό ρόλο σε πολλαπλές λειτουργίες οι οποίες δεν είναι απαραίτητο να συνεργάζονται, οπότε ένα γονίδιο μπορεί να ανήκει σε πολλές συστάδες, ή και σε καμία. Μία σειρά μεθόδων κατηγοριοποίησης υποδιαστημάτων προτάθηκε πρόσφατα, ώστε να κάνει κατανοητή την σχέση που παρουσιάζεται στα μπλοκ -ένας υποπίνακας που ορίζεται από ένα υποσύνολο γονιδίων και ένα υποσύνολο δειγμάτων- μέσα στον πίνακα έκφρασης γονιδίων. Η κατηγοριοποίηση υποδιαστημάτων αρχικά προτάθηκε από τον Agrawal (1998) [43] για να βρει υποσύνολα αντικειμένων ώστε τα αντικείμενα να εμφανίζονται σαν μία συστάδα σε έναν υποχώρο που ορίζεται από ένα υποσύνολο των χαρακτηριστικών.



Εικόνα 2.10. – Υποδιάστημα συστάδων A & B [45]

Στην παραπάνω Εικόνα 2.10 παρουσιάζεται ένα παράδειγμα του υποσυνόλου των συστάδων (A, B) απλωμένου στον πίνακα έκφρασης γονιδίων. Στην κατηγοριοποίηση αυτή είναι δυνατόν το υποσύνολο των χαρακτηριστικών να διαφέρει από συστάδα σε συστάδα. Δύο υποσύνολα συστάδων μπορούν να έχουν ως γνώρισμα τον διαμοιρασμό κοινών αντικειμένων και χαρακτηριστικών, ενώ από την άλλη κάποια αντικείμενα μπορεί να μην ανήκουν σε κανένα υποσύνολο συστάδας. Για ένα πίνακα έκφρασης γονιδίων που έχει n γονίδια και m δείγματα, η πολυπλοκότητα για πλήρη συνδυασμό τους είναι 2^{n+m} οπότε το πρόβλημα της συνολικά βέλτιστης επιλογής μπλοκ είναι δυσκολίας NP. Οι μέθοδοι κατηγοριοποίησης υποδιαστημάτων

συνήθως ορίζουν μοντέλα για να περιγράψουν το επιθυμητό μπλοκ και στην συνέχεια υιοθετούν κάποια μέθοδο να ψάξει στον χώρο των δειγμάτων-γονιδίων.[33]

2.6. Τεχνικές Κατηγοριοποίησης

Σε έναν πίνακα γονιδιακής έκφρασης, αν δύο γονίδια σχετίζονται (έχουν παρόμοιες λειτουργίες), τα γονιδιακά τους προφίλ παρουσιάζουν ομοιότητα. Η ομοιότητα υπολογίζεται με τη χρήση της ευκλείδειας απόστασης ή της υψηλής συσχέτισης. Η κλασσική ομαδοποίηση ή διαφορετικά clustering ομαδοποιεί τα στοιχεία ενός πίνακα, ώστε τα στοιχεία της ίδιας ομάδας να είναι παραπλήσια ενώ από την άλλη, οι ομάδες να είναι σαφώς χωρισμένες αναμεταξύ τους. Στη μέθοδο ομαδοποίησης, οι δύο κύριες τεχνικές βασίζονται στα ακόλουθα δυο κριτήρια: **στην κατάσταση**, όπου εδώ παίζει ρόλο η ομοιότητα της γονιδιακής έκφρασης, και **στο γονίδιο**, όπου η ομαδοποίηση είναι αλληλεξάρτηση της γονιδιακής έκφρασης.

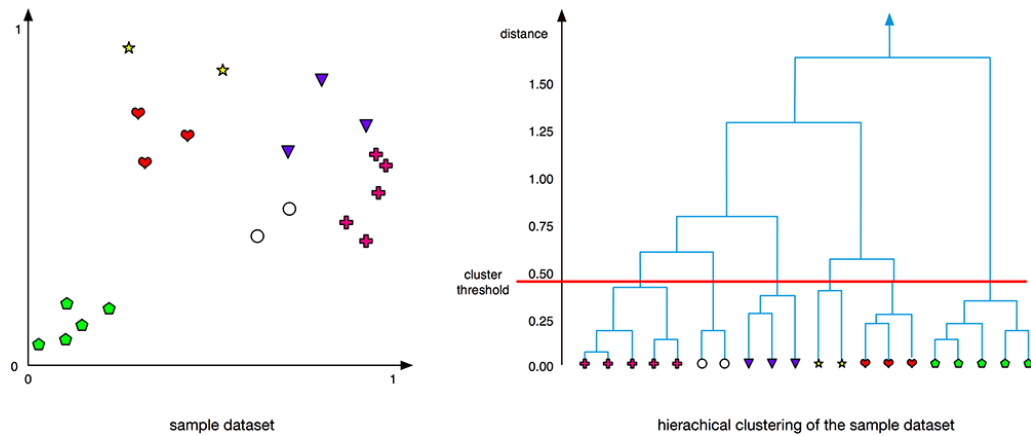
Μια από τους πιο δημοφιλείς τύπους αλγορίθμων clustering είναι:

Η ιεραρχική ομαδοποίηση (Hierarchical clustering) - Η ιεραρχική ομαδοποίηση εφαρμόζεται σχετικά εύκολα και περιλαμβάνει δύο είδη clustering, την ενωτική (συσσωρευτική) agglomerative και τη διαιρετική (διαχωριστική) (divisive). Στην αρχή της εφαρμογής ο αριθμός των ομάδων δεν είναι γνωστός. Κάθε αντικείμενο μπορεί να ομαδοποιηθεί μόνο μια φορά και τα γονίδια ομαδοποιούνται όλα, παρόλο που κάποιες ομάδες μπορεί να είναι πιο αδύναμες και λιγότερο ακριβείς.

Τα κύρια μειονεκτήματα αυτής της τεχνικής ομαδοποίησης είναι δύο. Αρχικά τα γονίδια ομαδοποιούνται με βάση την έκφρασή τους σε όλες τις καταστάσεις, αν και μια χρήσιμη πληροφορία μπορεί να κρύβεται μόνο σε ένα επιμέρους κομμάτι των καταστάσεων. Δεύτερον, κάθε γονίδιο ομαδοποιείται μόνο σ' ένα cluster, παρ' όλα αυτά ένα γονίδιο μπορεί να συμμετέχει σε πάνω από μία κυτταρικές διεργασίες.[46], [47]

Ιεραρχική Ομαδοποίηση (Hierarchical Clustering)

Αναλυτικότερα, η βασική ιδέα της τεχνικής της ιεραρχικής ομαδοποίησης (Hierarchical clustering), συνδέεται με την δημιουργία μιας ιεραρχίας από clusters η οποία παρουσιάζει τις σχέσεις (οι οποίες βασίζονται στην ομοιότητα) μεταξύ ανεξάρτητων μελών και συγχωνευμένων clusters των δεδομένων.[48]



Εικόνα 2.11. - Ιεραρχική ομαδοποίηση (Hierarchical clustering) [48]

Στην **ιεραρχική ομαδοποίηση** (Hierarchical clustering) πληρούνται οι ακόλουθοι κανόνες:

- Το συνολικό σύνολο δεδομένων (dataset) αντιπροσωπεύεται από μία ρίζα.
- Ένα συγκεκριμένο στοιχείο του συνόλου δεδομένων (dataset) αναπαρίσταται από ένα φύλλο.
- Ο κόμβος αντιπροσωπεύει την ένωση των στοιχείων στο υποδέντρο.
- Το ύψος ενός κόμβου αναπαριστά την απόσταση δύο παιδιών κόμβων.[49]

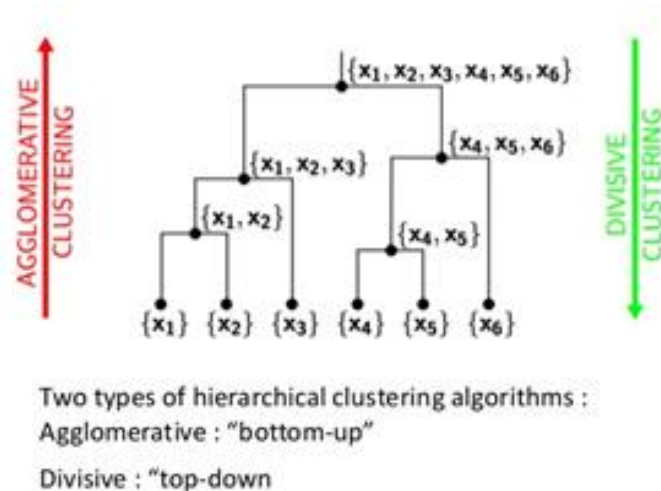
Όπως αναφέρθηκε, οι ιεραρχικές μέθοδοι χωρίζονται σε ενωτικές και διαχωριστικές:

Ενωτικές ή από κάτω προς τα πάνω (Agglomerative ή bottom-up):

- Ξεκινούν τοποθετώντας κάθε στοιχείο του πίνακα σε ένα ξεχωριστό cluster.
 - Σε κάθε βήμα συγχωνεύουν τα δύο πιο όμοια clusters.
 - Σταματούν όταν όλα τα στοιχεία ανήκουν σε ένα μόνο cluster ή όταν
 - Τερματιστεί κάποιο συγκεκριμένο κριτήριο.

Διαχωριστικές ή από πάνω προς τα κάτω (Divisive ή top-down):

- Ξεκινούν τοποθετώντας όλα τα στοιχεία του πίνακα σε ένα cluster.
 - Σε κάθε επανάληψη χωρίζουν ένα cluster σε δύο καινούργια.
- Ολοκληρώνονται όταν όλα τα στοιχεία ανήκουν στο δικό τους cluster ή όταν
 - Τερματιστεί κάποιο συγκεκριμένο κριτήριο.[50],[51]



Εικόνα 2.12. – Ενωτική (Agglomerative) και Διαχωριστική (Divisive) Ομαδοποίηση [52]

Συσταδοποίηση δυο-τρόπων (Biclustering)

Η συσταδοποίηση δυο-τρόπων ή αλλιώς **διπλή κατηγοριοποίηση** αποτελεί μια μέθοδο ομαδοποίησης που επιτρέπει την ταυτόχρονη κατηγοριοποίηση των γραμμών και των στηλών ενός πίνακα. Ο όρος εισήχθη για πρώτη φορά από τον Mirkin (1996),[7] αν και η συγκεκριμένη τεχνική κατηγοριοποίησης παρουσιάστηκε πολύ πιο νωρίς από τον Hartigan (1975).[6] Οι περισσότεροι από αυτούς τους αλγορίθμους χρησιμοποιούνται σε εφαρμογές της βιοπληροφορικής.

Δημιουργία συστάδων δυο-τρόπων

Δοθέντος ενός πίνακα μεγέθους X, Y , ο αλγόριθμος διπλής κατηγοριοποίησης δημιουργεί συστάδες με τον εξής τρόπο: ένα υποσύνολο γραμμών του πίνακα που παρουσιάζει παρόμοια συμπεριφορά κατά πλάτος ενός υποσυνόλου στηλών και το αντίθετο θα τοποθετηθεί σε μια νέα συστάδα.

Είδη αλγορίθμων διπλής κατηγοριοποίησης

Υπάρχουν διάφορα είδη αλγορίθμων διπλής κατηγοριοποίησης τα οποία παρουσιάζονται οπτικά και μέσω των σχημάτων που ακολουθούν. Συγκεκριμένα τα είδη αυτών των αλγορίθμων είναι τα εξής: 1. διπλή κατηγοριοποίηση με σταθερές τιμές, 2. διπλή κατηγοριοποίηση με σταθερές τιμές σε γραμμές και στήλες, και 3. διπλή κατηγοριοποίηση με συναφείς τιμές.[53]

2.0	2.0	2.0	2.0	2.0
2.0	2.0	2.0	2.0	2.0
2.0	2.0	2.0	2.0	2.0
2.0	2.0	2.0	2.0	2.0
2.0	2.0	2.0	2.0	2.0

α

1.0	1.0	1.0	1.0	1.0
2.0	2.0	2.0	2.0	2.0
3.0	3.0	3.0	3.0	3.0
4.0	4.0	4.0	4.0	4.0
5.0	5.0	5.0	5.0	5.0

β

1.0	2.0	3.0	4.0	5.0
1.0	2.0	3.0	4.0	5.0
1.0	2.0	3.0	4.0	5.0
1.0	2.0	3.0	4.0	5.0
1.0	2.0	3.0	4.0	5.0

γ

1.0	4.0	5.0	0.0	1.5
4.0	7.0	8.0	3.0	4.5
3.0	6.0	7.0	2.0	3.5
5.0	8.0	9.0	4.0	5.5
2.0	5.0	6.0	1.0	7.5

δ

1.0	0.5	2.0	0.2	0.8
2.0	1.0	4.0	0.4	1.6
3.0	1.5	6.0	0.6	2.4
4.0	2.0	8.0	0.8	3.2
5.0	2.5	10	1.0	4.0

ε

Εικόνα 2.13. - Είδη αλγορίθμων διπλής κατηγοριοποίησης : α) διπλή κατηγοριοποίηση σταθερής τιμής β) διπλή κατηγοριοποίηση σταθερής σειράς γ) διπλή κατηγοριοποίηση σταθερής στήλης δ) διπλή κατηγοριοποίηση συνεκτικής τιμής (προσθετική) ε) διπλή κατηγοριοποίηση συνεκτικής τιμής (μοντέλο) [54]

Αλγόριθμοι διπλής κατηγοριοποίησης

Το εύρος των αλγορίθμων αυτής της κατηγορίας που χρησιμοποιούνται σε εφαρμογές της βιοπληροφορικής είναι αρκετά μεγάλο. Μερικά παραδείγματα τέτοιων αλγορίθμων είναι τα ακόλουθα: ο αλγόριθμος συσταδοποίησης μπλοκ, ο CTWC, ο ITWC, ο δ -Bicluster, ο d-pCluster, ο d-pattem, ο FLOC, ο OPC, ο Plaid Model, ο OPSMs, ο Gibbs, ο SAMBA, ο RoBA, ο Crossing Minimization, ο cMonkey, ο PRMs και ο DCC. Ένα μεγάλο πρόβλημα που προκύπτει με όλους τους αλγορίθμους που ανήκουν σε αυτή την κατηγορία είναι ο βαθμός αξιοπιστίας των αποτελεσμάτων που δίνουν καθώς η συσταδοποίηση 2-τρόπων επιτρέπει την επικάλυψη μεταξύ διαφορετικών συστάδων.

Οι περισσότεροι από τους αναφερόμενους αλγορίθμους βασίζονται σε μια από τις πέντε (5) παρακάτω μεθόδους για την εύρεση των Biclusters:

- Επαναληπτικός συνδυασμός γραμμών και στηλών.
- Διαίρει και βασίλευε.
- Άπληστη επαναληπτική αναζήτηση.
- Εξαντλητική απαρίθμηση.
- Διανομή παραμέτρων αναγνώρισης.

Το γεγονός ότι κάποιος από αυτούς τους αλγορίθμους δεν είναι ντετερμινιστικός και συνάμα ότι το πρόβλημα εμπεριέχει μη επιβλεπόμενη κατηγοριοποίηση, δυσκολεύουν τον εντοπισμό σφαλμάτων στα αποτελέσματα. Μια προσέγγιση για την επίλυση του συγκεκριμένου προβλήματος είναι χρήση πολλών αλγορίθμων διπλής κατηγοριοποίησης από τους οποίους, η πλειονότητα θα ψηφίζει για να αποφασιστεί ποιο είναι το καλύτερο αποτέλεσμα κατηγοριοποίησης. Παρ' όλα αυτά οι συγκεκριμένοι αλγόριθμοι βρίσκουν αρκετά καλή εφαρμογή στο χώρο της βιολογίας και γενικότερα της βιοπληροφορικής.[3],[6],[10]

Δεδομένου ότι η παρούσα εργασία επικεντρώνεται στη μέθοδο της διπλής κατηγοριοποίησης, οι αλγόριθμοι και τα είδη της αυτά θα περιγραφούν αναλυτικότερα στο κεφάλαιο που ακολουθεί (Κεφάλαιο 3).

3

Αλγόριθμοι και Τεχνικές Διπλής Κατηγοριοποίησης (Biclustering)

3.1. Τύποι Διπλής Κατηγοριοποίησης

Ένα από τα βασικά κριτήρια για την ορθή επιλογή του κατάλληλου αλγορίθμου διπλής κατηγοριοποίησης είναι το είδος των Biclusters που αναμένεται ως αποτέλεσμα. Έτσι, ανάλογα με τον τύπο του αλγόριθμου, οι τέσσερις κύριοι τύποι είναι:

α) Biclusters με σταθερές τιμές,

β) Biclusters με σταθερές τιμές σε γραμμές

ή στήλες,

γ) Biclusters με συναφείς τιμές, και

δ) Biclusters με συναφείς εξελίξεις.

Στη διπλανή εικόνα βλέπουμε την α, τη β και τη δ περίπτωση.

<table><tr><td>a</td><td>a</td><td>a</td><td>a</td></tr><tr><td>a</td><td>a</td><td>a</td><td>a</td></tr><tr><td>a</td><td>a</td><td>a</td><td>a</td></tr><tr><td>a</td><td>a</td><td>a</td><td>a</td></tr></table> <p>Constant values</p>	a	a	a	a	a	a	a	a	a	a	a	a	a	a	a	a	<table><tr><td>a</td><td>a</td><td>a</td><td>a</td></tr><tr><td>a+i</td><td>a+i</td><td>a+i</td><td>a+i</td></tr><tr><td>a+j</td><td>a+j</td><td>a+j</td><td>a+j</td></tr><tr><td>a+k</td><td>a+k</td><td>a+k</td><td>a+k</td></tr></table> <p>Constant values on rows</p>	a	a	a	a	a+i	a+i	a+i	a+i	a+j	a+j	a+j	a+j	a+k	a+k	a+k	a+k	<table><tr><td>a</td><td>a+i</td><td>a+j</td><td>a+k</td></tr><tr><td>a</td><td>a+i</td><td>a+j</td><td>a+k</td></tr><tr><td>a</td><td>a+i</td><td>a+j</td><td>a+k</td></tr><tr><td>a</td><td>a+i</td><td>a+j</td><td>a+k</td></tr></table> <p>Constant values on cols</p>	a	a+i	a+j	a+k	a	a+i	a+j	a+k	a	a+i	a+j	a+k	a	a+i	a+j	a+k
a	a	a	a																																															
a	a	a	a																																															
a	a	a	a																																															
a	a	a	a																																															
a	a	a	a																																															
a+i	a+i	a+i	a+i																																															
a+j	a+j	a+j	a+j																																															
a+k	a+k	a+k	a+k																																															
a	a+i	a+j	a+k																																															
a	a+i	a+j	a+k																																															
a	a+i	a+j	a+k																																															
a	a+i	a+j	a+k																																															
<table><tr><td>a</td><td>b</td><td>c</td><td>d</td></tr><tr><td>a+i</td><td>b+i</td><td>c+i</td><td>d+i</td></tr><tr><td>a+j</td><td>b+j</td><td>c+j</td><td>d+j</td></tr><tr><td>a+k</td><td>b+k</td><td>c+k</td><td>d+k</td></tr></table> <p>Coherent evolutions (additive)</p>	a	b	c	d	a+i	b+i	c+i	d+i	a+j	b+j	c+j	d+j	a+k	b+k	c+k	d+k	<table><tr><td>a</td><td>b</td><td>c</td><td>d</td></tr><tr><td>a x i</td><td>b x i</td><td>c x i</td><td>d x i</td></tr><tr><td>a x j</td><td>b x j</td><td>c x j</td><td>d x j</td></tr><tr><td>a x k</td><td>b x k</td><td>c x k</td><td>d x k</td></tr></table> <p>Coherent evolutions (multiplicative)</p>	a	b	c	d	a x i	b x i	c x i	d x i	a x j	b x j	c x j	d x j	a x k	b x k	c x k	d x k																	
a	b	c	d																																															
a+i	b+i	c+i	d+i																																															
a+j	b+j	c+j	d+j																																															
a+k	b+k	c+k	d+k																																															
a	b	c	d																																															
a x i	b x i	c x i	d x i																																															
a x j	b x j	c x j	d x j																																															
a x k	b x k	c x k	d x k																																															

Εικόνα 3.1. - Τύποι Biclusters

Στους τρεις πρώτους τύπους σημαντικό ρόλο έχουν οι αριθμητικές τιμές του πίνακα και στοχεύουν στον εντοπισμό των υποομάδων γραμμών ή στηλών με παρόμοια συμπεριφορά όπως φαίνεται στις παρακάτω εικόνες [Εικόνες 3.1, 3.2, 3.3]. Στον τέταρτο τύπο η εύρεση ομάδων επιτυγχάνεται όχι βάσει των αριθμητικών τιμών του πίνακα αλλά βάσει των συμβόλων. Τα σύμβολα αυτά μπορεί να είναι κάποιοι συγκεκριμένοι χαρακτήρες, να ανταποκρίνονται σε μια συγκεκριμένη σειρά, είτε να αναπαριστούν αλλαγές (θετικές και αρνητικές) που συνδέονται με μία τιμή.[10]

Στη συνέχεια παρουσιάζεται μια σύντομη σημειογραφία χρήσιμη για την αναλυτικότερη περιγραφή των 4 αυτών τύπων. Δεδομένου ενός πίνακα $A=(X,Y)$ όπου X το σύνολο των γραμμών και Y το σύνολο των στηλών του, ένα Bicluster είναι ένας υποπίνακας (I,J) όπου I ένα υποσύνολο του X και J ένα υποσύνολο του Y αντίστοιχα και a_{ij} η τιμή κάθε στοιχείου του πίνακα. Ακόμη ορίζουμε τις μέσες τιμές των γραμμών και των στηλών και τη μέση τιμή όλου του πίνακα (I,J) :

$$\alpha_{iJ} = 1/|J| \sum_{j \in J} a_{ij}$$

$$\alpha_{iJ} = 1/|I| \sum_{i \in I} a_{ij} \text{ και}$$

$$\alpha_{IJ} = 1/|I||J| \sum_{i \in I, j \in J} a_{ij} = 1/|I| \sum_{i \in I} \alpha_{iJ} = 1/|J| \sum_{j \in J} \alpha_{IJ}$$

Biclusters με σταθερές τιμές

Οι αλγόριθμοι από τους οποίους προκύπτουν Biclusters με σταθερές τιμές τις πιο πολλές φορές αναδιατάσσουν τις γραμμές και τις στήλες ενός πίνακα με σκοπό να ομαδοποιήσουν γραμμές και στήλες με παραπλήσιες τιμές. Ένα τέλειο σταθερό Bicluster είναι ένας υποπίνακας (I,J) όπου όλες οι τιμές του είναι ίσες και για κάθε $i \in I$ και $j \in J$: $a_{ij}=\mu$. Τέτοια παραδείγματα συναντάμε σε κάποιους πίνακες, όμως τις περισσότερες φορές τα δεδομένα προς επεξεργασία εμπεριέχουν θόρυβο, πράγμα που σημαίνει ότι ένα σταθερό Bicluster παρουσιάζεται με τιμές $n_{ij}+\mu$, όπου το n_{ij} είναι ο σχετικός θόρυβος της τιμής μ του a_{ij} . [10]

1.0	1.0	1.0	1.0
1.0	1.0	1.0	1.0
1.0	1.0	1.0	1.0
1.0	1.0	1.0	1.0

Εικόνα 3.2. - Bicluster με σταθερές τιμές

Biclusters με σταθερές τιμές σε γραμμές ή στήλες

Οι τιμές ενός βέλτιστου Bicluster με σταθερές γραμμές σε έναν υποπίνακα (I,J) είναι της μορφής $a_{ij}=\mu + \alpha_i$ ή $a_{ij}=\mu \times \alpha_i$ όπου μ είναι μια σταθερή τιμή εντός του Bicluster και α_i αποτελεί μία προσαρμοσμένη τιμή για κάθε γραμμή $i \in I$, η οποία προκύπτει είτε από πρόσθεση ή πολλαπλασιασμό. Ανάλογα, οι τιμές ενός βέλτιστου Bicluster με σταθερές στήλες σε έναν υποπίνακα (I,J) είναι της μορφής $a_{ij}=\mu + \beta_j$ ή $a_{ij}=\mu \times \beta_j$, όπου μ είναι μια σταθερή τιμή εντός του Bicluster και β_j αποτελεί μία προσαρμοσμένη τιμή για κάθε στήλη $j \in J$, η οποία προκύπτει είτε από πρόσθεση ή πολλαπλασιασμό. [10]

1.0	1.0	1.0	1.0
2.0	2.0	2.0	2.0
3.0	3.0	3.0	3.0
4.0	4.0	4.0	4.0

1.0	2.0	3.0	4.0
1.0	2.0	3.0	4.0
1.0	2.0	3.0	4.0
1.0	2.0	3.0	4.0

Εικόνα 3.3. - Bicluster με σταθερές τιμές σε γραμμές και στήλες

Στο αριστερό σκέλος της Εικόνας 3.3 το Bicluster εμφανίζει σταθερές τιμές στις γραμμές του ακολουθώντας το μοντέλο $a_{ij}=\mu + \alpha_i$, όπου $\mu=0.0$ και κάθε φορά προστίθεται σε κάθε νέα γραμμή η τιμή

$\alpha_i=1.0$. Στο δεύτερο σκέλος της Εικόνας 3.3 το Biclustερ εμφανίζει σταθερές τιμές στις στήλες ακολουθώντας το μοντέλο $\alpha_{ij}=\mu + \beta_j$, όπου $\mu=0.0$ και κάθε φορά προστίθεται σε κάθε νέα γραμμή η τιμή $\beta_j=1.0$.

Biclusters με συναφείς (coherent) τιμές

Η κατηγορία των Biclustερ με συναφείς τιμές δεν ορίζεται με κάποιο απλό μοντέλο πρόσθεσης ή πολλαπλασιασμού που εφαρμόζεται στα σταθερά Biclustερ. Υπάρχουν πιο σύνθετες τεχνικές που σχετίζονται με την ανάλυση της διακύμανσης μεταξύ των υποομάδων του πίνακα. Αυτές οι μέθοδοι χρησιμοποιούν μια συγκεκριμένη μορφή συνδιακύμανσης μεταξύ γραμμών και στηλών σε ένα Biclustερ για την αξιολόγηση της αποτελεσματικότητας της ομαδοποίησης.

Χρησιμοποιώντας λοιπόν ένα πιο σύνθετο μοντέλο πρόσθεσης, ένα τέλειο συνεκτικό Biclustερ χαρακτηρίζεται από μία υποομάδα γραμμών και στηλών που οι τιμές της α_{ij} επιλέγονται με βάση την ακόλουθη έκφραση: $\alpha_{ij}=\mu + \alpha_i + \beta_j$, όπου μ είναι μια σταθερή τιμή εντός του Biclustερ και α_i μία προσαρμοσμένη τιμή για κάθε γραμμή $i \in I$, και β_j μία προσαρμοσμένη τιμή για κάθε στήλη $j \in J$, οι οποίες προστίθενται στην μ . Ένα άλλο μοντέλο αναφέρει ότι τα Biclustερ με συναφείς τιμές μπορούν να περιγραφούν από ένα μοντέλο πολλαπλασιασμού, όπου $\alpha_{ij}=\mu' \times \alpha_i' \times \beta_j'$. [10]

1.0	2.0	5.0	0.0
2.0	3.0	6.0	1.0
4.0	5.0	8.0	3.0
5.0	6.0	9.0	4.0

1.0	2.0	0.5	1.5
2.0	4.0	1.0	3.0
4.0	8.0	2.0	6.0
3.0	6.0	1.5	4.5

Εικόνα 3.4. - Biclustερ με συναφείς τιμές(αριστερά μοντέλο πρόσθεσης, δεξιά μοντέλο πολλαπλασιασμού)

Για να γίνει πιο κατανοητή η μορφή των Biclustερ με συναφείς τιμές παρουσιάζεται το παράδειγμα του Πίνακα 3.1 για το μοντέλο της πρόσθεσης: Γνωρίζουμε ότι $\alpha_{ij}=\mu + \alpha_i + \beta_j$, άρα για $\mu=0$

Πίνακας 3.1. - Μοντέλο Πρόσθεσης

	$\beta_1=0$	$\beta_2=1$	$\beta_3=4$	$\beta_4=-1$
$\alpha_1=1$	1	2	5	0
$\alpha_2=2$	2	3	6	1
$\alpha_3=4$	4	5	8	3
$\alpha_4=5$	5	6	9	4

Και για το μοντέλο του πολλαπλασιασμού βάσει ότι $\alpha_{ij}=\mu' \times \alpha_i' \times \beta_j'$, άρα για $\mu=1$

Πίνακας 3.2. - Μοντέλο Πολλαπλασιασμού

	$\beta 1=1$	$\beta 2=2$	$\beta 3=0,5$	$\beta 4=3/2$
$\alpha 1=1$	1	2	0,5	1,5
$\alpha 2=2$	2	4	3	3
$\alpha 3=4$	4	8	2	6
$\alpha 4=3$	3	6	1,5	4,5

Biclusters με συναφείς εξελίξεις (coherent evolutions)

Αυτός ο τύπος Bicluster έχει ως στόχο την ομαδοποίηση γραμμών ή στηλών ανεξάρτητα από τις ακριβείς τιμές των στοιχείων του πίνακα.[10] Τα Biclusters με συναφείς εξελίξεις, σε αντίθεση με τα Biclusters με συναφείς τιμές, προσδιορίζουν υποσύνολα των γραμμών (γονιδίων) που οι τιμές τους είτε αυξάνονται είτε μειώνονται σταδιακά κατά μήκος μιας ομάδας στηλών ή γραμμών.

Σε αντίθεση με τους παραπάνω τύπους, αυτός ο τύπος είναι πολύπλοκο να μοντελοποιηθεί με τη χρήση μαθηματικής εξίσωσης.

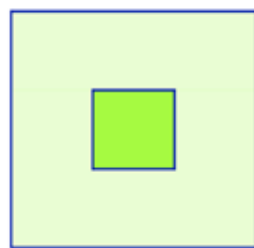
S1	S1	S1	S1	S1	S1	S1	S1
S1	S1	S1	S1	S2	S2	S2	S2
S1	S1	S1	S1	S3	S3	S3	S3
S1	S1	S1	S1	S4	S4	S4	S4
S1	S2	S3	S4	70	13	19	10
S1	S2	S3	S4	49	40	49	35
S1	S2	S3	S4	40	20	27	15
S1	S2	S3	S4	90	15	20	12

Εικόνα 3.5. - Biclusters με συναφείς εξελίξεις

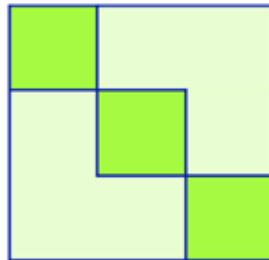
3.2. Δομή Διπλής Κατηγοριοποίησης

Μετά την ολοκλήρωση εκτέλεσης των αλγορίθμων Bicluster, οι δομές οι οποίες προκύπτουν ταξινομούνται σε 2 κατηγορίες: ένα μοναδικό Bicluster ή ένας μεγαλύτερος αριθμός από Biclusters με διάφορες μορφές, οι οποίες είναι οι παρακάτω:

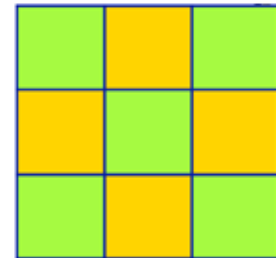
- Μονό.
- Αποκλειστικής γραμμής και στήλης.
- Μη επικαλυπτόμενο με τη μορφή σκακιάρας.
- Αποκλειστικά με σειρές.
- Αποκλειστικά με στήλες.
- Μη-επικαλυπτόμενα με δένδροειδή μορφή.
- Μη-επικαλυπτόμενα, μη αποκλειστικά.
- Επικαλυπτόμενα με ιεραρχική δομή.
- Αυθαίρετα τοποθετημένα-επικαλυπτόμενα.



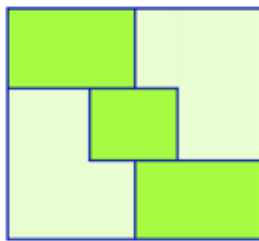
α



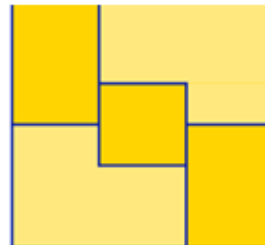
β



γ



δ



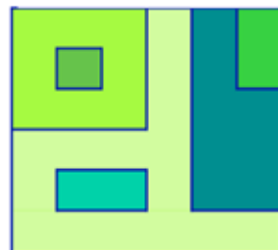
ε



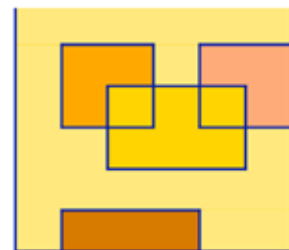
στ



ζ



η



θ

Εικόνα 3.6. - Δομές Bicluster: α. Μονό, β. Αποκλειστικής γραμμής και στήλης, γ. Μη επικαλυπτόμενο με τη μορφή σκακιάρας, δ. Αποκλειστικά με σειρές, ε. Αποκλειστικά με στήλες, στ. Μη-επικαλυπτόμενα με δένδροειδή μορφή, ζ. Μη-επικαλυπτόμενα, μη αποκλειστικά, η. Επικαλυπτόμενα με ιεραρχική δομή, θ. Αυθαίρετα τοποθετημένα-επικαλυπτόμενα [55]

Μία βοηθητική απεικόνιση για την αναζήτηση πολλών Biclusters σε έναν πίνακα είναι η χρήση του ίδιου χρώματος για κάθε στοιχείο με την ίδια τιμή a_{ij} . Με αυτόν τον τρόπο είναι πιο απλή η εναλλαγή γραμμών και στηλών και πιο εύκολη η ομαδοποίηση γραμμών και στηλών που σχηματίζουν «μπλοκς» με παρόμοια χρώματα. Δηλαδή τα Biclusters είναι τα «μπλοκς», δηλ. οι υποπίνακες του αρχικού πίνακα με παραπλήσιες τιμές έκφρασης. Στην ιδεατή περίπτωση η εικόνα που δημιουργείται αποτελείται από K τετράγωνα «μπλοκς» στη διαγώνιο του πίνακα (3.6.β). Αυτή η ιδέα στηρίζεται στο ότι κάθε γραμμή και στήλη του πίνακα μπορεί να ανήκει αποκλειστικά σε ένα από τα K Biclusters. Τις πιο πολλές φορές όμως πολλές γραμμές και στήλες του πίνακα μπορεί να ανήκουν σε περισσότερα από ένα Biclusters και είτε υπάρχει επικάλυψη μεταξύ των στοιχείων (3.6.η,θ), είτε όχι (3.6.γ,δ,ε,στ,ζ).[55]

3.3. Ποιοτικά Μέτρα Διπλής Κατηγοριοποίησης

Αρκετά ποιοτικά μέτρα για τα Biclusters έχουν προταθεί παρ'όλα αυτά, κανένα από τα προτεινόμενα ποιοτικά μέτρα δεν είναι σε θέση να αναγνωρίσει ένα τέλειο πρότυπο μετατόπισης και κλιμάκωσης σε ένα Bicluster.

Στην ενότητα αυτή θα **παρουσιαστούν** μερικά από τα πιο γνωστά μέτρα αξιολόγησης για Biclusters, όπως περιγράφονται στη συγκριτική μελέτη των Pontes και συνεργατών (2015).[56]

Διακύμανση (VAR)

Ο Hartigan χρησιμοποίησε τη διακύμανση Bicluster, όπως δίνεται στην εξίσωση 3.3.1 ως συνεχές μέτρο, όπου ο στόχος του αλγόριθμού του ήταν να ελαχιστοποιήσει το άθροισμα των διακυμάνσεων του Bicluster. Ο ορισμός της Διακύμανσης είναι ο εξής:

$$\text{VAR} (A) = \frac{1}{|I|*|J|} \sum_{i=1}^{|I|} \sum_{j=1}^{|J|} (a_{ij} - a_{IJ})^2 \quad (3.3.1)$$

Μπορεί να σημειωθεί ότι η διακύμανση ανιχνεύει μόνο σταθερά Biclusters. Ως εκ τούτου, οι προσεγγίσεις biclustering με βάση την διακύμανση των μεθόδων ελαχιστοποίησης περιλαμβάνουν επίσης και άλλα κριτήρια ομοιογένειας για να ανιχνεύουν άλλους τύπους Biclusters.[56]

Μέσο τετράγωνο υπολείμματος (Mean Square Residue, MSR)

Ο Cheng και Church ήταν οι πρώτοι που εφάρμοσαν biclustering σε δεδομένα γονιδιακής έκφρασης. Εισήγαγαν έναν από τους πιο δημοφιλείς αλγόριθμους biclustering που συνδυάζει μια άπληστη αναζήτηση για την εύρεση Biclusters με ένα μέτρο για την αξιολόγηση της ποιότητας των εν λόγω Biclusters. Προκειμένου να αξιολογηθεί η ποιότητα των Biclusters, όπως ήδη έχει παρουσιαστεί, ο αλγόριθμος χρησιμοποιεί το μέσο τετράγωνο υπολείμματος (MSR), το οποίο αποτελεί τη διακύμανση της ομάδας όλων των στοιχείων του εκάστοτε Bicluster. Το μέτρο αυτό αποσκοπεί στο να αξιολογήσει τη συνοχή των γονιδίων

και τις συνθήκες ενός Bicluster A που αποτελείται από σειρές I, και J στήλες. Ο ορισμός του MSR είναι ο εξής:

$$MSR(A) = \frac{1}{|I|*|J|} \sum_{i=1}^{|I|} \sum_{j=1}^{|J|} (a_{ij} - a_{iJ} - a_{IJ} + a_{IJ})^2 \quad (3.3.2)$$

Στόχος είναι η εύρεση μεγάλων Biclusters με χαμηλό αποτέλεσμα MSR και συγκεκριμένα μικρότερο από ένα συγκεκριμένο **κατώφλι** (threshold) δ - bicluster. Εάν ένα Bicluster έχει MSR ίσο με το μηδέν, αυτό σημαίνει ότι τα γονίδια του κυμαίνονται κατά τον ίδιο ακριβώς τρόπο κάτω από το υποσύνολο των πειραματικών συνθηκών, και έτσι μπορεί να θεωρηθεί ένα τέλειο Bicluster.[56]

Κλιμάκωση μέσου τετραγώνου υπολείμματος (Scaling Mean Squared Residue, SMSR)

Οι Mukhopadhyay, Maulik και Bandyopadhyay ανέπτυξαν ένα μέτρο αξιολόγησης για Biclusters που είναι σε θέση να αναγνωρίζει την κλιμάκωση προτύπων. Στο έργο τους, αναλύουν τους λόγους για τους οποίους το MSR είναι σε θέση να αναγνωρίσει τη μετατόπιση προτύπων (μοτίβων) σε Biclusters αλλά **όχι** την κλιμάκωση προτύπων. Χρησιμοποιώντας τον μαθηματικό τύπο για την κλιμάκωση προτύπων, ορίζουν ένα μέτρο που στη συνέχεια αποδείχθηκε ότι μπορεί να προσδιορίσει την κλιμάκωση προτύπων. Αυτό το νέο μέτρο καλείται SMSR, από το Scalling MSR. Παρ'όλα αυτά, το SMSR **δεν** είναι ικανό να ανιχνεύσει τη μετατόπιση προτύπων.[56] **Ο ορισμός του SMSR είναι ο εξής:**

$$SMSR(A) = \frac{1}{|I|*|J|} \sum_{i=1}^{|I|} \sum_{j=1}^{|J|} \frac{(a_{ij}*a_{IJ} - a_{iJ}*a_{IJ})^2}{a_{ij}^2*a_{IJ}^2} \quad (3.3.3)$$

Δείκτης Σχετικότητας ή Συνάφειας (Relevance Index, RI)

Οι Yip, Ng και Cheung πρότειναν ένα μέτρο αξιολόγησης που διαφέρει ελαφρώς από το MSR, στο οποίο η ποιότητα ενός Bicluster υπολογίζεται ως το άθροισμα των δεικτών σχετικότητας των στηλών. **Ο δείκτης σχετικότητας R_{Ij} για στήλη $j \in J$ ορίζεται ως :**

$$R_{Ij} = 1 - \frac{\sigma_{Ij}^2}{\sigma_j^2} \quad (3.3.4)$$

Όπου το σ_{Ij}^2 (τοπική διακύμανση) είναι η διακύμανση των τιμών στη στήλη j για το Bicluster και σ_j^2 (καθολική διακύμανση) είναι το σύνολο των δεδομένων. Εφόσον ο δείκτης δίνει μια υψηλή τιμή όταν η τοπική διακύμανση είναι μικρή σε σχέση με την ολική διακύμανση, ο σχετικός δείκτης, για μια στήλη, μεγιστοποιείται αν η τοπική διακύμανση του είναι μηδέν, υπό την προϋπόθεση ότι η συνολική διακύμανση δεν είναι. Βασιζόμενοι στο δείκτη σχετικότητας, η ποιότητα ενός Bicluster μετριέται ως το άθροισμα των τιμών του δείκτη από όλες τις επιλεγμένες καταστάσεις.[56]

Συσχέτιση βασισμένη σε μέτρα (Correlation-based Measures)

Οι συντελεστές συσχέτισης έχουν χρησιμοποιηθεί εκτενώς σε διαφορετικά είδη αναλύσεων μικροσυστοιχιών, όπως το clustering. Αυτές οι στατιστικές μετρούν την συνολική ομοιότητα των γονιδίων

χωρίς τοποθέτηση, με οποιαδήποτε έμφαση σε συγκεκριμένα μεγέθη, λαμβάνοντας υπόψη τόσο τις αρνητικές συσχετίσεις, όσο και τις θετικές. Παρακάτω εξετάζονται οι προσεγγίσεις βάσει συσχέτισης (correlation-based) που προτείνονται για την αξιολόγηση ενός Biclusters.[56]

Συντελεστής συσχέτισης Pearson (Pearson's Correlation Coefficient, PCC)

Ο PCC μεταξύ δύο μεταβλητών ορίζεται ως, η συνδιακύμανση των δύο μεταβλητών διαιρούμενη από το προϊόν των τυπικών αποκλίσεων τους. Πρόκειται για ένα μέτρο της γραμμικής εξάρτησης μεταξύ δύο μεταβλητών, και δίνει μια τιμή μεταξύ +1 και -1, συμπεριλαμβανομένων και των δύο. Η τιμή 1 υποδηλώνει ότι μία γραμμική εξίσωση περιγράφει τη σχέση μεταξύ των δύο μεταβλητών τέλεια (θετική συσχέτιση), ενώ η τιμή -1 υποδηλώνει ότι όλα τα σημεία δεδομένων βρίσκονται σε μια γραμμή, για την οποία μια μεταβλητή μειώνεται καθώς η άλλη αυξάνεται (αρνητική συσχέτιση). Η τιμή 0 σημαίνει ότι δεν υπάρχει γραμμική συσχέτιση μεταξύ των μεταβλητών.

Ο PCC είναι ένα πολύ αποτελεσματικό μέτρο, που ποσοτικοποιεί την από κοινού ρύθμιση μεταξύ των ζευγαριών των γονιδίων και επιτρέπει τη σύλληψη της **μετατόπισης** όσο και της **κλιμάκωσης** προτύπων, τα οποία θα προσδιοριστούν χωριστά από πρόσθετα και πολλαπλασιαστικά μοντέλα, αντίστοιχα. Παρ ' όλα αυτά, ο PCC **δεν είναι αποτελεσματικός** για την αναγνώριση

σταθερών Biclusters ή σταθερών μοτίβων σειράς, δεδομένου ότι αυτά τα είδη των προτύπων θα κάνουν τον παρονομαστή μηδέν.

Ο PCC μεταξύ δύο σειρών (γονίδια) $i_1, i_2 \in I$ σε σχέση με τις στήλες (καταστάσεις) $j \in J$ ορίζεται ως:

$$PCC(i_1, i_2) = \frac{\sum_{j=1}^{|J|} (a_{i_1 j} - a_{i_1 J})(a_{i_2 j} - a_{i_2 J})}{\sqrt{\sum_{j=1}^{|J|} (a_{i_1 j} - a_{i_1 J})^2 \sum_{j=1}^{|J|} (a_{i_2 j} - a_{i_2 J})^2}} \quad (3.3.5)$$

Όπου $a_{i_1 j}$ και $a_{i_2 j}$ συμβολίζουν τα στοιχεία στις γραμμές i_1, i_2 και τη στήλη j και $a_{i_1 J}, a_{i_2 J}$ αντιπροσωπεύουν τα μέσα των γραμμών i_1 και i_2 αντίστοιχα.

Ο PCC ποσοτικοποιεί τη συνοχή μεταξύ των ζευγαριών των γονιδίων. Ως εκ τούτου, προκειμένου να μετρηθεί η συνοχή του Biclusters, πρέπει κανείς να υπολογίσει όλα τα ζεύγη τιμών PCC μεταξύ των σειρών στο ίδιο Biclusters. Επιπροσθέτως, ο PCC έχει νόημα μόνο για τη μέτρηση της συνοχής μεταξύ των γραμμών, διότι είναι υπερβολικά περιοριστικός αν χρησιμοποιηθεί για τη μέτρηση της συνοχής μεταξύ των στηλών ταυτόχρονα. Ο PCC έχει χρησιμοποιηθεί ως έχει (μεταξύ των γραμμών) σε αρκετές έρευνες, ενώ σε άλλες μελέτες οι συγγραφείς ορίζουν ένα μέτρο που βασίζεται στον PCC (PCC-based), το οποίο ονομάζεται **Μέση Συσχέτιση (Average Correlation(-AC), με τη μορφή της εξίσωσης (3.3.6) :**

$$AC(A) = \frac{\sum_{i_1=1}^{|I|-1} \sum_{i_2=i_1+1}^{|I|} PCC(i_1, i_2)}{\binom{|I|}{2}} \quad (3.3.6)$$

Για να χρησιμοποιηθεί ο PCC για τη μέτρηση της συνοχής μεταξύ των στηλών, καθώς και για να ξεπεραστεί το θέμα της περιοριστικότητας, τόσο οι Yang, Dai ,και Yan στο έργο τους, όσο και οι Teng και Chan έχουν ορίσει ένα μέτρο που βασίζεται στον PCC (PCC-based), το οποίο είναι μέτρο αξιολόγησης του Biclustor όσον αφορά τη βελτίωση συσχέτισης, είτε στις γραμμές είτε στις στήλες.[56]

Βαθμός συσχέτισης υποπίνακα (Sub-Matrix Correlation Score, SCS)

Οι Yang, Dai και Yan χρησιμοποιούν το Pearson correlation score ως τη βάση για να καθορίσουν το μέτρο τους, υποθέτοντας ότι ένα τέλειο συσχετιζόμενο πρότυπο ικανοποιεί τέλεια γραμμική συσχέτιση σε διανύσματα γραμμής και στήλης. Σχέσεις, συσχετίσεις για γραμμές και στήλες παρουσιάζονται παρακάτω:

$$S_{row} = \min_{i_1 \in I} (S_{i_1 J}), S_{i_1 J} = 1 - \frac{\sum_{i_2 \neq i_1, i_2 \in I} |cor(x_{i_1 J}, x_{i_2 J})|}{|I| - 1} \quad (3.3.7)$$

$$S_{col} = \min_{j_1 \in J} (S_{I j_1}), S_{I j_1} = 1 - \frac{\sum_{j_2 \neq j_1, j_2 \in J} |cor(x_{I j_1}, x_{I j_2})|}{|J| - 1} \quad (3.3.8)$$

Όπου $cor(x_{i_1 J}, x_{i_2 J})$ και $cor(x_{I j_1}, x_{I j_2})$ αντιπροσωπεύουν τον PCC του κάθε ζεύγους γονιδίων ή καταστάσεων στο Biclustor, αντίστοιχα.

Το S_{row} score αντικατοπτρίζει το βαθμό συσχέτισης σχετικά με τις σειρές του Biclustor, ενώ το S_{col} ανακλά το βαθμό συσχέτισης σχετικά με τις στήλες, όπου και οι δύο ορισμοί είναι ασύμμετροι.[56]

Το σκορ συσχέτισης (correlation score) του υποπίνακα ορίζεται ως η ελάχιστη τιμή των δύο αυτών scores:

$$S(A) = \min(S_{row}(I, J), S_{col}(I, J)) \quad (3.3.9)$$

Μέση τιμή συσχέτισης Average Correlation Value (ACV)

Η ACV προτάθηκε από τους Teng και Chan για την αξιολόγηση της ομοιογένειας ενός Biclustor ή ενός πίνακα δεδομένων με τον ακόλουθο τρόπο:

$$ACV = \max \left\{ \frac{\sum_{i_1=1}^{|I|} \sum_{i_2=1}^{|I|} |r_{row_{i_1 i_2}}| - |I|}{|I|^2 - |I|}, \frac{\sum_{j_1=1}^{|J|} \sum_{j_2=1}^{|J|} |r_{col_{j_1 j_2}}| - |J|}{|J|^2 - |J|} \right\} \quad (3.3.10)$$

Όπου $r_{row_{i_1 i_2}}$, $r_{col_{j_1 j_2}}$ αναφέρονται στη συσχέτιση μεταξύ κάθε ζεύγους των σειρών i_1, i_2 ή στηλών j_1, j_2 σύμφωνα με τον συντελεστή Pearson.

Η ACV (A) έχει αποδειχθεί ότι είναι στο διάστημα $[0,1]$, όπου μία τιμή ίση με 1 σημαίνει ότι οι γραμμές ή οι στήλες του Biclustor έχουν σε μεγάλο βαθμό συνεκφραστεί, ενώ μια χαμηλή ACV σημαίνει ότι ούτε οι καταστάσεις (συνθήκες) ούτε τα γονίδια είναι παρόμοια. Επομένως, οι υψηλότερες τιμές της ACV

προτιμώνται. Οι συγγραφείς, στην εργασία τους, απέδειξαν ότι η ACV δίνει πάντα την επιθυμητή τιμή τόσο για το πρόσθετο όσο και για το πολλαπλασιαστικό μοντέλο, σε αντίθεση με το MSR.

Αν και η ACV παρουσιάζεται ως το κριτήριο για την αξιολόγηση Biclusters, δεν έχει χρησιμοποιηθεί για να καθοδηγήσει την αναζήτηση στον αλγόριθμό τους. Αντ' αυτού, ο αλγόριθμος βασίζεται στη χρήση ενός σταθμισμένου συντελεστή συσχέτισης.[56]

Μέσος Όρος Spearman's Rho (Average Spearman's Rho, ASR)

Ο ASR προτάθηκε για πρώτη φορά από τους Ayadi, Elloumi και Hao και βασίζεται στη χρήση του βαθμού συσχέτισης Spearman, η οποία μετρά την στατιστική εξάρτηση μεταξύ δύο μεταβλητών, αξιολογώντας το πόσο καλά η σχέση τους μπορεί να περιγραφεί χρησιμοποιώντας μια μονότονη συνάρτηση. Η τιμή του κυμαίνεται μεταξύ -1 και +1.

Η πιο σημαντική διαφορά μεταξύ των συσχετίσεων του Spearman και Pearson είναι ότι ο Pearson αξιολογεί γραμμικές σχέσεις, ενώ ο Spearman αξιολογεί μονοτονικές σχέσεις. Ως εκ τούτου, η συσχέτιση Spearman είναι λιγότερο ευαίσθητη σε ακραίες τιμές, που βρίσκονται στις ουρές και των δύο δειγμάτων.

Παρ' όλα αυτά, όταν τα δεδομένα διανέμονται περίπου ελλειπτικά και δεν υπάρχουν εξαιρετικά ακραίες τιμές, οι συσχετίσεις Spearman και Pearson δίνουν παρόμοιες τιμές.[56]

Ο ASR υπολογίζεται όπως στην παρακάτω εξίσωση και εξάγει μια τιμή στο διάστημα [-1,1], όπου η υψηλή /χαμηλή τιμή κοντά στο 1 / -1 δείχνει ότι τα γονίδια ή οι συνθήκες του Bicluster συσχετίζονται ισχυρώς, είτε θετικά είτε αρνητικά, αντίστοιχα.

$$ASR = 2 * \max \left\{ \frac{\sum_{i \in I} \sum_{j \geq i+1, j \in I} \rho_{ij}}{|I|(|I|-1)}, \frac{\sum_{k \in J} \sum_{l \geq k+1, l \in J} \rho_{kl}}{|J|(|J|-1)} \right\} \quad (3.3.11)$$

Όπου ρ_{ij} , ρ_{kl} αναφέρονται στη συσχέτιση Spearman μεταξύ δύο γονιδίων ή καταστάσεων, αντίστοιχα.[56]

Μέτρο διπλής κατηγοριοποίησης Spearman (Spearman's Biclustering Measure, SBM)

Το SBM έχει προταθεί πρόσφατα από τους Flores, Inza, Larranaga και Calvo, τονίζοντας την ικανότητά του να αναγνωρίζει πολύπλοκη συνοχή προτύπων στα Biclusters, όπως η μετατόπιση και η κλιμάκωση, όπως και, αρνητικές συσχετίσεις.

Ο υπολογισμός του συντελεστή αυτού επιτυγχάνεται σε δύο στάδια. Πρώτον, τα δεδομένα μετατρέπονται σε τάξεις, προκειμένου να υπολογιστεί ο συντελεστής Spearman, που συμβολίζεται ως $r_{x,y}$ για κάθε ζεύγος γονιδίων ή καταστάσεων. Χρησιμοποιώντας αυτούς τους συντελεστές, το SBM ορίζεται ως εξής:

$$SBM(A_{ij}) = \alpha(A_{ij}) * r_{A_{ij}}^{\bar{G}} * \beta(A_{ij}) * r_{A_{IJ}}^{\bar{G}} \quad (3.3.12)$$

Όπου $r_{A_{ij}}^{\bar{G}}$ και $r_{A_{IJ}}^{\bar{G}}$ υποδηλώνουν την συνοψιζόμενη έκφραση των τάσεων που παρατηρούνται στα γονίδια και στις καταστάσεις του Bicluster, αντίστοιχα. Αυτά υπολογίζονται όπως στις εξισώσεις που ακολουθούν:

$$r_{A_{IJ}}^{\bar{G}} = \frac{2}{|I| * (|I| - 1)} * \sum_{i=1}^{|I|} \sum_{i'=i+1}^{|I|} |r_{|ii|'}^G|, \quad (3.3.13)$$

$$r_{A_{IJ}}^{\bar{C}} = \frac{2}{|J| * (|J| - 1)} * \sum_{j=1}^{|J|} \sum_{j'=j+1}^{|J|} |r_{|jj|'}^C|$$

Όπου $|r_{|ii|'}^G|$ και του $|r_{|jj|'}^C|$ αντιστοιχούν στις απόλυτες τιμές του συντελεστή συσχέτισης Spearman μεταξύ ενός ζεύγους γονιδίων i, i' και ενός ζεύγους καταστάσεων j, j' , αντίστοιχα.

Οι όροι α (A_{ij}) και β (A_{ij}) στην εξίσωση (3.3.12) αντιπροσωπεύουν τους συντελεστές αξιοπιστίας, και το βάρος της επιρροής των τάσεων που παρατηρήθηκαν και στα δύο γονίδια και στις καταστάσεις, αντίστοιχα.[56]

Μέτρα βασισμένα σε τυποποίηση (Standardisation-based Measures)

Μία σημαντική παρατήρηση που μπορεί να εξαχθεί από την ανάλυση των Biclusters είναι ότι το εύρος των τιμών έκφρασης που λαμβάνονται από τα γονίδια μπορεί να ποικίλλει σημαντικά ανάλογα με τη μικροσυστοιχία που λαμβάνεται ως είσοδος. Συνεπώς, προκειμένου να γίνει σωστή σύγκριση μεταξύ κάθε γονιδίου και προτύπου, μια προηγούμενη διαδικασία τυποποίησης του Biclusters θα επιτρέψει στα επίπεδα έκφρασης να κλιμακωθούν σε ένα κοινό φάσμα. Ο μηχανισμός αυτός θα είναι επίσης υπεύθυνος για την ομαλοποίηση της συμπεριφοράς του γονιδίου, δεδομένου ότι η πιο σημαντική πτυχή είναι να χαρακτηρίσουμε το γονίδιο καλύτερα από την τάση του, παρά από τις αριθμητικές τιμές του.

Γονιδιακή τυποποίηση (**Gene Standardisation**) ενός A Biclusters αντιστοιχεί στο τυποποιημένο Biclusters \hat{A} , του οποίου το στοιχείο \hat{a}_{ij} , το εξάγουμε ως εξής:

$$\hat{a}_{ij} = \frac{\hat{a}_{ij} - \mu_{gi}}{\sigma_{gi}}, 1 \leq i \leq |I|, 1 \leq j \leq |J| \quad (3.3.14)$$

Όπου σ_{gi} είναι η τυπική απόκλιση όλων των τιμών έκφρασης του γονιδίου i και μ_{gi} είναι ο μέσος όρος της σειράς i στο A.

Μέσω της τυποποίησης, δύο ευδιάκριτα «έργα» εκτελούνται. Το πρώτο είναι να μετατοπιστεί το σύνολο των γονιδίων σε ένα παρόμοιο εύρος τιμών (κοντά στο 0 στην περίπτωση αυτή). Το δεύτερο είναι α) να ομογενοποιήσει τις τιμές έκφρασης για κάθε γονίδιο, τροποποιώντας με αυτόν τον τρόπο τις αξίες τους κάτω από όλες τις συνθήκες, και β) να εξομαλύνει τη γραφική αναπαράστασή τους, λόγω της διόρθωσης του καθολικού παράγοντα κλιμάκωσης στον παρονομαστή.[56]

Τρία διαφορετικά μέτρα έχουν οριστεί στη βιβλιογραφία, χρησιμοποιώντας τη διαδικασία της τυποποίησης για την αξιολόγηση Biclusters. Δύο από αυτά, MSA και VE, βασίζονται στη γονιδιακή τυποποίηση, ενώ το άλλο μέτρο βασίζεται σε μια κατάσταση τυποποίησης, όπως περιγράφεται παρακάτω.

Μέγιστη τυπική περιοχή (Maximal Standard Area, MSA)

Η ιδέα πίσω από τη MSA είναι να μετρηθεί η έκταση της περιοχής μεταξύ των μέγιστων και ελάχιστων τιμών των επιπέδων έκφρασης που λαμβάνουν τα γονίδια υπό τις καταστάσεις που περιέχονται στο Biclusters. Έτσι, αυτό που υπολογίζεται είναι η περιοχή που απεικονίζεται από τη μέγιστη διακύμανση των επιπέδων έκφρασης

για κάθε πειραματική κατάσταση. Για κάθε κατάσταση, λαμβάνονται οι ελάχιστες και οι μέγιστες τιμές του επιπέδου έκφρασης για όλα τα γονίδια που περιέχονται στο Biclust. Αυτά τα ζεύγη τιμών ορίζουν μια ζώνη σε όλες τις καταστάσεις στο Biclust, και η περιοχή της ζώνης αυτής είναι, ως εκ τούτου, το μέτρο MSA.

Τα ψηλότερα και χαμηλότερα όρια ενός Biclust για κάθε κατάσταση j ορίζονται ως $M_j(A)$ και $m_j(A)$ αντίστοιχα, με τον ακόλουθο τρόπο:

$$M_j(A) = \max_i a_{ij}, \forall i; m_j(A) = \min_i a_{ij}, \forall i \quad (3.3.15)$$

Χρησιμοποιώντας αυτά τα όρια, η MSA ορίζεται ως η οριοθετημένη περιοχή από τα όρια της κάθε κατάστασης στο τυποποιημένο Biclust:

$$MSA(A) = \sum_{j=1}^{|J|-1} \left| \frac{M_j(\hat{A}) - m_j(\hat{A}) + M_{j+1}(\hat{A}) - m_{j+1}(\hat{A})}{2} \right| \quad (3.3.16)$$

Αν τα γονίδια ενός Biclust έχουν μια απόλυτα συνεκτική τάση, τότε η MSA (A) είναι ίση με το μηδέν. Αντίθετα, η MSA θα πάρει υψηλότερες τιμές όταν τα γονίδια είναι λιγότερο συσχετισμένα μεταξύ τους, αυτό οφείλεται στο γεγονός ότι τα $M(\hat{A})$ και $m(\hat{A})$ είναι πιο απομακρυσμένα μεταξύ τους.[56]

Εικονικό σφάλμα (Virtual Error, VE)

Το Virtual Error ακολουθεί παρόμοιες παραδοχές σχετικά με τη διαδικασία τυποποίησης στο πλαίσιο της αξιολόγησης με τη MSA. Η βασική ιδέα πίσω από το VE είναι να μετρήσει πώς τα γονίδια ακολουθούν τη γενική τάση εντός του Biclust. Για να «πιάσει» αυτή τη γενική τάση των γονιδίων σε όλες τις καταστάσεις που εμπεριέχονται στο Biclust, μια νέα εικονική γραμμή υπολογίζεται από τα γονίδια του Biclust, που ονομάζεται εικονικό σχέδιο ή εικονικό γονίδιο ρ . Κάθε στοιχείο ρ_j του ρ υπολογίζεται ως μέση τιμή της j -στήλης στην παρακάτω εξίσωση:

$$\rho_j = \frac{1}{|I|} \sum_{i=1}^{|I|} a_{ij} \quad (3.3.17)$$

Μόλις το εικονικό γονίδιο ρ έχει υπολογιστεί, και για να εκτιμηθεί πόσο καλά ένα συγκεκριμένο γονίδιο g_i του Biclust, ακολουθεί τη γενική τάση, το VE υπολογίζει τις διαφορές μεταξύ των τιμών έκφρασης του επιπέδου g_i και τις τιμές του ρ , για κάθε πειραματική κατάσταση του Biclust. Αυτές οι διαφορές υπολογίζονται χρησιμοποιώντας τα τυποποιημένα γονίδια στο Biclust, καθώς και το τυποποιημένο εικονικό γονίδιο $\hat{\rho}_j$. Ο ορισμός του **VE(A)** είναι ο εξής:

$$VE(A) = \frac{1}{|I| * |J|} \sum_{i=1}^{|I|} \sum_{j=1}^{|J|} |a_{ij} - \hat{\rho}_j| \quad (3.3.18)$$

Το VE υπολογίζει τις διαφορές μεταξύ των πραγματικών γονιδίων και των εικονικών, τη στιγμή που έχουν τυποποιηθεί. Ως εκ τούτου, όσο πιο όμοια είναι τα γονίδια, τόσο χαμηλότερη είναι η τιμή για το VE. Πράγματι, το VE είναι μηδέν για εκείνα τα Biclusts που ακολουθούν είτε μία τέλεια μετατόπιση ή μία κλιμάκωση

προτύπου. Παρ' όλα αυτά, το VE δεν μπορεί να αποδειχθεί ότι αναγνωρίζει και τα δύο είδη των προτύπων ταυτόχρονα.[56]

Μετατόπιση εικονικού σφάλματος(Transposed Virtual Error ,VE^t)

Το **VE^t** υπολογίζεται παρομοίως με το VE αλλά, λαμβάνοντας υπόψη το μεταφερόμενο Bicluster. Η ιδέα εδώ είναι να δημιουργήσουμε το εικονικό μοτίβο στην “κατάσταση διάστασης” (condition dimension), η οποία έχει ονομαστεί εικονική κατάσταση. Στη συνέχεια, οι διαφορές μεταξύ των τυποποιημένων τιμών για κάθε κατάσταση και των τυποποιημένων εικονικών καταστάσεων υπολογίζονται με τον ίδιο τρόπο όπως στο VE.

Με αυτό τον τρόπο, η εικονική κατάσταση υπολογίζεται ως:

$$\rho_i = \frac{1}{|J|} \sum_{j=1}^{|J|} a_{ij} \quad (3.3.19)$$

και η τιμή Vet του bicluster A επιτυγχάνεται με τη χρήση της τυποποιημένης εικονικής κατάστασης $\hat{\rho}$ μαζί με το τυποποιημένο Bicluster όπως τα δείγματα:

$$VE^t(A) = \frac{1}{|I|*|J|} \sum_{i=1}^{|I|} \sum_{j=1}^{|J|} |\hat{a}_{ij} - \hat{\rho}_i| \quad (3.3.20)$$

Το **VE^t** έχει αποδειχθεί ότι είναι μηδέν για τα Biclusters με τέλεια μετατόπιση, με κλιμάκωση ή και σε συνδυασμό των δύο. Ως εκ τούτου, φαίνεται να αναγνωρίζει αποτελεσματικά τόσο τη μετατόπιση όσο και τη κλιμάκωση σε Biclusters είτε ταυτόχρονα είτε ανεξάρτητα.[56]

Βαθμολογίες ομοιότητας (Similarity Scores)

Ο Liu και ο Wang πρότειναν τη χρήση μιας βαθμολογίας ομοιότητας μεταξύ δύο γονιδίων και επίσης μια βαθμολογίας ομοιότητας για έναν υποπίνακα. Η πρώτη χρησιμοποιείται όταν το αναφερόμενο γονίδιο είναι γνωστό εκ των προτέρων. Όταν δεν είναι γνωστό, ένας αριθμός γονιδίων θα μπορούσε επίσης να επιλεγεί τυχαία όπως τα αναφερόμενα γονίδια.[56]

Βαθμός ομοιότητας μεταξύ γονιδίων (Similarity Score Between Genes)

Ο βαθμός ομοιότητας μεταξύ δύο γονιδίων (γονίδιο i και γονίδιο αναφοράς i^*) υπό τον όρο j υπολογίζεται σύμφωνα με την παρακάτω εξίσωση:

$$s_{ij} = \begin{cases} 0 & \text{if } d_{ij} > a * d_{avg} \\ 1 - \frac{d_{ij}}{a*d_{avg}} + \beta & \text{αλλιώς} \end{cases} \quad (3.3.21)$$

Όπου d_{avg} ορίζεται ως η μέση τιμή της απόστασης όλων των στοιχείων του πίνακα έκφρασης.

$$d_{avg} = \frac{\sum_{i \in I, j \in J} d_{ij}}{|I||J|} \quad (3.3.22)$$

και το d_{ij} είναι η απόλυτη τιμή της διαφοράς έκφρασης μεταξύ του γονιδίου I και του γονιδίου αναφοράς i^* για την κατάσταση j στον πίνακα έκφρασης α :

$$d_{ij} = |a_{ij} - a_{i*j}| \quad (3.3.23)$$

Το $a * d_{avg}$ χρησιμοποιείται ως το όριο για να αγνοήσει στοιχεία με μεγάλο d_{ij} , ώστε να βρεθούν σταθερά Biclusters, και το β είναι το κέρδος για τις μικρές d_{ij} . Με αυτό τον τρόπο, το β διευρύνει τον βαθμό ομοιότητας για τις μικρές d_{ij} και αγνοεί d_{ij} που είναι μεγαλύτερα από το όριο. Σύμφωνα με την εξίσωση s_{ij} s(3.3.21), η τιμή της s_{ij} θα είναι πάντα μεγαλύτερη από ή ίση με μηδέν, όπου είναι η χειρότερη αξία του βαθμού ομοιότητας.[56]

Βαθμός ομοιότητας για ένα Bicluster (Similarity Score For a Bicluster, SS)

Για κάθε γραμμή $i \in I$, ο βαθμός ομοιότητας της σειράς i στο A είναι:

$$s(i, J) = \sum_{j \in J} s_{ij} \quad (3.3.24)$$

Ομοίως για κάθε στήλη $j \in J$ ο βαθμός ομοιότητας στην A αντιστοιχεί σε :

$$s(I, j) = \sum_{i \in I} s_{ij} \quad (3.3.25)$$

Χρησιμοποιώντας αυτές τις εξισώσεις ο βαθμός ομοιότητας για το Bicluster υπολογίζεται ως η ελάχιστη τιμή του βαθμού ομοιότητας και των δύο γονιδίων και της κατάστασης των Biclusters:

$$s(A) = s(I, J) = \min \{ \min_{i \in I} s(i, J), \min_{j \in J} s(I, j) \} \quad (3.3.26)$$

Ο στόχος όταν ψάχνουμε για Biclusters είναι να βρεθούν υπο-πίνακες με υψηλότερες τιμές για τον βαθμό ομοιότητας. Προκειμένου να βελτιωθεί η ποιότητα της παραγωγής, οι Liu και Wang εισήγαγαν επίσης ως δεύτερο κριτήριο τον μέσο όρο βαθμού ομοιότητας. Αυτός αποτελείται από τον μέσο όρο όλων των όμοιων τιμών της εξίσωσης d_{ij} για όλα τα στοιχεία του Bicluster.

Παρόλο που ο τύπος του Bicluster που βρέθηκε χρησιμοποιώντας το βαθμό ομοιότητας, εξαρτάται από τις τιμές για τα διαφορετικά κατώτατα όρια (α , β , και γ για το μέσο όρο), αναγνωρίζονται μόνο τα σταθερά και τα πρόσθετα Biclusters.[56]

Στην εργασία αυτή βασιστήκαμε σε όλα τα παραπάνω, δηλαδή τα πιο γνωστά υφιστάμενα μέτρα αξιολόγησης για Biclusters, αλλά παράλληλα ακολουθήσαμε μια δική μας προσέγγιση, η οποία παρουσιάζεται στο επόμενο κεφάλαιο (Κεφάλαιο 4).

Δεδομένα και Προτεινόμενη Μεθοδολογία

4.1. Ανάπτυξη Μεθοδολογίας

Η μέθοδος διπλής κατηγοριοποίησης εφαρμόστηκε σε ένα σύνολο δεδομένων, το οποίο χορηγήθηκε ευγενικά στο Εργαστήριο Ψηφιακής Επεξεργασίας Σήματος και Εικόνας από την Δρ. Obermayr (Ιατρικό Πανεπιστήμιο Βιέννης, Τμήμα Γενικής Γυναικολογίας και Γυναικολογικής Ογκολογίας), που αφορούσε σε τέσσερις διαφορετικούς τύπους καρκινικών κυτταρικών σειρών: α) κυτταρικές σειρές καρκίνου του μαστού, β) κυτταρικές σειρές καρκίνου του τραχήλου της μήτρας, γ) κυτταρικές σειρές καρκίνου του ενδομητρίου, και δ) κυτταρικές σειρές καρκίνου των ωοθηκών. Στην πρωτότυπη εργασία της Δρ. Obermayr, χρησιμοποιήθηκε τεχνολογία μικροσυστοιχιών (Applied Biosystems) για την μέτρηση της γονιδιακής έκφρασης των 38 καρκινικών κυτταρικών σειρών (του μαστού, των ωοθηκών, του τραχήλου της μήτρας και του ενδομητρίου) και των 10 δειγμάτων μονοπύρηνων κυττάρων περιφερικού αίματος (PBMCs) από υγιείς γυναίκες δότριες. Ο στόχος της εργασίας των Obermayr και συνεργατών (2010) ήταν ο προσδιορισμός νέων γονιδιακών δεικτών για την ανίχνευση κυκλοφορούντων καρκινικών κυττάρων (CTCs) στο περιφερικό αίμα των θηλέων ασθενών με καρκίνο.[57],[12]

Καρκινικές κυτταρικές σειρές καλούνται τα καρκινικά κύτταρα που συνεχίζουν να διαιρούνται και να αυξάνονται με το χρόνο, κάτω από ορισμένες συνθήκες, σε ένα εργαστήριο. Σε σύγκριση με τα κανονικά κύτταρα, τα καρκινικά κύτταρα έχουν υποστεί κάποιες μεταλλάξεις. Αυτές είναι υπεύθυνες για την ανάπτυξη, τις μεταστάσεις και την αντίσταση που έχει στα φάρμακα ο καρκίνος. Οπότε, η χρησιμότητα των μοντέλων εξαρτάται από το κατά ποσό το μοντέλο μπορεί να περιγράψει τις αλλαγές στους πρωτεΐοντες/κύριους όγκους. Γι' αυτό το λόγο στην έρευνα για τον καρκίνο χρησιμοποιούνται συλλογές από κυτταρικές σειρές που έχουν ληφθεί από όγκους επειδή αυτές οι κυτταρικές σειρές περιέχουν εκατοντάδες ή και χιλιάδες από αυτές τις μεταλλάξεις των κυττάρων που προέκυψαν στον όγκο από τον οποίο ελήφθησαν. Οι καρκινικές κυτταρικές σειρές χρησιμοποιούνται στην έρευνα για τη μελέτη της βιολογίας του καρκίνου και για τη δοκιμή των θεραπειών του καρκίνου.[58] Η *Εγκυκλοπαίδεια Καρκινικής Κυτταρικής Σειράς* (CCLE) έχει συσταθεί ως μια προσπάθεια λεπτομερούς γενετικού χαρακτηρισμού ενός μεγάλου αριθμού ανθρώπινων καρκινικών κυτταρικών σειρών (~1000), παρέχοντας δημόσια πρόσβαση ανάλυσης και απεικόνισης του αριθμού αντιγράφων του DNA, της έκφρασης του mRNA, των δεδομένων μετάλλαξης, κ.ά.[59]

Διαφορετικές μέθοδοι biclustering χρησιμοποιούν διαφορετικές έννοιες και αντιμετωπίζουν το πρόβλημα Bicluster από άλλη οπτική ο καθένας (Υποκεφάλαια 2.8, 3.1, 3.2). Γι' αυτό το λόγο καθίσταται αρκετά δύσκολη

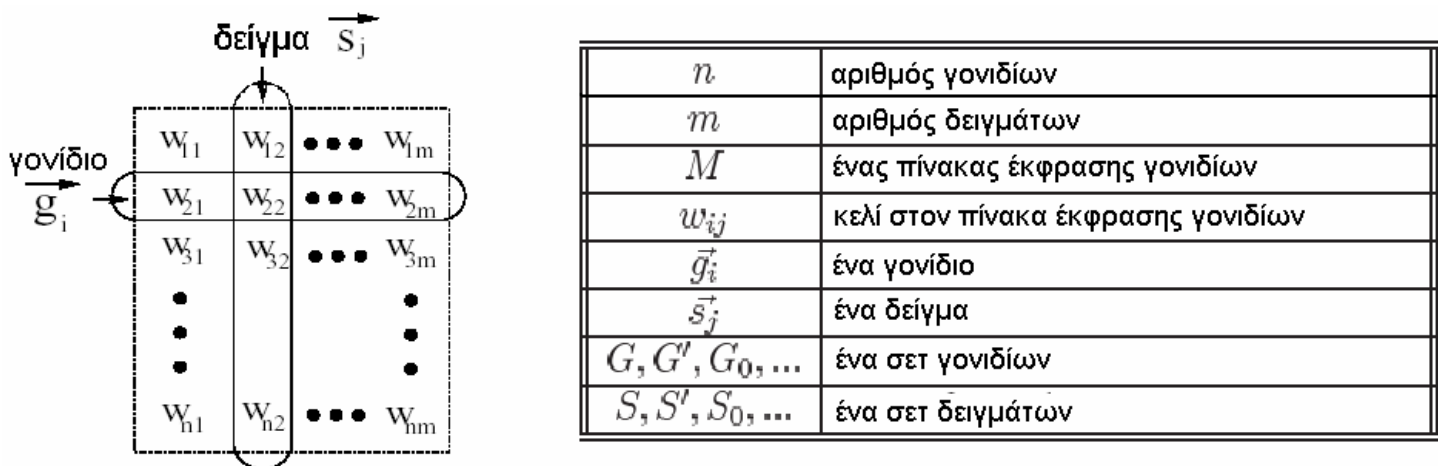
η επιλογή του πιο κατάλληλου αλγορίθμου για εφαρμογή. Ο αλγόριθμος Cheng και Church (2000) είναι ένας από τους πιο παλιούς και αξιόπιστους αλγορίθμους.[4] Τον επιλέξαμε γιατί αποτελεί μία μέθοδο ικανή να αντιμετωπίσει το μεγάλο σύνολο δεδομένων μας, παρέχει την δυνατότητα μη-επικάλυσης και τα Biclusters που εξάγει έχουν συναφείς (coherent) coherent τιμές, γεγονός που μας επιτρέπει να μελετήσουμε τη συμπεριφορά των γονιδίων κατά μήκος των σειρών χωρίς οι τιμές να είναι απαραίτητα σταθερές.

4.2. Δεδομένα

Το αρχείο προς επεξεργασία που μας δόθηκε αποτελείται από έναν πίνακα 33096 γραμμών και 38 στηλών. Οι γραμμές αντιπροσωπεύουν τα γονίδια και οι στήλες τις καρκινικές κυτταρικές σειρές.[57]

Η τεχνολογία των μικροσυστοιχιών έχει πρόσβαση σε μεγάλο μέγεθος ακολουθιών DNA κάτω από πολλές συνθήκες. Ένα σύνολο δεδομένων έκφρασης γονιδίων εκφράζεται ως πίνακας έκφρασης πραγματικών τιμών,

$$M = \{w_{ij} \mid 1 \leq i \leq n, 1 \leq j \leq m\}$$



Εικόνα 4.1. - Σύνολο δεδομένων εκφρασμένο ως πίνακα πραγματικών τιμών

όπου οι γραμμές είναι εκφράσεις προτύπων στα γονίδια και οι στήλες τα προφίλ των δειγμάτων. Ο αρχικός πίνακας δεδομένων περιέχει θόρυβο, χαμένες τιμές και αποκλίσεις που προκύπτουν από τις πειραματικές διαδικασίες. Άρα, απαιτείται ένα στάδιο προεπεξεργασίας, ώστε να λυθούν αυτά τα προβλήματα, πριν εφαρμόσουμε τις τεχνικές κατηγοριοποίησης. Με το στάδιο της προεπεξεργασίας δεν ασχολούμαστε σε αυτό το σημείο, δεδομένου ότι τα δεδομένα μας έχουν περάσει από αυτό το στάδιο.[32],[33]

Όπως θα αναφερθεί στη συνέχεια του κειμένου, ο αρχικός πίνακας των 33096 γραμμών και 38 στηλών διαμορφώθηκε σε 1000 γραμμές και 38 στήλες.

4.2.1. Επιλογή αλγορίθμου Cheng και Church

Στη δική μας εργασία, όπως αναφέρθηκε ήδη, επιλέξαμε να εφαρμόσουμε τον αλγόριθμο διπλής κατηγοριοποίησης Cheng και Church (2000).[4] Κατά την εφαρμογή του αλγορίθμου εμφανίστηκαν προβλήματα που κληθήκαμε να αντιμετωπίσουμε με μια σειρά βελτιώσεων ώστε να πετύχουμε την καλύτερη ομαδοποίηση των γονιδίων ή καταστάσεων, με βάση την πλειοψηφία των περιπτώσεων παθολογίας. Στόχοι μας ήταν: να πετύχουμε: 1) τη μείωση των αποστάσεων μεταξύ των γονιδίων, αλλά και 2) παρόμοιες διακυμάνσεις μεταξύ των γονιδίων που ομαδοποιούνται στο Bicluster.

Στην ανάλυση των δεδομένων γονιδιακής έκφρασης, πέραν της ομαδοποίησης των γονιδίων που βασίζονται στη συνολική ομοιότητα, κρίνεται κάποιες φορές απαραίτητο να διασωθούν πληροφορίες που μπορεί να χάνονταν κατά τη διάρκεια του υπολογισμού των ομάδων. Στόχος ενός τέτοιου εγχειρήματος αποτελεί η αποκάλυψη της συμμετοχής ενός γονιδίου ή μίας κατάστασης σε περισσότερες από ή υποομάδες. Η διπλή κατηγοριοποίηση επιτυγχάνει αυτή την προσπάθεια ομαδοποιώντας υποσύνολα γονιδίων και καταστάσεων με σκορ υψηλής ομοιότητας, το οποίο αποτελεί μέτρο συνοχής για αυτά.

Ένα συγκεκριμένο σκορ που αντιπροσωπεύει τα λογαριθμισμένα δεδομένα έκφρασης αποτελεί τη διακύμανση της ομάδας όλων των στοιχείων του Bicluster (mean squared residue score). Το υπόλειμμα (residue) του στοιχείου a_{ij} σε ένα Bicluster που αποτελείται από τα υποσύνολα I και J είναι το $a_{ij} - a_{iJ} - a_{iJ} + a_{IJ}$ όπου a_{ij} είναι η μέση τιμή της i -οστής σειράς του Bicluster, a_{iJ} η μέση τιμή της j -οστής στήλης, και a_{IJ} η μέση τιμή όλων των στοιχείων του Bicluster. Η τιμή «mean squared residue score» αποτελεί τη διακύμανση της ομάδας όλων των στοιχείων του Bicluster και τη μέση διακύμανση γραμμής και στήλης αντίστοιχα. Ο βασικός στόχος είναι η εύρεση μεγάλων Biclusters με χαμηλό mean squared residue score και συγκεκριμένα μικρότερο από ένα συγκεκριμένο κατώφλι (threshold). Μια ειδική περίπτωση για ένα βέλτιστο σκορ (μηδενικό «mean squared residue score») είναι ένα σταθερό Bicluster (constant bicluster) με στοιχεία με ίδια, συγκεκριμένη τιμή. Όταν ένα Bicluster έχει μη-μηδενικό σκορ είναι πάντα πιθανό αφαιρώντας κάποια γραμμή ή στήλη, το σκορ να μειωθεί, μέχρι το Bicluster που θα απομείνει να είναι σταθερό (constant).

Στην περίπτωση της γονιδιακής έκφρασης, περισσότερο ενδιαφέρον δεν παρουσιάζει η εύρεση ενός μεγάλου -αναφορικά με τον αριθμό γονιδίων- Bicluster, αλλά η ανακάλυψη ενός μεγαλύτερου αριθμού Biclusters αποτελούμενα από ομάδες γονιδίων που εμφανίζουν παρόμοια ανοδική ή καθοδική συμπεριφορά κατά μήκος μιας σειράς καταστάσεων. Ένα χαμηλό «mean squared residue score » και μια μεγάλη παραλλαγή από την σταθερή μορφή αποτελούν καλά κριτήρια για την σωστή επιλογή αυτών των γονιδίων και καταστάσεων.[4]

4.2.2. Εισαγωγή στη χρήση του αλγορίθμου Cheng και Church

Ένας πίνακας γονιδίων-καταστάσεων αποτελείται από πραγματικούς αριθμούς, με πιθανές μηδενικές τιμές σε κάποια από τα στοιχεία του. Κάθε στοιχείο-κελί του πίνακα ανακλά το επίπεδο έκφρασης ενός γονιδίου κάτω από μια συγκεκριμένη κατάσταση και αντιπροσωπεύεται από έναν πραγματικό αριθμό που αποτελεί τον λογάριθμο της σχετικής έκφρασης του mRNA του γονιδίου σ' αυτή την κατάσταση. Ο λογάριθμος χρησιμοποιείται για να μειώσει το δυναμικό εύρος τιμών των γονιδίων.[4]

Έστω X ο αριθμός των γονιδίων και Y ο αριθμός των καταστάσεων, a_{ij} ο αριθμός έκφρασης του εκάστοτε γονιδίου, και I, J τα υποσύνολα των X και Y αντίστοιχα. Το σετ (I, J) ορίζει έναν υποπίνακα A_{IJ} με το ακόλουθο «mean squared residue score »:

$$H(I, J) = \frac{1}{|I||J|} \sum_{i \in I, j \in J} (a_{ij} - \alpha_{iJ} - \alpha_{IJ} + \alpha_{IJ})^2$$

Όπου

$$\alpha_{iJ} = \frac{1}{|J|} \sum_{j \in J} a_{ij}$$

$$\alpha_{IJ} = \frac{1}{|I|} \sum_{i \in I} \alpha_{iJ}$$

$$\alpha_{IJ} = \frac{1}{|I||J|} \sum_{i \in I, j \in J} a_{ij} = \frac{1}{|I|} \sum_{i \in I} \alpha_{iJ} = \frac{1}{|J|} \sum_{j \in J} \alpha_{IJ}$$

είναι οι μέσες τιμές των γραμμών και των στηλών και η μέση τιμή όλου του πίνακα (I, J) . Ο υποπίνακας A_{ij} ονομάζεται δ -bicluster, αν $H(I, J) \leq \delta$ για κάποιο $\delta \geq 0$.

Αλγόριθμος 0-Διαγραφή Κόμβου

Κάθε πίνακας γονιδιακής έκφρασης περιέχει έναν υποπίνακα με βέλτιστο σκορ ($H(I, J)=0$) και κάθε ξεχωριστό στοιχείο είναι ένας τέτοιος υποπίνακας. Συγκεκριμένα όμως το είδος των Biclusters που αναζητούνται πρέπει να έχουν ένα μέγιστο μέγεθος όσον αφορά και τα γονίδια και τις καταστάσεις.

Ξεκινώντας από τον συνολικό πίνακα, το ερώτημα είναι πως θα γίνει η σωστή επιλογή ενός υποπίνακα με χαμηλό H score. Μια άπληστη μέθοδος είναι να αφαιρεθούν η γραμμή ή η στήλη που θα πετύχει την μεγαλύτερη μείωση του σκορ. Αυτό απαιτεί τον υπολογισμό των σκορ όλων των υποπινάκων που θα αποτελούν αποτέλεσμα κάποιας αφαίρεσης γραμμής ή στήλης. Αυτή η μέθοδος απαιτεί $O((n+m)nm)$ χρόνο για την εύρεση ενός Biclusters, όπου n και m είναι τα μεγέθη των γραμμών και των στηλών του πίνακα αντίστοιχα.[4]

Η μέθοδος που ακολουθεί ο αρχικός αυτός αλγόριθμος είναι ο υπολογισμός score H για κάθε πιθανή πρόσθεση/αφαίρεση στήλης/γραμμής και επιλογή αυτών που μειώνουν όσο το δυνατόν περισσότερο το H . Ο αλγόριθμος σταματάει όταν καμία άλλη κίνηση δεν επηρεάζει το H ή αν $H \leq \delta$.

Μετέπειτα προτάθηκε ο Αλγόριθμος 1 (Μονή Διαγραφή Κόμβου (Single Node Deletion)) με χρονική πολυπλοκότητα $O(nm)$ και ο Αλγόριθμος 2 (Πολλαπλή Διαγραφή Κόμβου (Multiple Node Deletion)) με χρονική πολυπλοκότητα $O(m \log n)$, ο συνδυασμός των οποίων αποτελεί μια αρκετά αποτελεσματική μέθοδο εύρεσης Biclusters με χαμηλό σκορ. Η ακρίβεια και η ορθότητα αυτών των αλγορίθμων στηρίζεται σε έναν αριθμό λημμάτων, όπου οι γραμμές και οι στήλες αντιμετωπίζονται σαν σημεία σε ένα διάστημα όπου η απόσταση είναι ορισμένη.[4]

Αλγόριθμος 1: Μονή Διαγραφή Κόμβου (Single Node Deletion) [4]

Είσοδος: Ένας πίνακας A, μία τιμή για το $\delta \geq 0$ που αποτελεί το μέγιστο αποδεκτό mean squared residue score.

Έξοδος: A_{IJ} , ένα δ -Bicluster υποσύνολο του αρχικού πίνακα με score όχι μεγαλύτερο του δ .

Μέθοδος:

1. Υπολογίζει a_{iJ} για όλα τα i που ανήκουν στο I , a_{iJ} για όλα τα j που ανήκουν στο J , a_{IJ} και H . Αν $H \leq \delta$ επιστρέφει τον υποπίνακα A_{IJ} αλλιώς:

2. Βρίσκει την γραμμή i που ανήκει στο I με το μεγαλύτερο

$$d(i) = \frac{1}{|J|} \sum_{j \in J} (a_{ij} - a_{iJ} - a_{Ij} + a_{IJ})^2$$

και τη στήλη j που ανήκει στο J με το μεγαλύτερο

$$d(j) = \frac{1}{|I|} \sum_{i \in I} (a_{ij} - a_{iJ} - a_{Ij} + a_{IJ})^2$$

και αφαιρεί αυτήν (στήλη ή γραμμή) με τη μεγαλύτερη ποσότητα d ανανεώνοντας τα I , J .

Η ορθότητα του βήματος 1 φαίνεται από το θεώρημα το οποίο αναφέρει ότι το σετ των γραμμών που μπορούν ολικώς ή μερικώς να αφαιρεθούν με επίδραση στη μείωση του σκορ ενός Bicluster A_{IJ} , είναι:

$R = \{i \in I, \frac{1}{|J|} \sum_{j \in J} (a_{ij} - a_{iJ} - a_{Ij} + a_{IJ})^2 > H(I, J)\}$, [60] με την έννοια ότι κάθε αφαίρεση μειώνει το τελικό σκορ. Επειδή υπάρχει πεπερασμένος αριθμός γραμμών και στηλών προς αφαίρεση, ο αλγόριθμος τερματίζει σε όχι περισσότερες από $n+m$ επαναλήψεις.

Το πρώτο στάδιο του αλγορίθμου σε κάθε επανάληψή του απαιτεί $O(nm)$ χρόνο και ο συνολικός υπολογισμός όλων των d τιμών στο στάδιο 2 χρειάζεται επίσης $O(nm)$ χρόνο. Η επιλογή της καλύτερης γραμμής και στήλης προς αφαίρεση παίρνει $O(\log n + \log m)$ χρόνο. [4]

Αλγόριθμος 2: Πολλαπλή Διαγραφή Κόμβου (Multiple Node Deletion)

Είσοδος: Ένας πίνακας A, μία τιμή για το $\delta \geq 0$ που αποτελεί το μέγιστο αποδεκτό mean squared residue score και μια τιμή για το $\alpha > 1$, ένα κατώφλι (threshold) για τη μέθοδο “Multiple node deletion”.

Έξοδος: A_{IJ} , ένα δ -Bicluster υποσύνολο του αρχικού πίνακα με score όχι μεγαλύτερο του δ .

Μέθοδος:

1. Υπολογίζει a_{iJ} για όλα τα i που ανήκουν στο I , a_{Ij} για όλα τα j που ανήκουν στο J , a_{IJ} και H . Αν $H \leq \delta$ επιστρέφει τον υποπίνακα A_{IJ} αλλιώς:

2. Αφαιρεί όλες τις γραμμές για τις οποίες ισχύει ότι:

$$\alpha H(I, J) \geq \frac{1}{|J|} \sum_{j \in J} (a_{ij} - a_{iJ} - a_{Ij} + a_{IJ})^2$$

3. Υπολογίζει ξανά τα μεγέθη a_{Ij} , a_{IJ} και H .

4. Αφαιρεί όλες τις στήλες για τις οποίες ισχύει ότι

$$\alpha H(I,J) \geq \frac{1}{|I|} \sum_{i \in I} (\alpha_{ij} - \alpha_i - \alpha_j + \alpha) ^2$$

5. Αν δεν έχει αφαιρεθεί τίποτα εκτελείται ο αλγόριθμος 1.

Αλγόριθμος 3: Προσθήκη Κόμβων (Node Addition)

Μετά τη διαγραφή κόμβων, το δ -Bicluster μπορεί να μην αποτελεί το μέγιστο δυνατό, με την έννοια ότι θα μπορούσαν να προστεθούν σε αυτό κάποιες γραμμές ή στήλες που δεν θα αυξήσουν το σκορ.[4]

Είσοδος: Ένας πίνακας A, και I, J που αποτελούν ένα δ -Bicluster.

Έξοδος: I' και J' όπου I' ανήκει στο I και J' ανήκει στο J με την ιδιότητα ότι $H(I',J') \leq H(I,J)$.

Μέθοδος:

1. Υπολογίζει α_{iJ} για όλα τα i που ανήκουν στο I, α_{iJ} για όλα τα j που ανήκουν στο J, α_{IJ} και H.

2. Προσθέτει όλες τις στήλες που δεν ανήκουν στο J για τις οποίες:

$$H(I,J) \geq \frac{1}{|I|} \sum_{i \in I} (\alpha_{ij} - \alpha_i - \alpha_j + \alpha) ^2$$

3. Υπολογίζει ξανά τα μεγέθη α_{iJ} , α_{IJ} και H.

$$H(I,J) \geq \frac{1}{|J|} \sum_{j \in J} (\alpha_{ij} - \alpha_i - \alpha_j + \alpha) ^2$$

4. Προσθέτει όλες τις γραμμές που δεν ανήκουν στο I για τις οποίες:

5. Αν σταματάει να προστίθεται κάτι, δίνονται τα I' και J' στην έξοδο.

Από τη μαθηματική ανάλυση βγαίνει το συμπέρασμα ότι η πρόσθεση γραμμών και στηλών στον αλγόριθμο 3 δεν θα αυξήσουν το σκορ. Παρόλα' αυτά το δ -Bicluster που δίνεται στην έξοδο μπορεί να μην είναι μέγιστο. Ο λόγος είναι ότι προσθέτοντας γραμμές και στήλες το σκορ μπορεί να μειωθεί κατά πολύ σε σχέση με το δ , και αυτό συμβαίνει γιατί σε κάθε επανάληψη του βήματος η πρόσθεση γίνεται με βάση το σκορ εκείνη τη δεδομένη στιγμή και όχι με βάση το ορισμένο από τη αρχή δ .

Ο αλγόριθμος 3 είναι πολύ αποτελεσματικός. Η χρονική του πολυπλοκότητα είναι συγκρίσιμη με αυτήν του αλγορίθμου 2 και είναι περίπου της τάξης του $O(nm)$. [4]

Οι αλγόριθμοι που περιγράφηκαν παραπάνω, περικλείονται στον αλγόριθμο 4 (Εύρεση συγκεκριμένου αριθμού biclusters). Οι τιμές των παραμέτρων δ , α και n καθορίζονται πριν την έναρξη του αλγορίθμου αυθαίρετα.

Σε περίπτωση μηδενικών τιμών του πίνακα γίνεται αντικατάστασή τους με τυχαίους αριθμούς που θα αποτελούν τους πρώτους υποψηφίους προς αφαίρεση στη διαδικασία αφαίρεσης κόμβων.

Επισημαίνεται ότι, επειδή ο αλγόριθμος Cheng και Church είναι ντετερμινιστικός, οι επαναληπτικές εφαρμογές του δεν θα δώσουν διαφορετικά αποτελέσματα, εκτός αν χρησιμοποιηθεί κάποιου είδους μάσκα στα αποτελέσματα. Κάθε φορά επομένως που εξάγεται ένα Bicluster, τα στοιχεία του υποπίνακα αντικαθίστανται από τυχαίους αριθμούς.[4],[44]

Τα προβλήματα που τέθηκαν προς αντιμετώπιση, όπως αναφέρθηκαν παραπάνω, είναι: **α)** ο εξαιρετικά μεγάλος αριθμός των ομάδων διπλής κατηγοριοποίησης που προκύπτουν κατά την αρχική εκτέλεση του αλγορίθμου για τα 33096 γονίδια για τον καθένα από τους τέσσερις καρκινικούς τύπους ξεχωριστά, παραδείγματος χάριν το πλήθος των Biclusters για τον καρκίνο του μαστού ανερχόταν σε 2411, (κεφάλαιο 5), **γι' αυτό τον λόγο χρειάστηκε να μειώσουμε το πλήθος των γονιδίων που θα μελετούσαμε σε 1000 γονίδια ανά καρκινικό τύπο, β)** το μεγάλο πλήθος των γονιδίων που ομαδοποιούνται σε καθεμιά ομάδα διπλής κατηγοριοποίησης κατά την αρχική εκτέλεση του αλγορίθμου για τα 33096 γονίδια, **γι' αυτό τον λόγο χρειάστηκε να μειώσουμε το πλήθος των εξεταζόμενων Biclusters σε τρία για κάθε μία από τις τρεις εκτελέσεις του αλγορίθμου ανά καρκινικό τύπο και γ)** η ανομοιομορφία μεταξύ των γονιδίων που συγκεντρώνονται σε μία ομάδα διπλής κατηγοριοποίησης, δηλαδή τα γονίδια διέφεραν στην τάση παρουσιάζοντας μεγάλες διακυμάνσεις, αλλά και στη μεταξύ τους απόσταση καλύπτοντας ένα μεγάλο εύρος τιμών, **για αυτό το λόγο πραγματοποιήσαμε κάποιες τροποποιήσεις στον υπολογισμό του di (mean squared residue της γραμμής i) και dj (mean squared residue της στήλης j),** όπως παρουσιάζονται αναλυτικά παρακάτω.

Στην παρούσα μελέτη, ο αλγόριθμος Cheng και Church εφαρμόστηκε 18 φορές, όπως αποτυπώνεται στο σχεδιάγραμμα ροής (Σχήμα 4.1).

Συγκεκριμένα, ο αλγόριθμος εκτελέστηκε τρεις φορές:

- ❖ Στην αρχική του μορφή, όπως έχει προταθεί από τους Cheng & Church → 1η Εκτέλεση – $didj$.
- ❖ Με αλλαγή του di → 2η Εκτέλεση- di_new, dj .
- ❖ Με αλλαγή του di και του dj → 3η Εκτέλεση - di_new, dj_new .
 - Για τα 1000 πρώτα γονίδια για α) τον καρκίνο του μαστού, β) τον καρκίνο του τραχήλου της μήτρας, γ) τον καρκίνο των ωοθηκών, και δ) τον καρκίνο του ενδομητρίου.
 - Για τα 33096 γονίδια για τον καρκίνο του μαστού.
 - Για τα 33096 γονίδια για τον καρκίνο του τραχήλου της μήτρας.

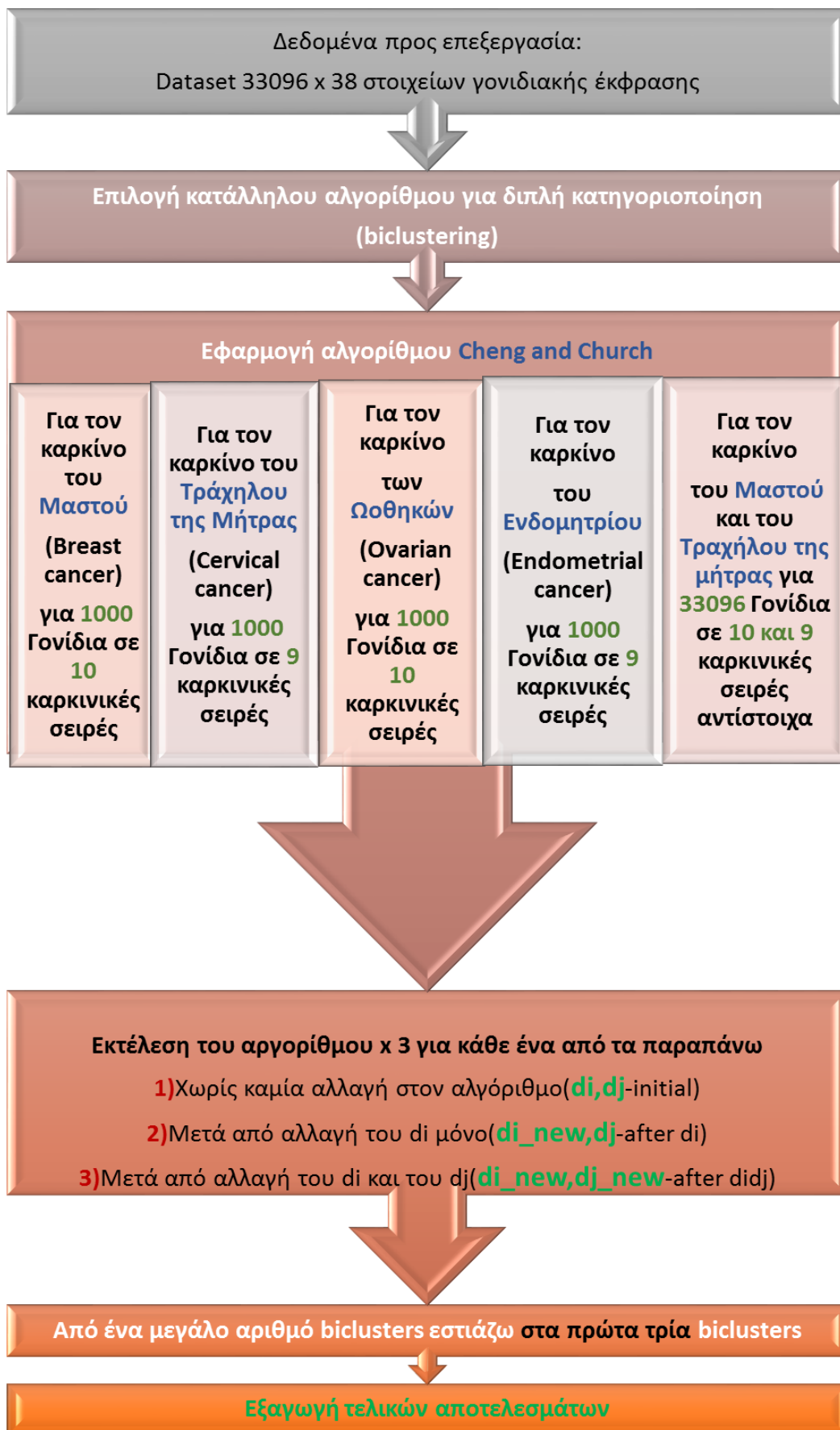
Χρειάστηκαν όλες οι παραπάνω εφαρμογές στον αλγόριθμο ώστε να επιλεγούν και να μελετηθούν τα τρία πρώτα Biclusters προκειμένου να καταλήξουμε στις πιο ενδεικτικές τιμές.

Όπως αναφέρθηκε, οι τύποι των καρκίνων από τους οποίους προέρχονται οι κυτταρικές σειρές είναι:

καρκίνος του ενδομητρίου (endometrial cancer), καρκίνος του τραχήλου της μήτρας (cervical cancer) και στον καθένα από αυτούς ανήκουν 10, 9, 9 και 10 σειρές αντίστοιχα.

Η εξαγωγή ενός μεγάλου αριθμού ομάδων διπλής κατηγοριοποίησης (Biclusters) σε πραγματικά δεδομένα μπορεί να οδηγήσει σε αποτελέσματα, τα οποία είναι δύσκολο να ερμηνευθούν βιολογικά. Ως εκ τούτου, εστιάσαμε στην μελέτη των 3 πρώτων ομάδων διπλής κατηγοριοποίησης (Biclusters).

Σχήμα 4.1. - Σχεδιάγραμμα ροής



Σ' αυτό το σημείο παρουσιάζεται ένας συγκεντρωτικός πίνακας με τον αριθμό των Biclusters που δημιουργούνται κατά την εκτέλεση του αλγορίθμου τόσο για τα 33096 γονίδια, όσο και για τα 1000 γονίδια για τα τρία πρώτα Biclusters κάθε φορά (Πίνακας 4.1).

Πίνακας 4.1. – Συγκεντρωτικά Αποτελέσματα Αριθμού Γονιδίων για Κάθε Ένα Από τα Τρία Biclusters που μελετήθηκαν

ΚΑΡΚΙΝΟΣ ΤΟΥ ΜΑΣΤΟΥ(BREAST CANCER)-ΓΙΑ 33096 ΓΟΝΙΔΙΑ									
ΑΡΧΙΚΗ ΚΑΤΑΣΤΑΣΗ(INITIAL)			ΑΛΛΑΓΗ DI(AFTER DI)			ΑΛΛΑΓΗ DI,DJ(AFTER DIDJ)			
BICLUSTER1	BICLUSTER2	BICLUSTER3	BICLUSTER1	BICLUSTER2	BICLUSTER3	BICLUSTER1	BICLUSTER2	BICLUSTER3	
300	126	105	59	11	5	35	11		17
ΚΑΡΚΙΝΟΣ ΤΟΥ ΜΑΣΤΟΥ(BREAST CANCER)-ΓΙΑ 1000 ΓΟΝΙΔΙΑ									
ΑΡΧΙΚΗ ΚΑΤΑΣΤΑΣΗ(INITIAL)			ΑΛΛΑΓΗ DI(AFTER DI)			ΑΛΛΑΓΗ DI,DJ(AFTER DIDJ)			
BICLUSTER1	BICLUSTER2	BICLUSTER3	BICLUSTER1	BICLUSTER2	BICLUSTER3	BICLUSTER1	BICLUSTER2	BICLUSTER3	
30	18	27	4	6	2	9	7		5
ΚΑΡΚΙΝΟΣ ΤΟΥ ΤΩΝ ΩΘΗΚΩΝ(OVARIAN CANCER)-ΓΙΑ 33096 ΓΟΝΙΔΙΑ									
ΑΡΧΙΚΗ ΚΑΤΑΣΤΑΣΗ(INITIAL)			ΑΛΛΑΓΗ DI(AFTER DI)			ΑΛΛΑΓΗ DI,DJ(AFTER DIDJ)			
BICLUSTER1	BICLUSTER2	BICLUSTER3	BICLUSTER1	BICLUSTER2	BICLUSTER3	BICLUSTER1	BICLUSTER2	BICLUSTER3	
539	158	160	200	54	23	86	23		35
ΚΑΡΚΙΝΟΣ ΤΟΥ ΤΩΝ ΩΘΗΚΩΝ(OVARIAN CANCER)-ΓΙΑ 1000 ΓΟΝΙΔΙΑ									
ΑΡΧΙΚΗ ΚΑΤΑΣΤΑΣΗ(INITIAL)			ΑΛΛΑΓΗ DI(AFTER DI)			ΑΛΛΑΓΗ DI,DJ(AFTER DIDJ)			
BICLUSTER1	BICLUSTER2	BICLUSTER3	BICLUSTER1	BICLUSTER2	BICLUSTER3	BICLUSTER1	BICLUSTER2	BICLUSTER3	
26	29	15	3	4	2	10	9		8
ΚΑΡΚΙΝΟΣ ΤΟΥ ΕΝΔΟΜΗΤΡΙΟΥ(ENDOMETRIAL CANCER)-ΓΙΑ 33096 ΓΟΝΙΔΙΑ									
ΑΡΧΙΚΗ ΚΑΤΑΣΤΑΣΗ(INITIAL)			ΑΛΛΑΓΗ DI(AFTER DI)			ΑΛΛΑΓΗ DI,DJ(AFTER DIDJ)			
BICLUSTER1	BICLUSTER2	BICLUSTER3	BICLUSTER1	BICLUSTER2	BICLUSTER3	BICLUSTER1	BICLUSTER2	BICLUSTER3	
940	204	233	279	49	107	172	30		70
ΚΑΡΚΙΝΟΣ ΤΟΥ ΕΝΔΟΜΗΤΡΙΟΥ(ENDOMETRIAL CANCER)-ΓΙΑ 1000 ΓΟΝΙΔΙΑ									
ΑΡΧΙΚΗ ΚΑΤΑΣΤΑΣΗ(INITIAL)			ΑΛΛΑΓΗ DI(AFTER DI)			ΑΛΛΑΓΗ DI,DJ(AFTER DIDJ)			
BICLUSTER1	BICLUSTER2	BICLUSTER3	BICLUSTER1	BICLUSTER2	BICLUSTER3	BICLUSTER1	BICLUSTER2	BICLUSTER3	
42	22	23	8	2	2	16	8		5
ΚΑΡΚΙΝΟΣ ΤΟΥ ΤΡΑΧΗΛΟΥ(CERVICAL CANCER)-ΓΙΑ 33096 ΓΟΝΙΔΙΑ									
ΑΡΧΙΚΗ ΚΑΤΑΣΤΑΣΗ(INITIAL)			ΑΛΛΑΓΗ DI(AFTER DI)			ΑΛΛΑΓΗ DI,DJ(AFTER DIDJ)			
BICLUSTER1	BICLUSTER2	BICLUSTER3	BICLUSTER1	BICLUSTER2	BICLUSTER3	BICLUSTER1	BICLUSTER2	BICLUSTER3	
544	295	250	207	75	125	110	25		85
ΚΑΡΚΙΝΟΣ ΤΟΥ ΤΡΑΧΗΛΟΥ(CERVICAL CANCER)-ΓΙΑ 1000 ΓΟΝΙΔΙΑ									
ΑΡΧΙΚΗ ΚΑΤΑΣΤΑΣΗ(INITIAL)			ΑΛΛΑΓΗ DI(AFTER DI)			ΑΛΛΑΓΗ DI,DJ(AFTER DIDJ)			
BICLUSTER1	BICLUSTER2	BICLUSTER3	BICLUSTER1	BICLUSTER2	BICLUSTER3	BICLUSTER1	BICLUSTER2	BICLUSTER3	
36	24	16	12	4	4	13	11		5

Παρατηρήσεις:

Στον συγκεντρωτικό Πίνακα 4.1 παρατηρούμε ότι κατά την εκτέλεση του αλγορίθμου για το σύνολο των 33096 γονιδίων, ο όγκος των γονιδίων που ομαδοποιούνται σε ένα Bicluster είναι μεγάλος. Γι' αυτό το λόγο, προκειμένου να διακρίνουμε και να εξάγουμε συμπεράσματα για τη συμπεριφορά των γονιδίων εκτελέσαμε τον αλγόριθμο για τα 1000 πρώτα γονίδια από τα 33096 και εστίασαμε στη μελέτη των τριών πρώτων -στη σειρά- Biclusters. Στην ακόλουθη ενότητα αναφέρονται αναλυτικά τα προβλήματα που κληθήκαμε να αντιμετωπίσουμε καθώς και η προτεινόμενη μεθοδολογία για την επίλυσή τους.

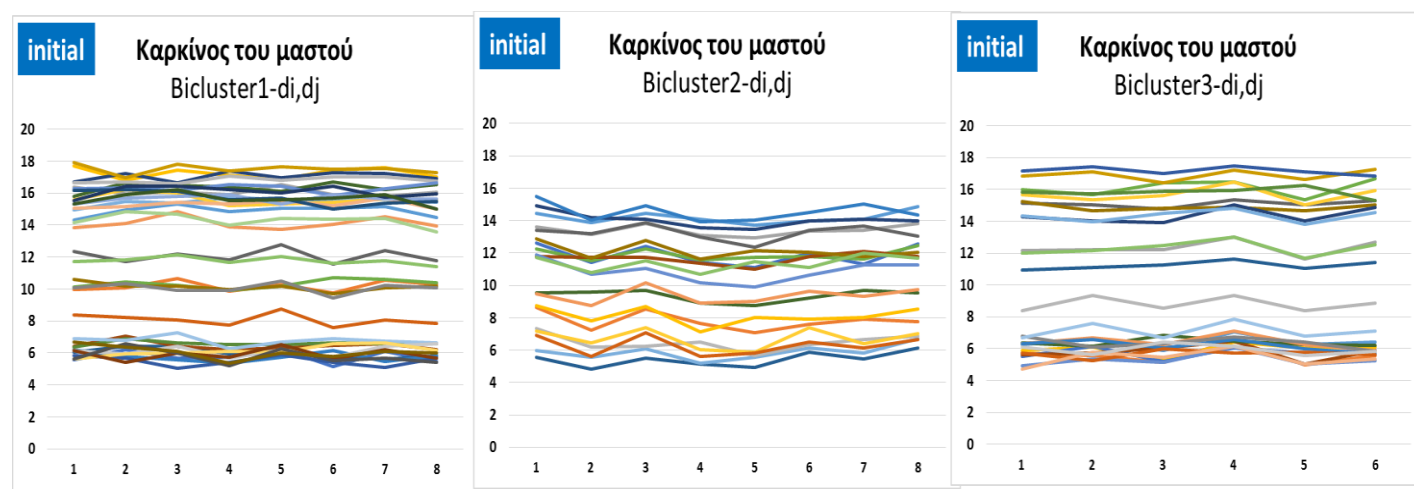
4.3. Προτεινόμενη Μεθοδολογία

Στην παρούσα διπλωματική εργασία, εφαρμόστηκε ο αλγόριθμος Cheng και Church, όπως αναλύθηκε στο προηγούμενο κεφάλαιο (με τις προτεινόμενες αλλαγές), τρεις φορές, για κάθε τύπο καρκίνου για συγκεκριμένες τιμές των παραμέτρων δ και α και με τροποποιήσεις στον υπολογισμό του d_i (mean squared residue της γραμμής i) και d_j (mean squared residue της στήλης j).

Αρχικά, τρέξαμε τον αλγόριθμο Cheng και Church για τα **1000 πρώτα γονίδια** από τα 33096 από τις **κυτταρικές σειρές του καρκίνου του μαστού (Breast cancer) και εστιάσαμε στα τρία πρώτα Biclusters**.

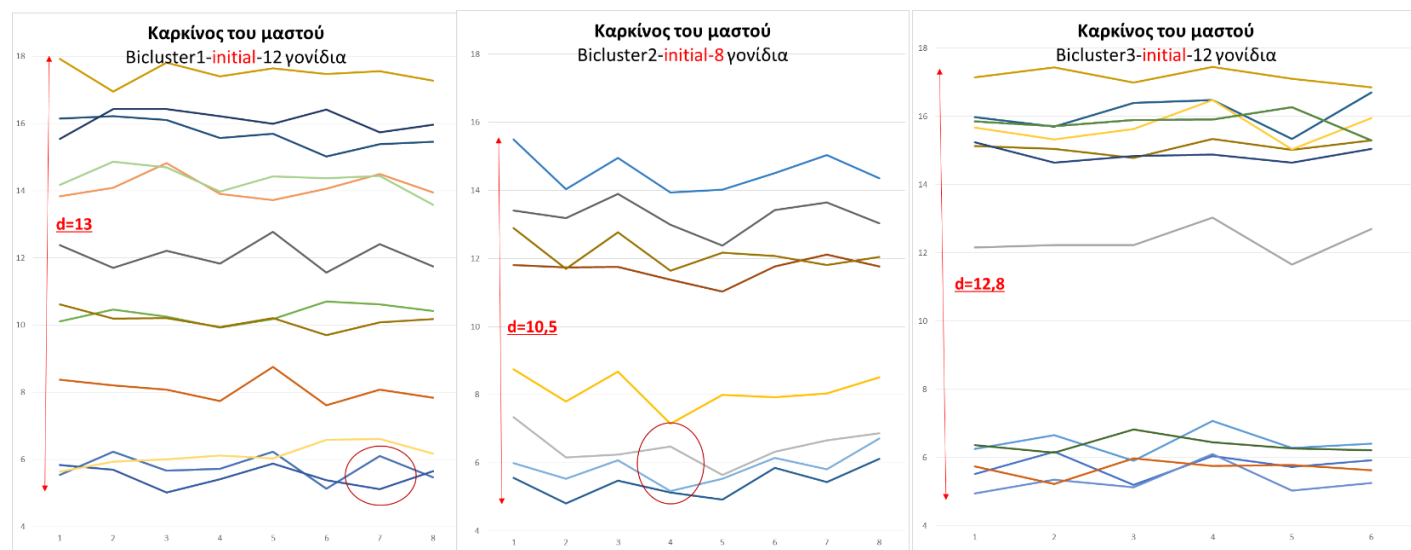
Οι παρακάτω Εικόνες (4.1.α, 4.1.β, 4.1.γ) παρουσιάζουν το συνολικό πλήθος των γονιδίων και τη συμπεριφορά τους σε κάθε ένα από τα τρία Biclusters μετά τη πρώτη εκτέλεση του αλγορίθμου (Εικόνα 4.1.α) καθώς και δυο σημαντικά προβλήματα που προκύπτουν κατά τη πρώτη εκτέλεση του αλγορίθμου (Εικόνες 4.1.β και 4.1.γ).

Εικόνα 4.1.α. - Πρώτη Εκτέλεση του Αλγορίθμου (d_i, d_j) για 1000 Γονίδια του Καρκίνου του Μαστού



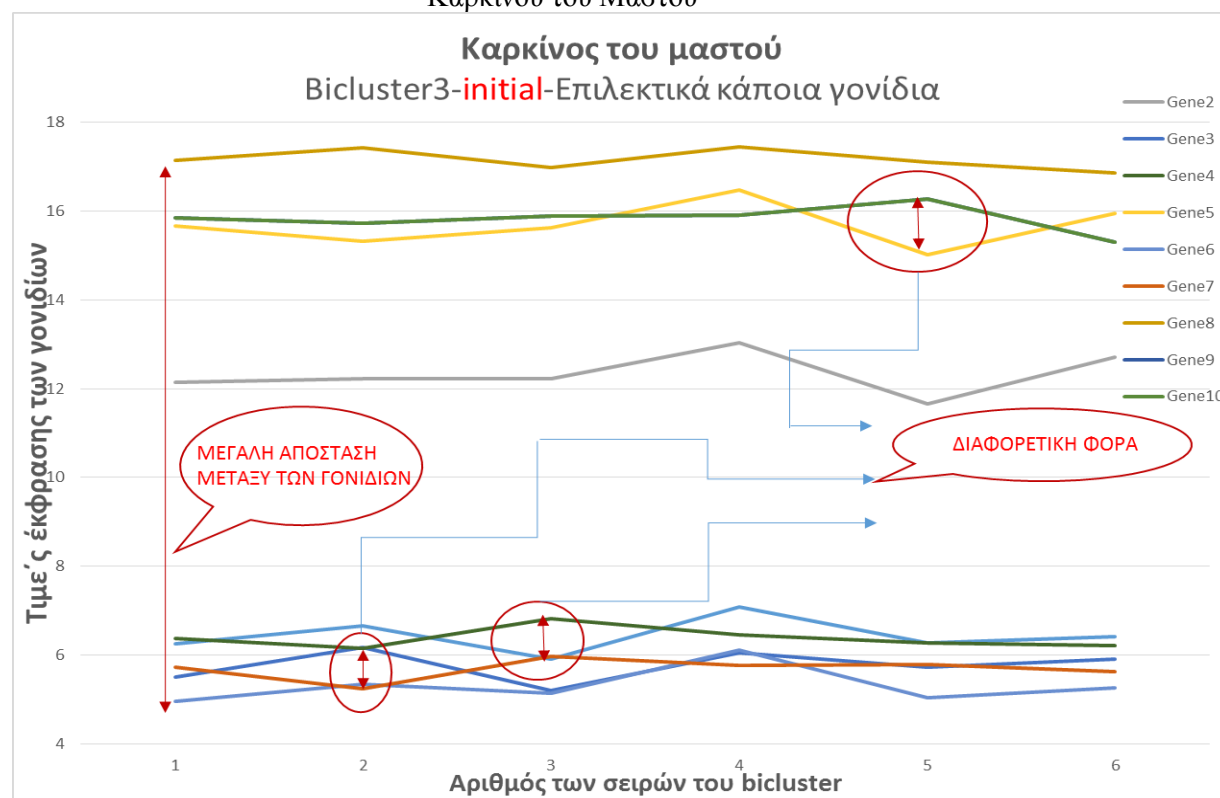
Στην Εικόνα 4.1.α παρατηρούμε το συνολικό πλήθος των γονιδίων και τη συμπεριφορά τους σε κάθε ένα από τα τρία Biclusters μετά τη πρώτη εκτέλεση του αλγορίθμου. Μικραίνοντας τα όρια του άξονα y (σε κλίμακα 4 έως 18 από 0 έως 20 που ήταν η αρχική) και αφαιρώντας ένα μεγάλο αριθμό γονιδίων (12 από 30 για το Bicluster 1, 8 από 18 για το Bicluster 2, και 12 από 27 για το Bicluster 3) από τα διαγράμματα, μπορούμε να διακρίνουμε καλύτερα δυο σημαντικά προβλήματα που προκύπτουν, όπως αποτυπώνονται στην Εικόνα 4.1.β.

Εικόνα 4.1.β. - Επιλεκτική Εμφάνιση Γονιδίων μετά την Πρώτη εκτέλεση του αλγορίθμου (d_i, d_j) για 1000 Γονίδια του Καρκίνου του Μαστού



Στοχεύοντας στη μελέτη της συμπεριφοράς των γονιδίων, εστίασαμε στην εμφάνιση 9 γονιδίων από το Biclust 3 της Εικόνας 1.β., όπως αποτυπώνεται στην Εικόνα 4.1.γ.

Εικόνα 4.1.γ. - Εστιασμένη Εμφάνιση Γονιδίων μετά την Πρώτη εκτέλεση του αλγορίθμου (d_i, d_j) για 1000 Γονίδια του Καρκίνου του Μαστού



Στην Εικόνα 4.1.γ διακρίνουμε τα εξής :

- ❑ Στα κυκλωμένα σημεία, τα οποία είναι ενδεικτικά, **παρατηρούμε διαφορετική τάση (different trend)** Αυτό υποδηλώνει ότι ο αλγόριθμος χρησιμοποιεί είτε απόλυτες τιμές είτε τετραγωνική ρίζα.

✓ Πρέπει να αποφεύγονται κανονικοποιήσεις ή αλλαγές προσήμου.

Το δεύτερο πρόβλημα που διαφαίνεται στην Εικόνα 4.1.γ είναι **η μεγάλη απόσταση μεταξύ των καμπυλών του διαγράμματος.**

- ❑ Ο αλγόριθμος σε κάποιο σημείο χρησιμοποιεί x_i, y_i όπου $i = 1, \dots, 10$

Και $\text{correlation}(x, y) = \frac{\langle x_i, y_i \rangle}{|x| * |y|}$. Ο αλγόριθμος δηλαδή αφαιρεί το μέτρο και διαιρεί με $|x| * |y|$.

✓ Πρέπει λοιπόν να αποφεύγονται κανονικοποιήσεις του μέτρου.

- ❑ Αντιμετώπιση “προβλημάτων”

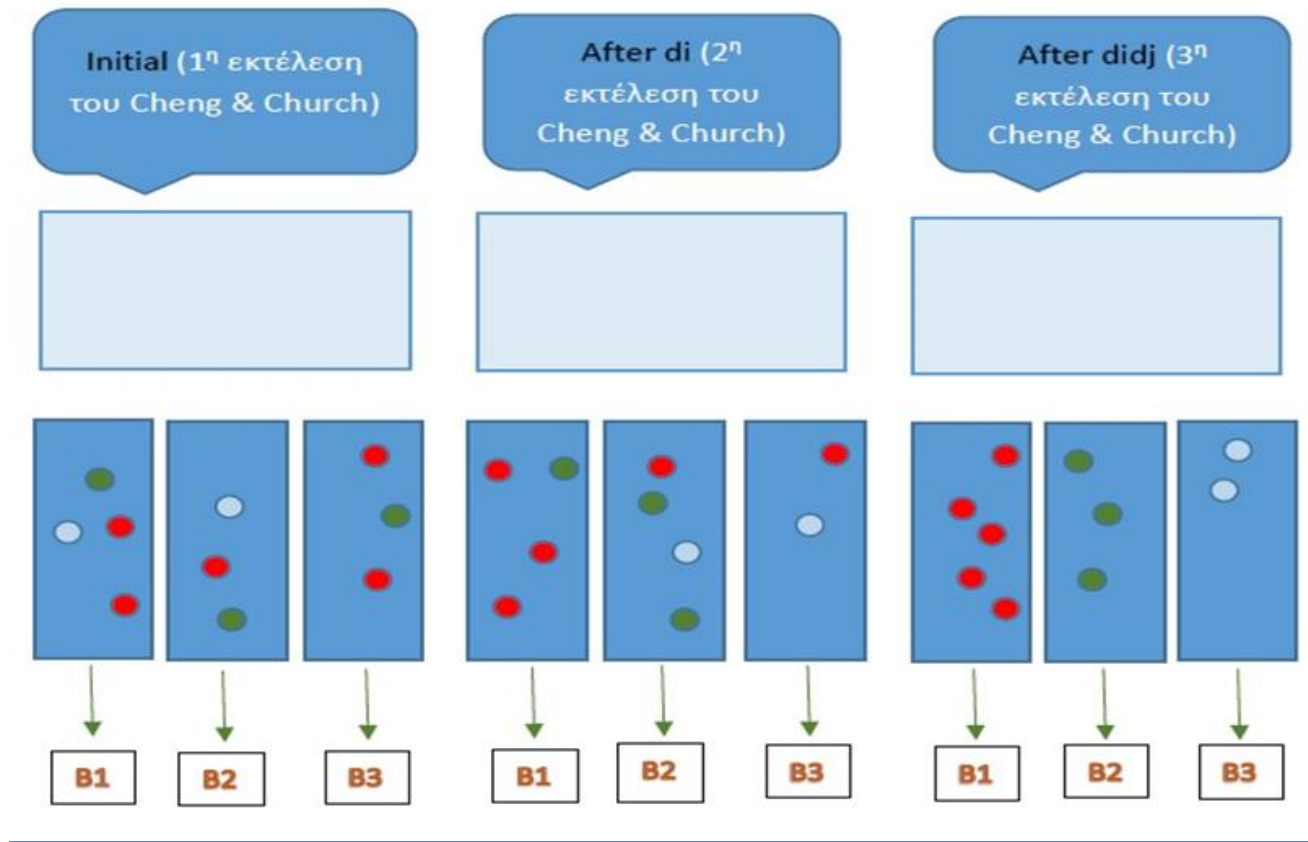
Με βάση τα παραπάνω προβλήματα, προκειμένου να εξαχθούν ταυτόχρονα ομάδες γονιδίων με ομοιόμορφη συμπεριφορά, πραγματοποιήθηκε η εφαρμογή της μεθόδου διπλής κατηγοριοποίησης δεδομένων γονιδιακής έκφρασης πέραν της αρχικής μορφής του αλγορίθμου και με προτεινόμενες τροποποιήσεις.

Ο αλγόριθμος εφαρμόστηκε, επομένως, τρεις φορές για κάθε τύπο καρκίνου για συγκεκριμένες τιμές των δ και α και με τροποποιήσεις στον υπολογισμό του di (mean squared της γραμμής i) και dj (mean squared της στήλης j).

Όπως αναφέρθηκε και παραπάνω τα δεδομένα γονιδιακής έκφρασης αναπαρίστανται με τη μορφή πινάκων, όπου οι γραμμές αποτελούν τα γονίδια και οι στήλες τις διάφορες πειραματικές συνθήκες ή καταστάσεις (π.χ. διαφορετικές κυτταρικές σειρές). Η ανάλυση των πινάκων τοποθετείται σε δύο βασικές κατευθύνσεις. Στην ομαδοποίηση των γονιδίων ή καταστάσεων με βάση την πλειοψηφία των περιπτώσεων παθολογίας και στην κατηγοριοποίηση και πρόβλεψη νέων γονιδίων ή δειγμάτων στηριζόμενοι σε ήδη γνωστή βιολογική πληροφορία. Στο παρακάτω Σχήμα 4.2 φαίνεται η μετακίνηση των γονιδίων από το ένα Biclust στο άλλο κατά την διαδικασία αλλαγής του di και dj με στόχο την καλύτερη ομαδοποίηση γονιδίων ως προς το μέγεθος έκφρασης και τάσης.

Στο Σχήμα 4.2 παρουσιάζεται με έναν απλό τρόπο τι πετυχαίνουμε με το να πραγματοποιούμε αλλαγές στο di και στο dj .

Σχήμα 4.2. - Σχηματική Αναπαράσταση της Συμπεριφοράς των Γονιδίων στις Τρεις Εκτελέσεις του Αλγορίθμου Cheng και Church



Με B1, B2, B3 αναφερόμαστε στο Bicluster1, Bicluster2, Bicluster3 αντίστοιχα. Ενδεικτικά, όπως αποτυπώνεται στο Σχήμα 4.2, τα γονίδια με κόκκινο χρώμα εμφανίζονται σε όλες τις ομάδες. Μετά την αλλαγή του d_i έχει φύγει ένα «κόκκινο» γονίδιο και έχει τοποθετηθεί στο Bicluster1. Τέλος, μετά την αλλαγή και του d_j φαίνεται πως όλα τα «κόκκινα» γονίδια έχουν ομαδοποιηθεί στο Bicluster1. Αντίστοιχη συμπεριφορά εμφανίζουν και τα «πράσινα» και τα «γαλάζια» γονίδια. Έτσι, μετά το τέλος των τριών εκτελέσεων του αλγορίθμου Cheng και Church παρατηρούμε πως έχουμε πετύχει την ομαδοποίηση γονιδίων με παρόμοια χαρακτηριστικά.

Πιο αναλυτικά, εκτελέστηκε ο αλγόριθμος Cheng και Church για κάθε καρκίνο στην κανονική του μορφή. Στη συνέχεια εκτελέστηκε ο αλγόριθμος C&C για τον κάθε τύπο καρκίνου με αλλαγή στον τρόπο υπολογισμού του d_i .

Αντί για τον αρχικό τύπο του d_i :

$$d(i) = \frac{1}{|J|} \sum_{j \in J} (a_{ij} - \bar{a}_{i\cdot} - \bar{a}_{\cdot j} + \bar{a}_{\cdot\cdot})^2$$

Υπολογίζει μια μέση απόσταση του γονιδίου αυτού, σε όλες τις καταστάσεις, και επιχειρεί να κρατήσει τη μέση αυτή απόσταση όσο το δυνατόν μικρότερη. Επομένως το νέο d_i ορίζεται ως :

2^η εκτέλεση
Εκτέλεση ανά
γονίδιο

$$d(i) = \frac{1}{|J|} \sum_{j \in J} (|a_{ij} - a_{iJ}|^2 + |a_{ij} - a_{Ij}|^2) \cdot$$

Το $d(i)$ αναφέρεται στην ΤΑΣΗ των γονιδίων. Με τον όρο $|a_{ij} - a_{iJ}|^2$ τα γονίδια, σε όλες τις χρονικές στιγμές, δεν κινούνται πολύ γρήγορα γύρω από το Μέσο Όρο του ίδιου του γονιδίου, δηλαδή δεν έχουμε σημαντικές αλλαγές. Με τον όρο $|a_{ij} - a_{Ij}|^2$ τα γονίδια μέσα στο Biclustet, σε όλες τις χρονικές στιγμές, συμπεριφέρονται σαν τον Μέσο Όρο του συνόλου. Ουσιαστικά, ο αλγόριθμος κρατάει γονίδια που έχουν ίδια ή παρόμοια τάση.

Στη συνέχεια εκτελέστηκε ο αλγόριθμος Cheng και Church για τον κάθε τύπο καρκίνου με αλλαγή στον τρόπο υπολογισμού του dj .

Αντί για τον αρχικό τύπο του dj :

$$d(j) = \frac{1}{|I|} \sum_{i \in I} (a_{ij} - a_{iJ} - a_{Ij} + a_{IJ})^2$$

Για κάθε κατάσταση ή κυτταρική σειρά j , το κριτήριο dj επιχειρεί να κρατήσει όλα τα γονίδια του Biclustet σε μικρή απόσταση μεταξύ τους. Επομένως το νέο dj ορίζεται ως :

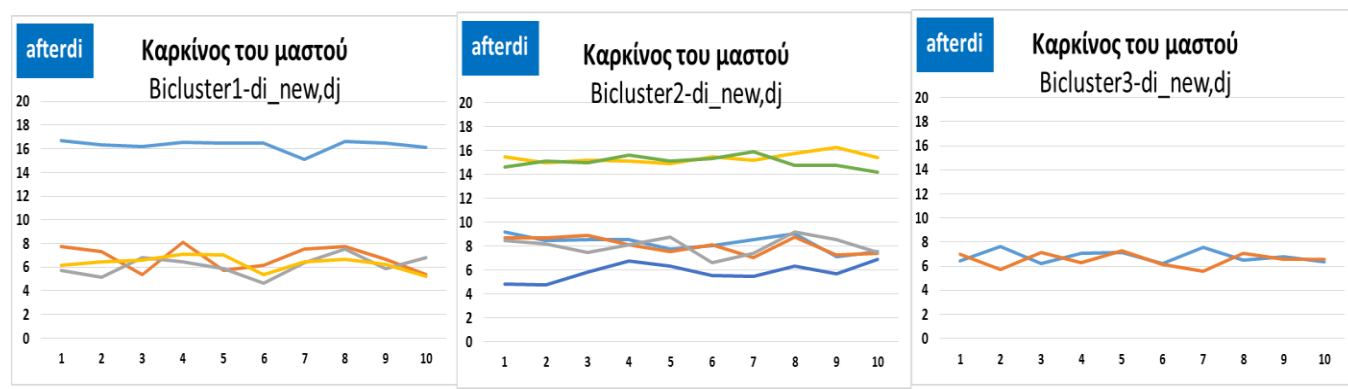
3^η εκτέλεση
Εκτέλεση ανά
κατάσταση

$$d(j) = \frac{1}{|I|} \sum_{i \in I} (|a_{ij} - a_{iJ}|^2 + |a_{ij} - a_{Ij}|^2) \cdot$$

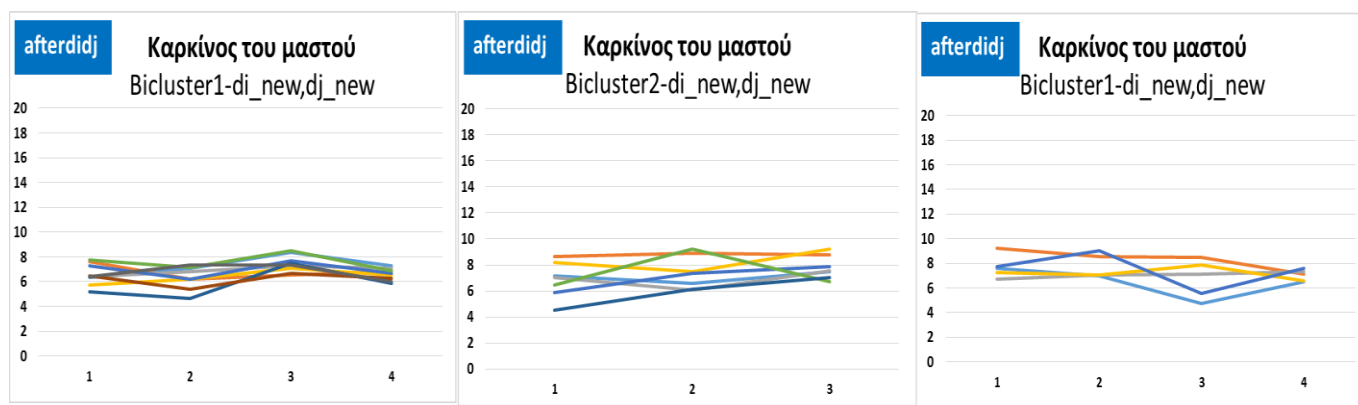
Το $d(j)$ αναφέρεται στην ΤΙΜΗ των γονιδίων. Με τον όρο $|a_{ij} - a_{iJ}|^2$ τα γονίδια, σε κάθε χρονική στιγμή, βρίσκονται κοντά στο Μέσο Όρο του ίδιου του γονιδίου. Με τον όρο $|a_{ij} - a_{Ij}|^2$ τα γονίδια μέσα στο Biclustet, σε κάθε χρονική στιγμή, είναι μαζεμένα σε ένα Μέσο Όρο. Ουσιαστικά, ο αλγόριθμος κρατάει καταστάσεις στις οποίες τα γονίδια έχουν όμοια συμπεριφορά γύρω από ένα γενικό Μέσο Όρο.

Οι Εικόνες 4.2 και 4.3 αναδεικνύουν τις βελτιώσεις της συμπεριφοράς των γονιδίων σε κάθε Biclustet μετά την εφαρμογή των παραπάνω τροποποιήσεων παρουσιάζονται τα αποτελέσματα των γονιδίων μετά την δεύτερη και στη συνέχεια τη τρίτη εκτέλεση του αλγορίθμου.

Εικόνα 4.2. - Δεύτερη Εκτέλεση του Αλγορίθμου (di_new, dj) για 1000 Γονίδια του Καρκίνου του Μαστού



Εικόνα 4.3. - Τρίτη Εκτέλεση του Αλγορίθμου (di_new, dj_new) για 1000 Γονίδια του Καρκίνου του Μαστού



Στο Κεφάλαιο 5 παρουσιάζονται αναλυτικά τα αποτελέσματα της εργασίας μας για όλες τις κυτταρικές σειρές των τεσσάρων καρκινικών τύπων που μελετήσαμε.

5

Αποτελέσματα

Παρακάτω παρουσιάζεται η σειρά με την οποία εξήχθησαν τα αποτελέσματα.

1. Αρχικά, τρέξαμε τον αλγόριθμο Cheng και Church για τα 1000 πρώτα γονίδια από τα 33096 για τους παρακάτω τύπους καρκίνου:

- α. Καρκίνος του μαστού (Breast cancer)
- β. Καρκίνος του τραχήλου της μήτρας (Cervical cancer)
- γ. Καρκίνος των ωοθηκών (Ovarian cancer)
- δ. Καρκίνος του ενδομητρίου (Endometrial cancer)

Στους τέσσερις αυτούς τύπους καρκίνου, μειώνοντας τον αριθμό των γονιδίων από 33096 σε 1000, μπορέσαμε να διακρίνουμε καλύτερα τη συμπεριφορά των γονιδίων αναφορικά με την τάση και την τιμή τους και να προβούμε σε κατάλληλες αλλαγές στον αλγόριθμο.

Τα αποτελέσματα παρουσιάζονται στη συνέχεια της εργασίας με την εξής σειρά:

❖ 1.α. Για τον καρκίνο του μαστού για 1000 Γονίδια, για τα γονίδια που προκύπτουν μετά την:

- α) πρώτη εκτέλεση του αλγορίθμου (di, dj ή initial),
- β) δεύτερη εκτέλεση του αλγορίθμου (di_new, dj ή after di) και
- γ) τρίτη εκτέλεση του αλγορίθμου (di_new, dj_new ή after di_dj)

❖ 1.β. Για τον καρκίνο του τραχήλου της μήτρας για 1000 Γονίδια που προκύπτουν μετά την:

- α) πρώτη εκτέλεση του αλγορίθμου (di, dj ή initial),
- β) δεύτερη εκτέλεση του αλγορίθμου (di_new, dj ή after di) και
- γ) τρίτη εκτέλεση του αλγορίθμου (di_new, dj_new ή after di_dj)

❖ 1.γ. Για τον καρκίνο των ωοθηκών για 1000 Γονίδια, για τα γονίδια που προκύπτουν μετά την:

Όλα τα γονίδια που προκύπτουν μετά την:

- α) πρώτη εκτέλεση του αλγορίθμου (di, dj ή initial),
- β) δεύτερη εκτέλεση του αλγορίθμου (di_new, dj ή after di) και
- γ) τρίτη εκτέλεση του αλγορίθμου (di_new, dj_new ή after di_dj)

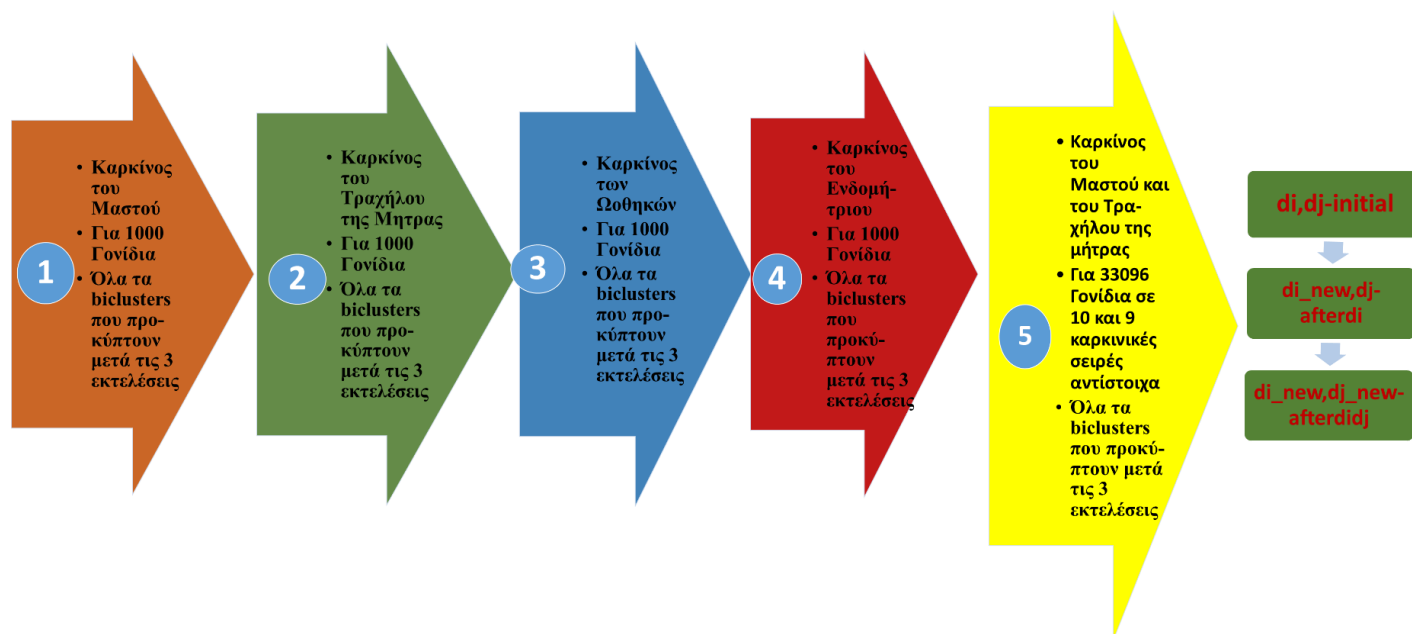
❖ 1.δ. Για τον καρκίνο του ενδομητρίου για 1000 Γονίδια, για τα γονίδια που προκύπτουν μετά την:

- α) πρώτη εκτέλεση του αλγορίθμου (di, dj ή initial),
- β) δεύτερη εκτέλεση του αλγορίθμου (di_new, dj ή after di) και
- γ) τρίτη εκτέλεση του αλγορίθμου (di_new, dj_new ή after di_dj)

2. Για τον καρκίνο του μαστού και τον καρκίνο του τραχήλου της μήτρας, για 33096 Γονίδια, για τα γονίδια που προκύπτουν μετά την:

- α) πρώτη εκτέλεση του αλγορίθμου (di, dj ή initial),
- β) δεύτερη εκτέλεση του αλγορίθμου (di_new, dj ή after di) και
- γ) τρίτη εκτέλεση του αλγορίθμου (di_new, dj_new ή after di_dj)

Πίνακας 5.1. – Συγκεντρωτικός πίνακας αποτύπωσης προτεινόμενης μεθοδολογίας

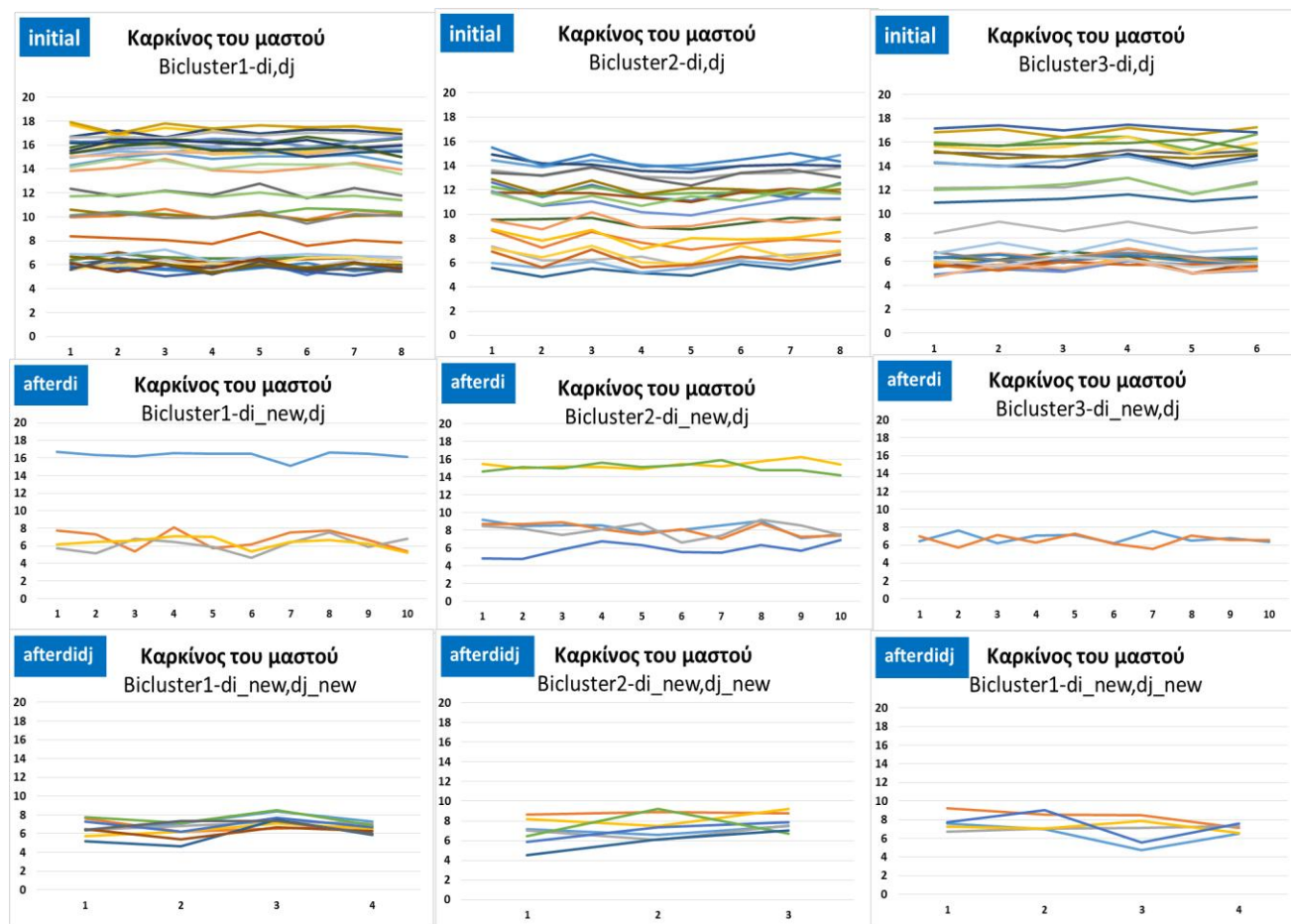


Στον Πίνακα 5.1 απεικονίζεται η σειρά με την οποία θα παρουσιαστούν τα αποτελέσματα στην ακόλουθη ενότητα.

5.1. Αποτελέσματα biclustering από κυτταρικές σειρές καρκίνου του μαστού

Η παρακάτω συγκεντρωτική Εικόνα 5.1 παρουσιάζει το συνολικό πλήθος των γονιδίων και τη συμπεριφορά τους για τον καρκίνο του μαστού για τα 1000 πρώτα γονίδια, σε κάθε ένα από τα τρία πρώτα Biclusters.

Εικόνα 5.1. - Συγκεντρωτική Παρουσίαση των Τριών εκτελέσεων του Αλγορίθμου Cheng και Church για 1000 Γονίδια από τις Κυτταρικές Σειρές Καρκίνου του Μαστού

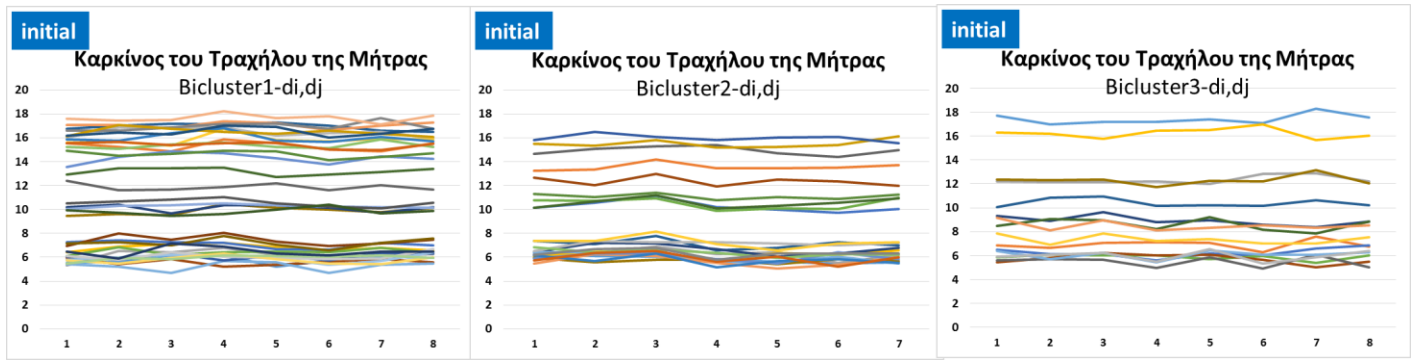


5.2. Αποτελέσματα biclustering από κυτταρικές σειρές του καρκίνου του τραχήλου της μήτρας

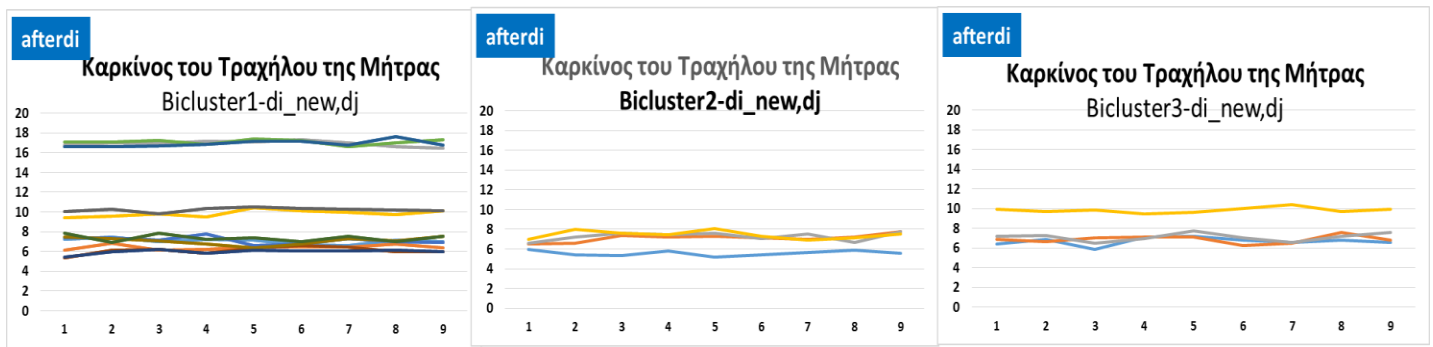
Στη συνέχεια, τρέξαμε τον αλγόριθμο Cheng και Church για τις κυτταρικές σειρές του καρκίνου του τραχήλου της μήτρας, για τα 1000 πρώτα γονίδια_από τα 33096 και εστίασαμε στα τρία πρώτα Biclusters (Cervical cancer).

Οι παρακάτω Εικόνες (5.2.α, 5.2.β, 5.2.γ) παρουσιάζουν το συνολικό πλήθος των γονιδίων και τη συμπεριφορά τους σε κάθε ένα από τα τρία Biclusters μετά τη πρώτη εκτέλεση του αλγορίθμου (Εικόνα 5.2.α), μετά τη δεύτερη εκτέλεση του αλγορίθμου (Εικόνες 5.2.β) και τέλος, μετά την Τρίτη εκτέλεση του αλγορίθμου (Εικόνες 5.2.γ).

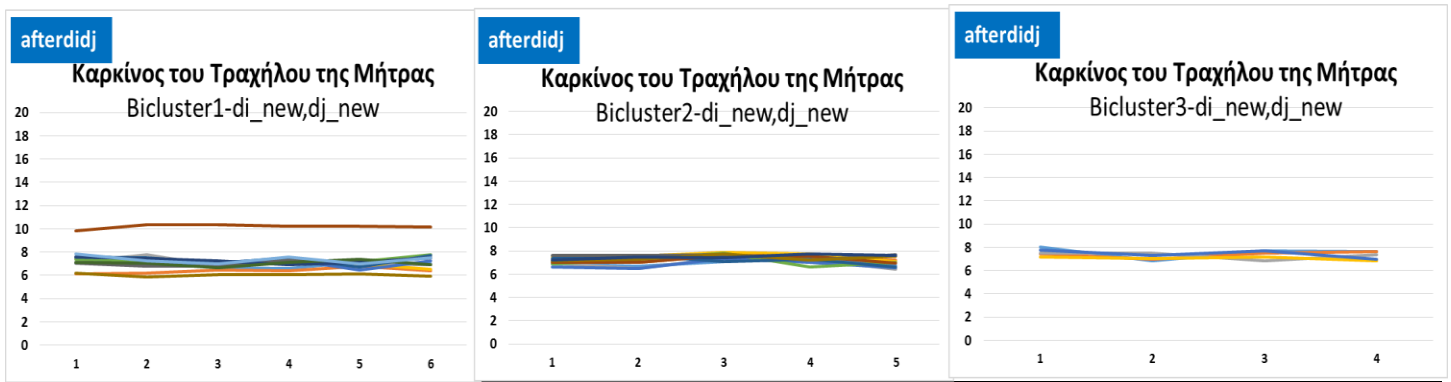
Εικόνα 5.2.α. - Πρώτη Εκτέλεση του Αλγορίθμου (di,dj) για 1000 Γονίδια του Καρκίνου του Τραχήλου της Μήτρας



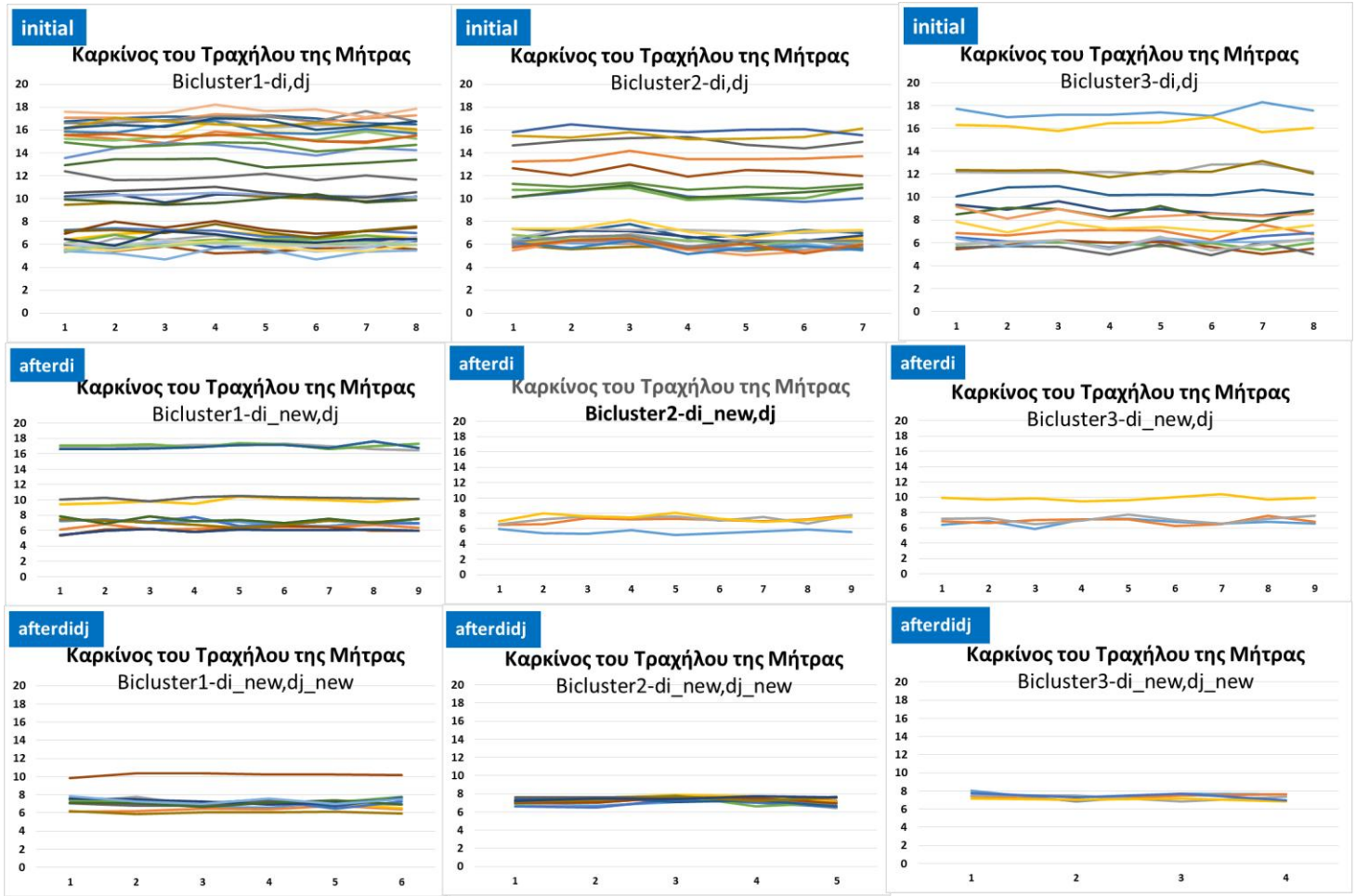
Εικόνα 5.2.β. - Δεύτερη Εκτέλεση του Αλγορίθμου (di_new,dj) για 1000 Γονίδια του Καρκίνου του Τραχήλου της Μήτρας



Εικόνα 5.2.γ. - Τρίτη Εκτέλεση του Αλγορίθμου (di_new,dj_new) για 1000 Γονίδια του Καρκίνου του Τραχήλου της Μήτρας



Εικόνα 5.3. - Συγκεντρωτική Παρουσίαση των Τριών εκτελέσεων του Αλγορίθμου Cheng και Church για 1000 Γονίδια από τις Κυτταρικές Σειρές Καρκίνου του Τραχήλου της Μήτρας

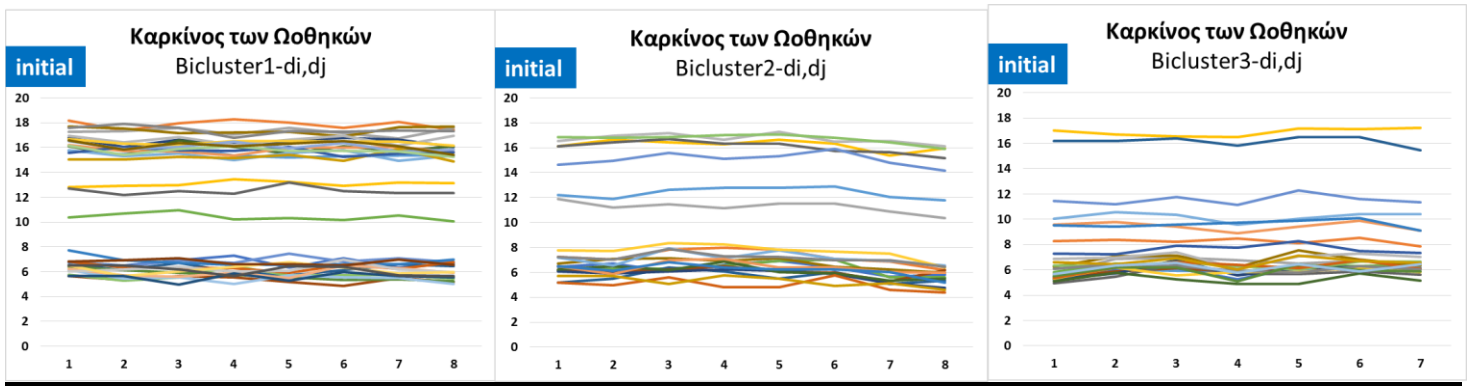


5.3. Αποτελέσματα biclustering από κυτταρικές σειρές καρκίνου των Ωοθηκών

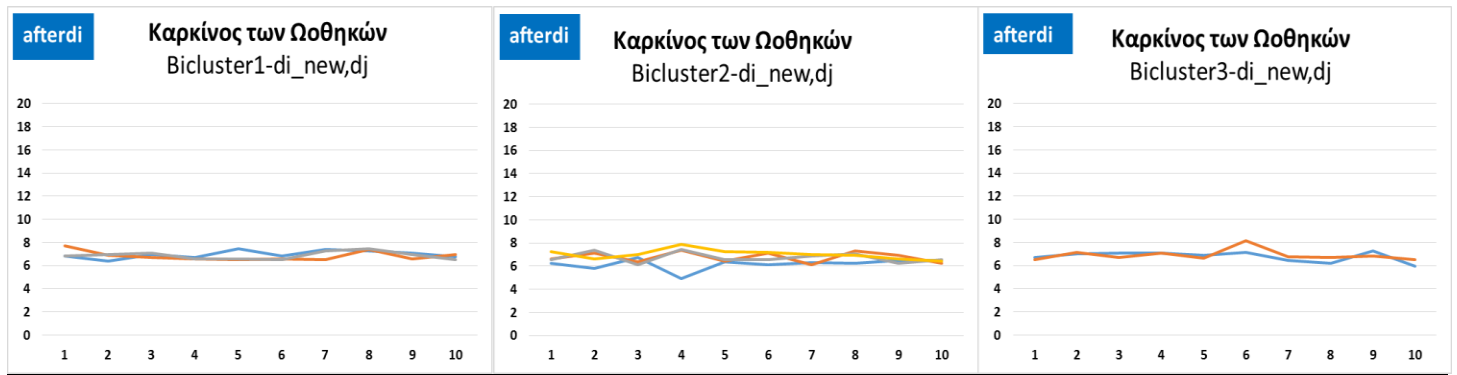
Στη συνέχεια, τρέξαμε τον αλγόριθμο Cheng και Church για τις κυτταρικές σειρές του καρκίνου των ωοθηκών, για τα 1000 πρώτα γονίδια από τα 33096, και εστίασαμε στα τρία πρώτα Biclusters.

Οι παρακάτω Εικόνες (5.3.α, 5.3.β, 5.23.γ) παρουσιάζουν το συνολικό πλήθος των γονιδίων και τη συμπεριφορά τους σε κάθε ένα από τα τρία Biclusters μετά τη πρώτη εκτέλεση του αλγορίθμου (Εικόνα 5.3.α), μετά τη δεύτερη εκτέλεση του αλγορίθμου (Εικόνες 5.3.β) και τέλος, μετά την Τρίτη εκτέλεση του αλγορίθμου (Εικόνες 5.3.γ).

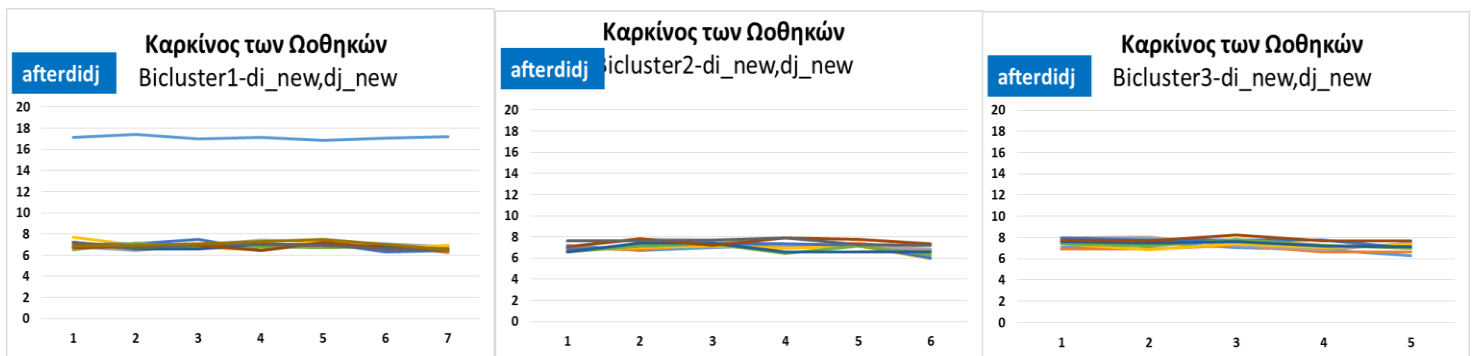
Εικόνα 5.3.α. - Πρώτη Εκτέλεση του Αλγορίθμου (di, dj) για 1000 Γονίδια του Καρκίνου των Ωοθηκών



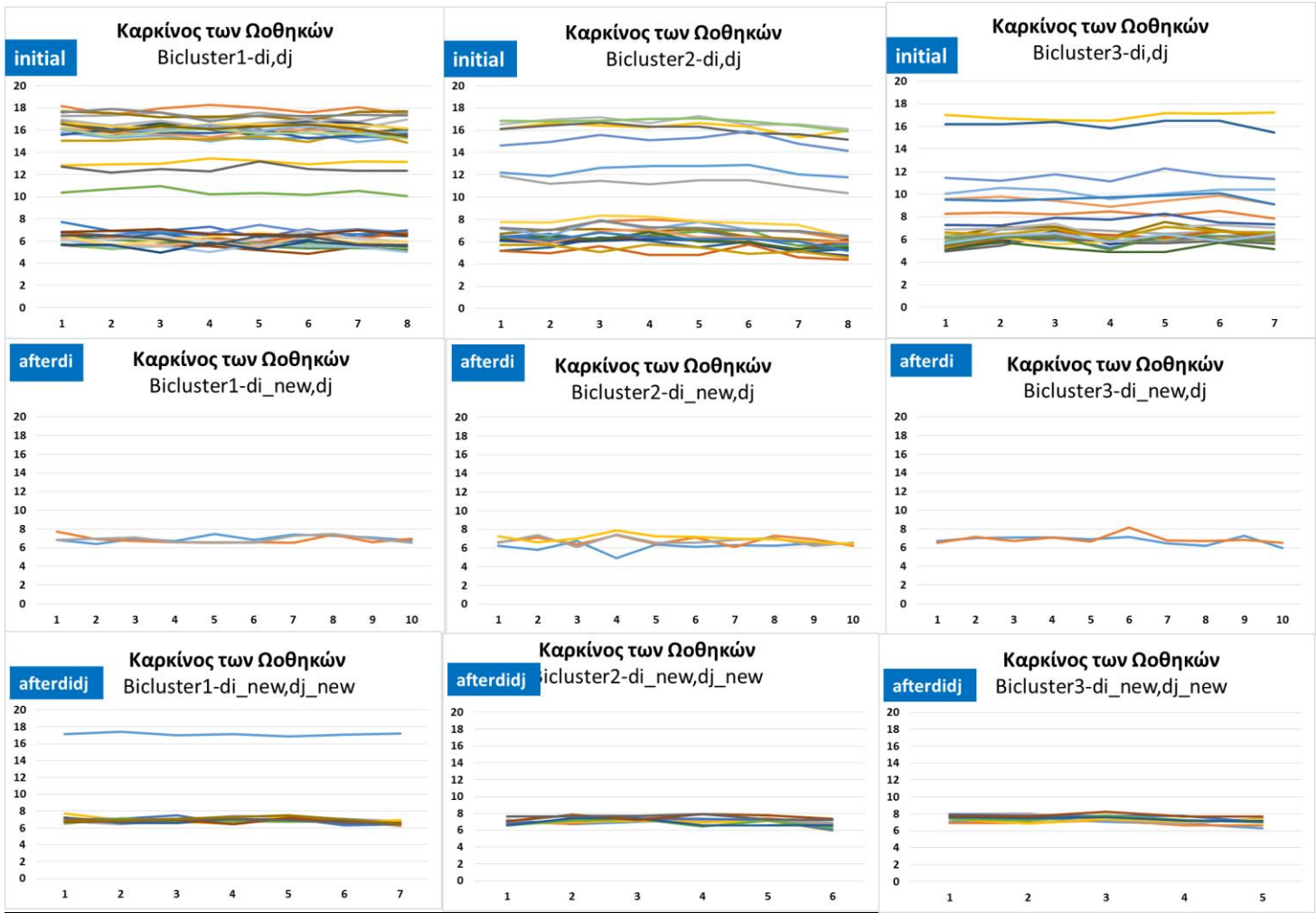
Εικόνα 5.3.β. - Δεύτερη Εκτέλεση του Αλγορίθμου (di_new, dj) για 1000 Γονίδια του Καρκίνου των Ωοθηκών



Εικόνα 5.3.γ. - Τρίτη Εκτέλεση του Αλγορίθμου (di_new, dj_new) για 1000 Γονίδια του Καρκίνου των Ωοθηκών



Εικόνα 5.4. - Συγκεντρωτική Παρουσίαση των Τριών εκτελέσεων του Αλγορίθμου Cheng και Church για 1000 Γονίδια από τις Κυτταρικές Σειρές Καρκίνου των Ωοθηκών

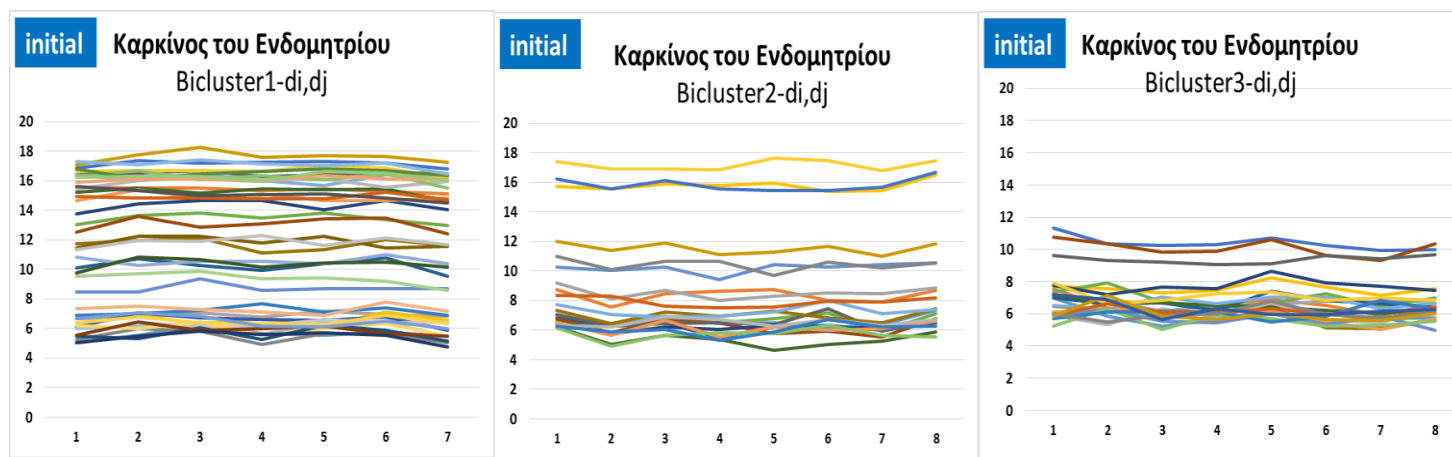


5.4. Αποτελέσματα biclustering από κυτταρικές σειρές καρκίνου του Ενδομητρίου

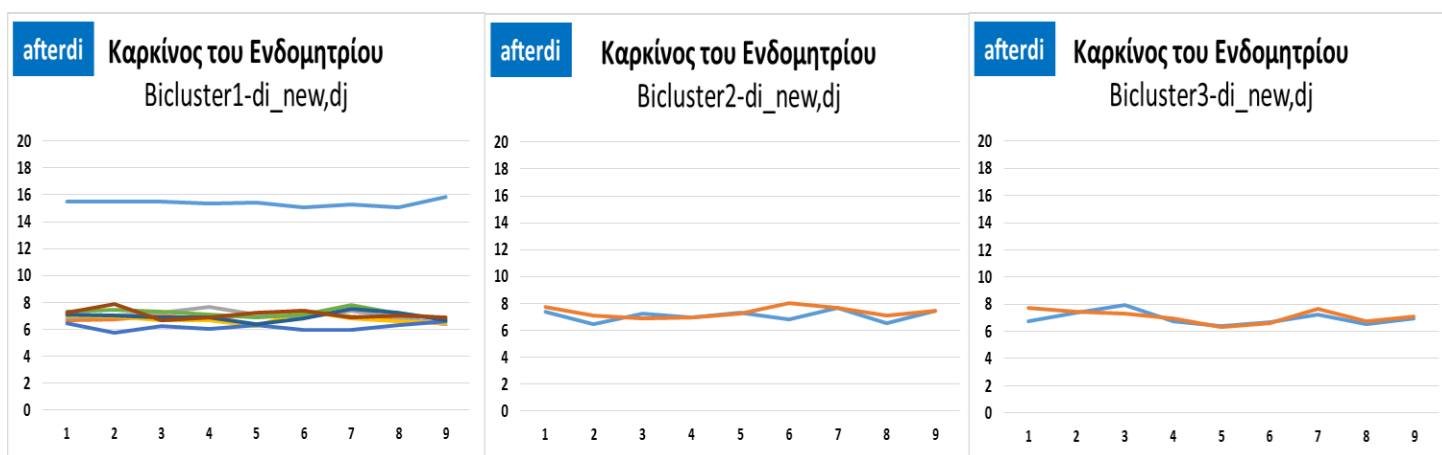
Στη συνέχεια, τρέξαμε τον αλγόριθμο Cheng και Church για τις κυτταρικές σειρές του καρκίνου του ενδομητρίου, για τα 1000 πρώτα γονίδια_από τα 33096 και εστίασαμε στα τρία πρώτα Biclusters.

Οι παρακάτω Εικόνες (5.4.α, 5.4.β, 5.4.γ) παρουσιάζουν το συνολικό πλήθος των γονιδίων και τη συμπεριφορά τους σε κάθε ένα από τα τρία Biclusters μετά τη πρώτη εκτέλεση του αλγορίθμου (Εικόνα 5.4.α), μετά τη δεύτερη εκτέλεση του αλγορίθμου (Εικόνες 5.4.β) και τέλος, μετά την Τρίτη εκτέλεση του αλγορίθμου (Εικόνες 5.4.γ).

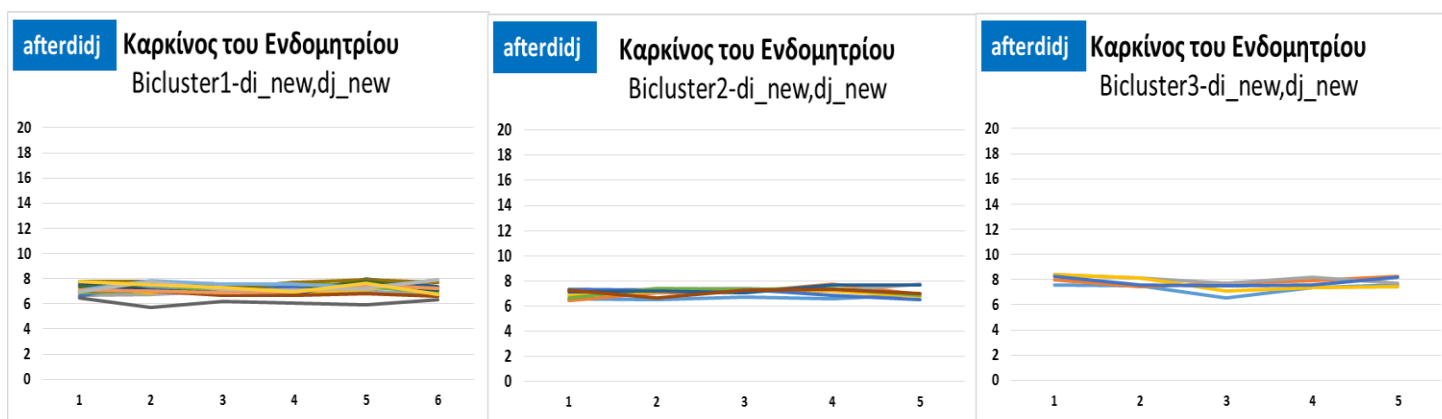
Εικόνα 5.4.α. - Πρώτη Εκτέλεση του Αλγορίθμου (di, dj) για 1000 Γονίδια του Καρκίνου του Ενδομητρίου



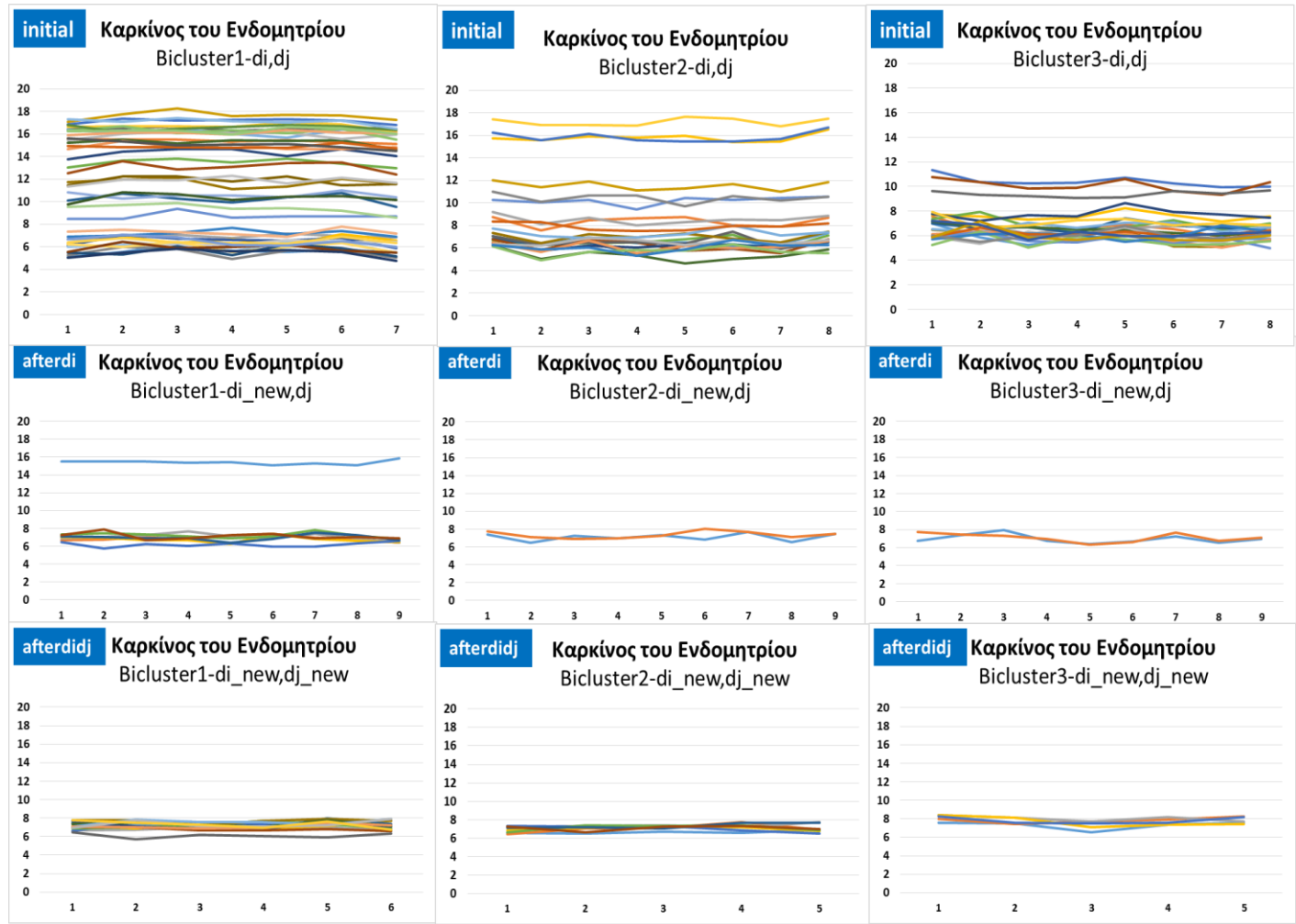
Εικόνα 5.4.β. - Δεύτερη Εκτέλεση του Αλγορίθμου (di_new, dj) για 1000 Γονίδια του Καρκίνου του Ενδομητρίου



Εικόνα 5.4.γ. - Τρίτη Εκτέλεση του Αλγορίθμου (di_new, dj_new) για 1000 Γονίδια του Καρκίνου του Ενδομητρίου



Εικόνα 5.5. - Συγκεντρωτική Παρουσίαση των Τριών εκτελέσεων του Αλγορίθμου Cheng και Church για 1000 Γονίδια από τις Κυτταρικές Σειρές Καρκίνου του Ενδομητρίου

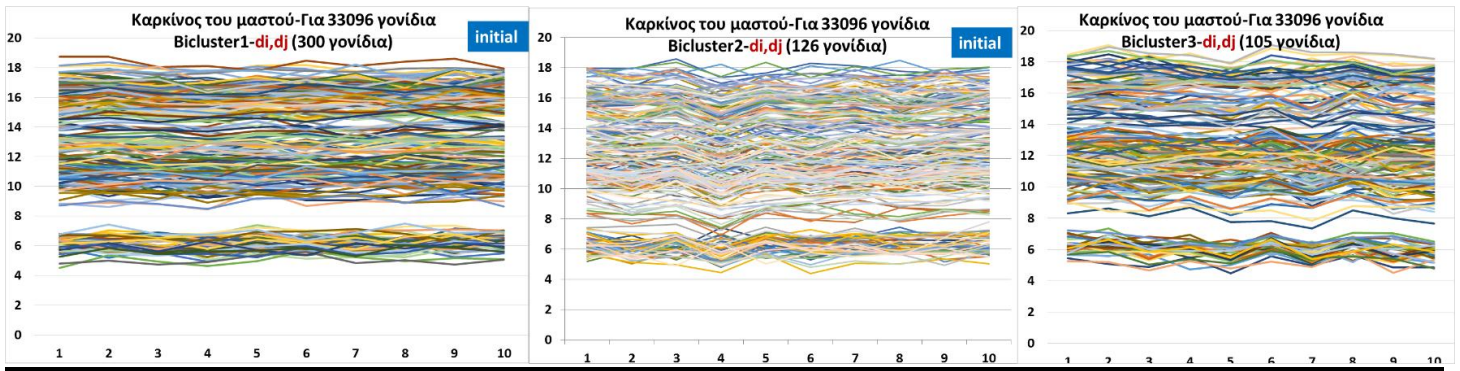


5.5. Αποτελέσματα Biclustering από Κυτταρικές Σειρές Καρκίνου του Μαστού για 33096 Γονίδια

Όπως έχει ειπωθεί και σε προηγούμενο κεφάλαιο, παρατηρήσαμε ότι κατά την εκτέλεση του αλγορίθμου για το σύνολο των 33096 γονιδίων, ο αριθμός τόσο των γονιδίων που ομαδοποιούνται σε ένα Bicluster, όσο και των Biclusters, είναι μεγάλος. Γι' αυτό το λόγο, προκειμένου να διακρίνουμε και να εξάγουμε συμπεράσματα για τη συμπεριφορά των γονιδίων εκτελέσαμε τον αλγόριθμο για τα 1000 πρώτα γονίδια από τα 33096 (Ενότητα 5.4.).

Σε αυτό το σημείο, τρέξαμε τον αλγόριθμο Cheng και Church για τα 33096 για τις κυτταρικές σειρές καρκίνου του μαστού, και εστίασαμε στα τρία πρώτα Biclusters. Οι παρακάτω Εικόνες (5.5.α, 5.5.β, 5.5.γ) παρουσιάζουν το συνολικό πλήθος των γονιδίων και τη συμπεριφορά τους σε κάθε ένα από τα τρία Biclusters μετά τη πρώτη εκτέλεση του αλγορίθμου (Εικόνα 5.8.α), μετά τη δεύτερη εκτέλεση του αλγορίθμου (Εικόνες 5.8.β) και τέλος, μετά την Τρίτη εκτέλεση του αλγορίθμου (Εικόνες 5.5.γ).

Εικόνα 5.5.α. - Πρώτη Εκτέλεση του Αλγορίθμου (di,dj) για 33096 Γονίδια του Καρκίνου του Μαστού

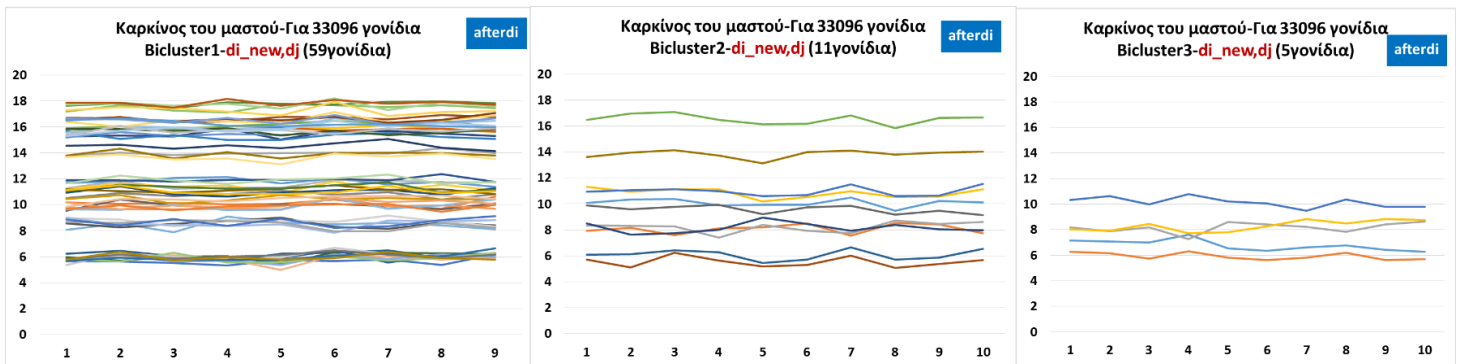


Αριστερό τμήμα: Το πρώτο διάγραμμα απεικονίζει 300 γονίδια. Λόγω του πλήθους των γονιδίων (300) που προκύπτουν στο Bicluster1 κατά την πρώτη εκτέλεση του αλγορίθμου. Στην εικόνα εμφανίζουμε ΜΟΝΟ τα 255 πρώτα γονίδια, εφόσον το πρόγραμμα του excel επιτρέπει την εμφάνιση έως 255 στοιχείων (Εικόνα 5.5.α - αριστερό τμήμα)

Μεσαίο τμήμα: Το δεύτερο διάγραμμα απεικονίζει 126 γονίδια. Εμφάνιση των 126 γονιδίων που προκύπτουν στο δεύτερο Bicluster κατά την πρώτη εκτέλεση του αλγορίθμου. (Εικόνα 5.5.α - μεσαίο τμήμα)

Δεξιό τμήμα: Το τρίτο διάγραμμα απεικονίζει 105 γονίδια. Εμφάνιση των 105 γονιδίων που προκύπτουν στο τρίτο Bicluster κατά την πρώτη εκτέλεση του αλγορίθμου. (Εικόνα 5.5.α - δεξιό τμήμα)

Εικόνα 5.5.β. - Δεύτερη Εκτέλεση του Αλγορίθμου (di_new,dj) για 33096 Γονίδια του Καρκίνου του Μαστού

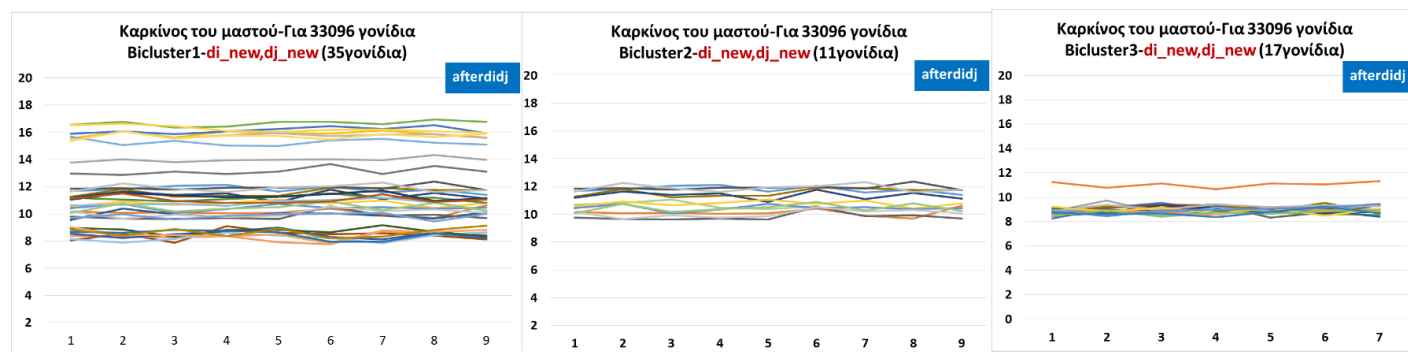


Αριστερό τμήμα: Το πρώτο διάγραμμα απεικονίζει 59 γονίδια. Εμφάνιση των 59 γονιδίων που προκύπτουν στο πρώτο Bicluster κατά τη δεύτερη εκτέλεση του αλγορίθμου. (Εικόνα 5.5.β - αριστερό τμήμα)

Μεσαίο τμήμα: Το δεύτερο διάγραμμα απεικονίζει 11 γονίδια. Εμφάνιση των 11 γονιδίων που προκύπτουν στο δεύτερο Bicluster κατά τη δεύτερη εκτέλεση του αλγορίθμου. (Εικόνα 5.5.β - μεσαίο τμήμα)

Δεξιό τμήμα: Το τρίτο διάγραμμα απεικονίζει 5 γονίδια. Εμφάνιση των 5 γονιδίων που προκύπτουν στο τρίτο Bicluster κατά τη δεύτερη εκτέλεση του αλγορίθμου. (Εικόνα 5.5.β- δεξιό τμήμα)

Εικόνα 5.5.γ. - Τρίτη Εκτέλεση του Αλγορίθμου (di_new, dj_new) για 33096 Γονίδια του Καρκίνου του Μαστού

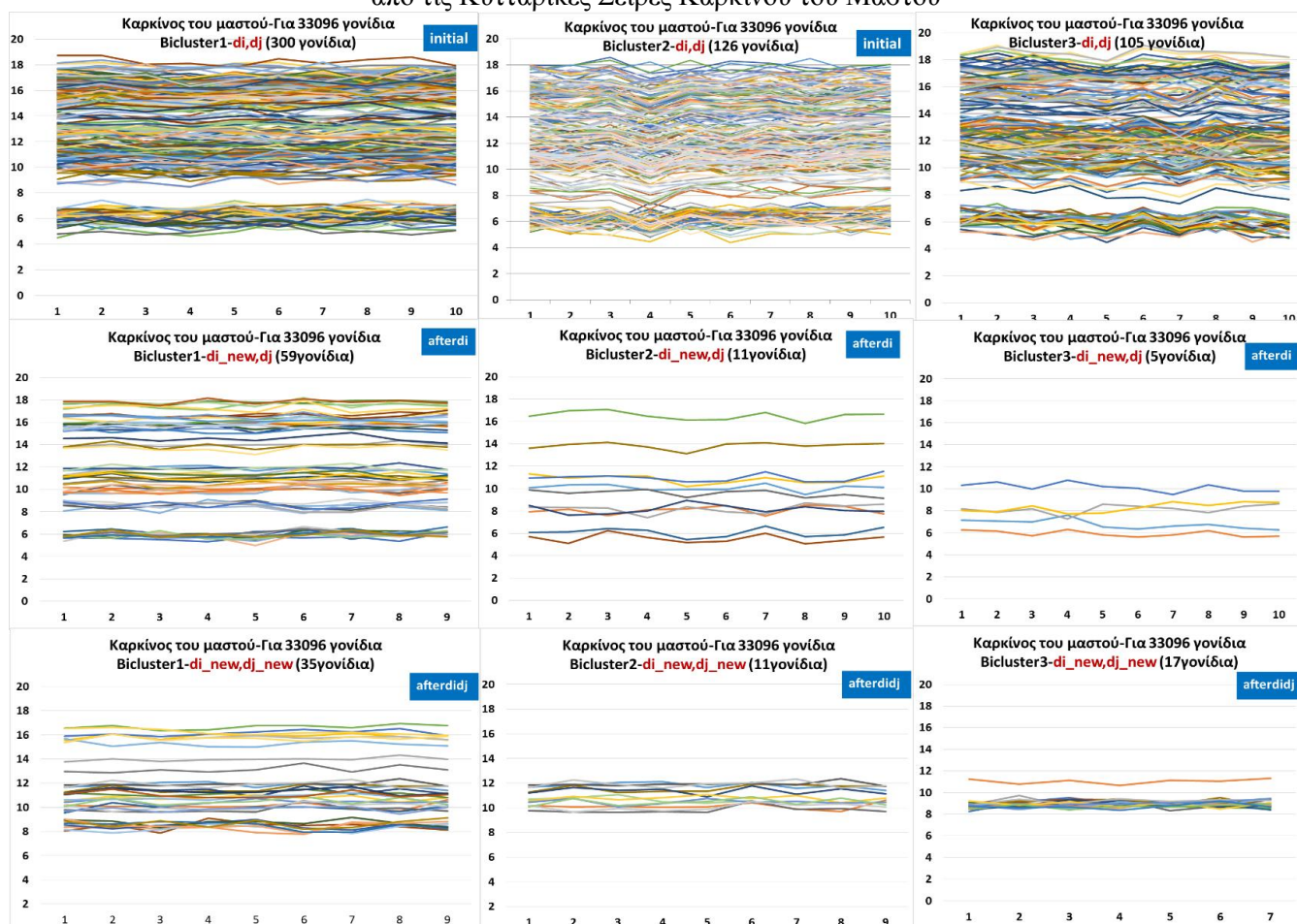


Αριστερό τμήμα: Το πρώτο διάγραμμα απεικονίζει 35 γονίδια. Εμφάνιση των 35 γονιδίων που προκύπτουν στο πρώτο Bicluster κατά την τρίτη εκτέλεση του αλγορίθμου. (Εικόνα 5.5.γ - αριστερό τμήμα)

Μεσαίο τμήμα: Το δεύτερο διάγραμμα απεικονίζει 11 γονίδια. Εμφάνιση των 11 γονιδίων που προκύπτουν στο δεύτερο Bicluster κατά την τρίτη εκτέλεση του αλγορίθμου. (Εικόνα 5.5.γ - μεσαίο τμήμα)

Δεξιό τμήμα: Το τρίτο διάγραμμα απεικονίζει 17 γονίδια. Εμφάνιση των 17 γονιδίων που προκύπτουν στο τρίτο Bicluster κατά την τρίτη εκτέλεση του αλγορίθμου. (Εικόνα 5.5.γ - δεξιό τμήμα)

Εικόνα 5.6. - Συγκεντρωτική Παρουσίαση των Τριών εκτελέσεων του Αλγορίθμου Cheng και Church για 33096 Γονίδια από τις Κυτταρικές Σειρές Καρκίνου του Μαστού



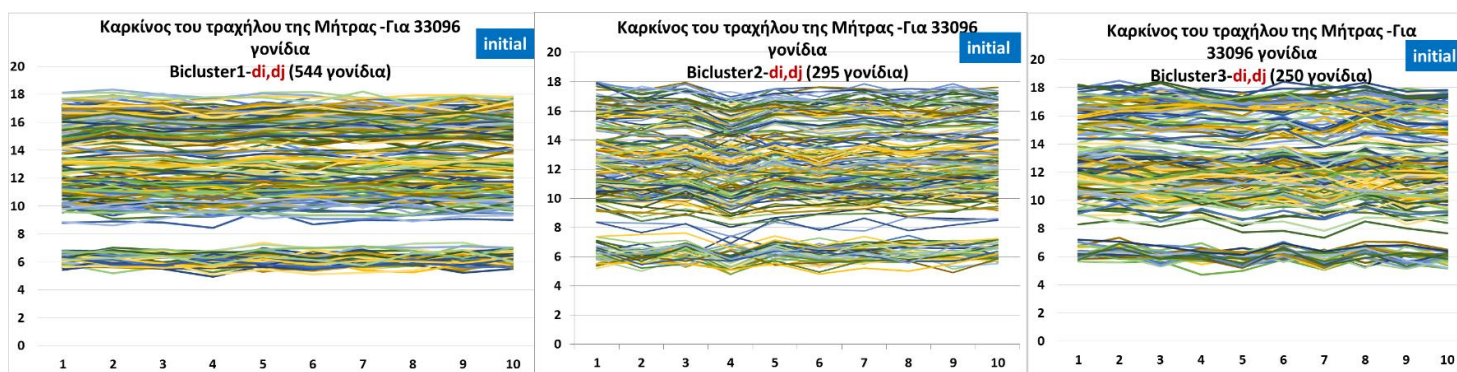
5.6. Αποτελέσματα Biclustering από Κυτταρικές Σειρές Καρκίνου του Τραχήλου της Μήτρας για 33096 Γονίδια

Όπως έχει ειπωθεί και σε προηγούμενο κεφάλαιο, παρατηρήσαμε ότι κατά την εκτέλεση του αλγορίθμου για το σύνολο των 33096 γονιδίων, ο όγκος των γονιδίων που ομαδοποιούνται σε ένα Bicluster και ο αριθμός των Biclusters είναι μεγάλος. Γι' αυτό το λόγο, προκειμένου να διακρίνουμε και να εξάγουμε συμπεράσματα για τη συμπεριφορά των γονιδίων εκτελέσαμε τον αλγόριθμο για τα 1000 πρώτα γονίδια από τα 33096, όπως παρουσιάστηκαν και τα αποτελέσματα αυτών των εκτελέσεων παραπάνω.

Σε αυτό το σημείο, τρέξαμε τον αλγόριθμο Cheng and Church για τα 33096 και εστίασαμε στα τρία πρώτα Biclusters από κυτταρικές σειρές καρκίνου του τραχήλου της Μήτρας .

Οι παρακάτω Εικόνες (5.6.α, 5.6.β, 5.6.γ) παρουσιάζουν το συνολικό πλήθος των γονιδίων και τη συμπεριφορά τους σε κάθε ένα από τα τρία Biclusters μετά τη πρώτη εκτέλεση του αλγορίθμου (Εικόνα 5.6.α), μετά τη δεύτερη εκτέλεση του αλγορίθμου (Εικόνες 5.6.β) και τέλος, μετά την Τρίτη εκτέλεση του αλγορίθμου (Εικόνες 5.6.γ).

Εικόνα 5.6.α. - Πρώτη Εκτέλεση του Αλγορίθμου (d_i, d_j) για 33096 Γονίδια του Καρκίνου του Τραχήλου της Μήτρας



Το πρώτο διάγραμμα λόγω του πλήθους των γονιδίων (544) που προκύπτουν στο Bicluster1 κατά την πρώτη εκτέλεση του αλγορίθμου, εμφανίζει ΜΟΝΟ τα 255 πρώτα γονίδια, εφόσον το πρόγραμμα του excel επιτρέπει την εμφάνιση έως 255 στοιχείων. (Εικόνα 5.6.α - αριστερό τμήμα)

Το δεύτερο διάγραμμα λόγω του πλήθους των γονιδίων (295) που προκύπτουν στο Bicluster2 κατά την πρώτη εκτέλεση του αλγορίθμου, εμφανίζει ΜΟΝΟ τα 255 πρώτα γονίδια, εφόσον το πρόγραμμα του excel επιτρέπει την εμφάνιση έως 255 στοιχείων. (Εικόνα 5.6.α - μεσαίο τμήμα)

Το τρίτο διάγραμμα απεικονίζει 250 γονίδια. Εμφάνιση των 250 γονιδίων που προκύπτουν στο τρίτο Bicluster κατά την πρώτη εκτέλεση του αλγορίθμου. (Εικόνα 5.6.α – δεξί τμήμα)

Εικόνα 5.6.β. - Δεύτερη Εκτέλεση του Αλγορίθμου (di_new, dj) για 33096 Γονίδια του Καρκίνου του Τραχήλου της Μήτρας

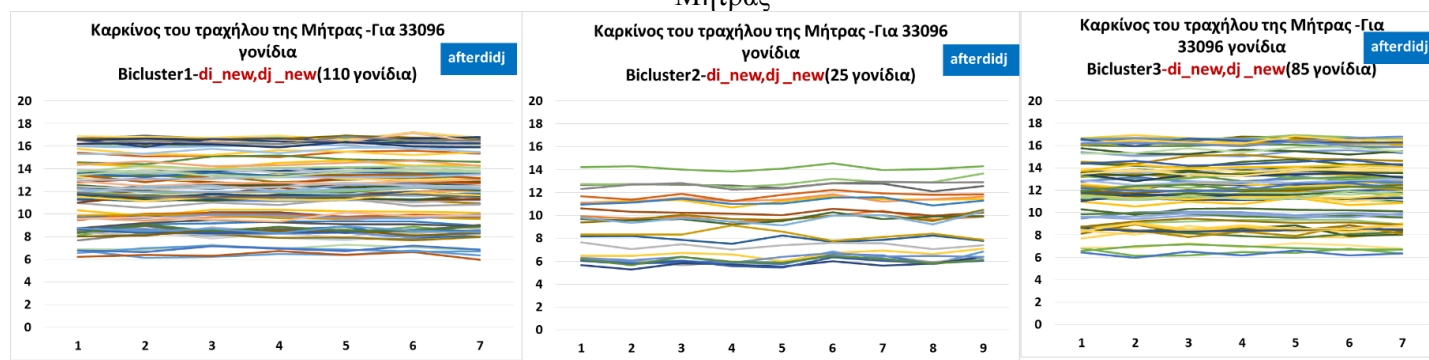


Το πρώτο διάγραμμα απεικονίζει 207 γονίδια. Εμφάνιση των 207 γονιδίων που προκύπτουν στο πρώτο Bicluster κατά τη δεύτερη εκτέλεση του αλγορίθμου. (Εικόνα 5.6.β - αριστερό τμήμα)

Το δεύτερο διάγραμμα απεικονίζει 75 γονίδια. Εμφάνιση των 75 γονιδίων που προκύπτουν στο δεύτερο Bicluster κατά τη δεύτερη εκτέλεση του αλγορίθμου. (Εικόνα 5.6.β – μεσαίο τμήμα)

Το τρίτο διάγραμμα απεικονίζει 125 γονίδια. Εμφάνιση των 125 γονιδίων που προκύπτουν στο τρίτο Bicluster κατά τη δεύτερη εκτέλεση του αλγορίθμου. (Εικόνα 5.6.β - δεξί τμήμα)

Εικόνα 5.6.γ. - Τρίτη Εκτέλεση του Αλγορίθμου (di_new, dj_new) για 33096 Γονίδια του Καρκίνου του Τραχήλου της Μήτρας

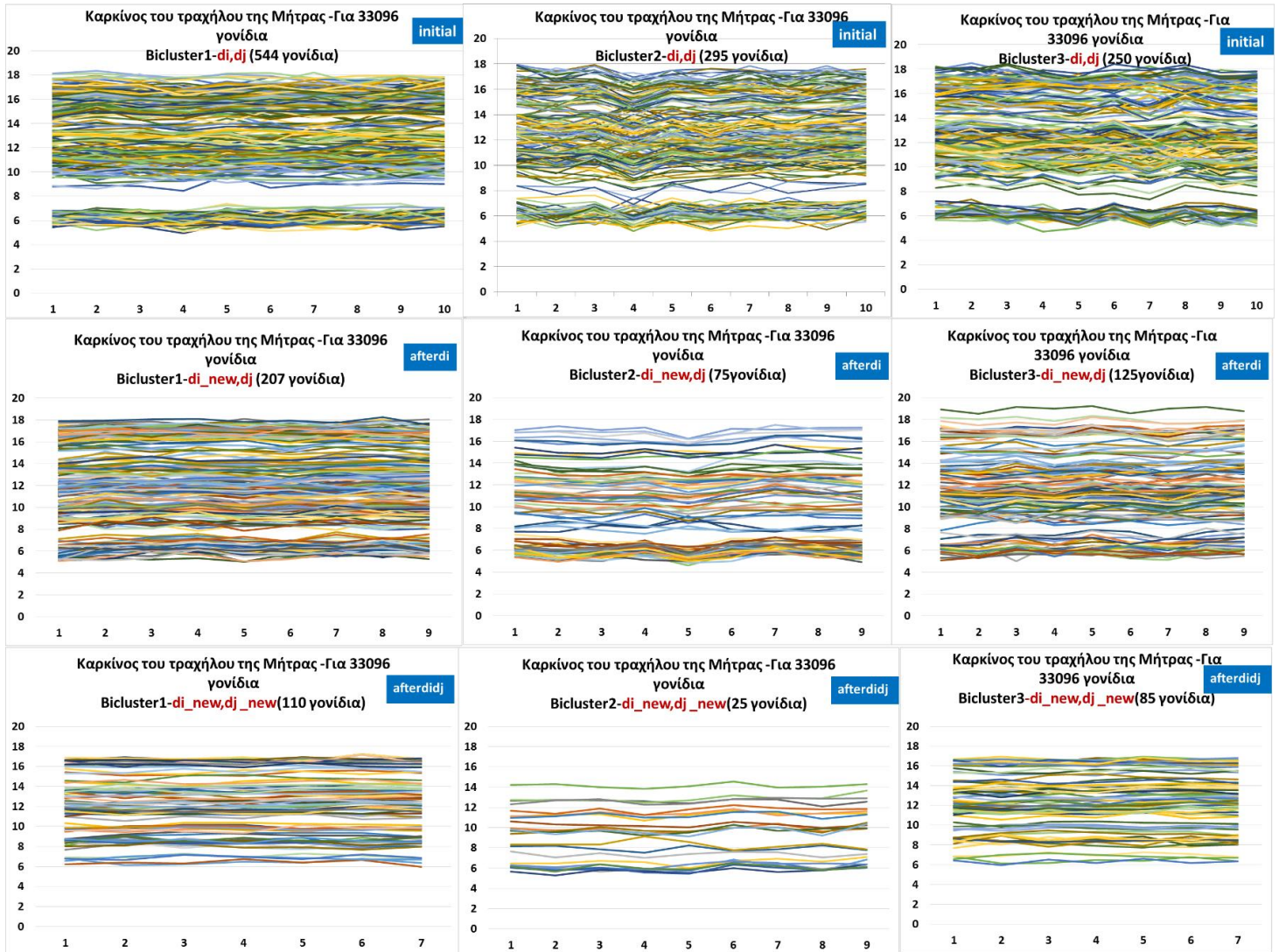


Το πρώτο διάγραμμα απεικονίζει 110 γονίδια. Εμφάνιση των 110 γονιδίων που προκύπτουν στο πρώτο Bicluster κατά την τρίτη εκτέλεση του αλγορίθμου. (Εικόνα 5.6.γ - αριστερό τμήμα)

Το δεύτερο διάγραμμα απεικονίζει 25 γονίδια. Εμφάνιση των 25 γονιδίων που προκύπτουν στο δεύτερο Bicluster κατά την τρίτη εκτέλεση του αλγορίθμου. (Εικόνα 5.6.γ - μεσαίο τμήμα)

Το τρίτο διάγραμμα απεικονίζει 85 γονίδια. Εμφάνιση των 85 γονιδίων που προκύπτουν στο τρίτο Bicluster κατά την τρίτη εκτέλεση του αλγορίθμου. (Εικόνα 5.6.γ - δεξί τμήμα)

Εικόνα 5.7. - Συγκεντρωτική Παρουσίαση των Τριών εκτελέσεων του Αλγορίθμου Cheng and Church για 33096 Γονίδια από τις Κυτταρικές Σειρές Καρκίνου του Τραχήλου της Μήτρας



6

Αξιολόγηση Των Αποτελεσμάτων

Η αξιολόγηση των αποτελεσμάτων ακολούθησε την ανάλυση της διπλής κατηγοριοποίησης σύμφωνα με την προτεινόμενη μεθοδολογία, που όπως ήδη αναφέρθηκε, εφαρμόστηκε σ' ένα σύνολο δεδομένων που περιείχε τις τιμές έκφρασης 33096 γονιδίων για 38 κυτταρικές σειρές τεσσάρων καρκινικών τύπων (μαστού, τραχήλου της μήτρας, ενδομητρίου, και ωοθηκών).[57]

Τα προβλήματα που τέθηκαν προς αντιμετώπιση, όπως αναφέρθηκαν, ήταν:

α) ο εξαιρετικά μεγάλος αριθμός των ομάδων διπλής κατηγοριοποίησης που προκύπτουν κατά την αρχική εκτέλεση του αλγορίθμου για τα 33096 γονίδια για τον καθένα από τους τέσσερις καρκινικούς τύπους ξεχωριστά, παραδείγματος χάριν το πλήθος των Biclusters για τον καρκίνο του μαστού ανερχόταν σε 2411, (κεφάλαιο 5),

β) το μεγάλο πλήθος των γονιδίων που ομαδοποιούνται σε καθεμιά ομάδα διπλής κατηγοριοποίησης κατά την αρχική εκτέλεση του αλγορίθμου για τα 33096 γονίδια και

γ) η ανομοιομορφία μεταξύ των γονιδίων που συγκεντρώνονται σε μία ομάδα διπλής κατηγοριοποίησης, δηλαδή τα γονίδια διέφεραν στην τάση παρουσιάζοντας μεγάλες διακυμάνσεις, αλλά και στη μεταξύ τους απόσταση καλύπτοντας ένα μεγάλο εύρος τιμών.

Η εφαρμογή της μεθόδου διπλής κατηγοριοποίησης στο σύνολο δεδομένων καρκινικών κυτταρικών σειρών, είχε ως στόχο τη μελέτη ενός μικρού αριθμού ομάδων διπλής κατηγοριοποίησης και κατ' επέκταση τη μελέτη ενός μικρού αριθμού γονιδίων ώστε να είναι πιο εύκολη η παρατήρηση της συμπεριφοράς τους και μετέπειτα η εξαγωγή συμπερασμάτων με βάση αυτή. Το αποτέλεσμα αυτής της προσέγγισης ήταν η ανάδειξη των 3 πρώτων ομάδων διπλής κατηγοριοποίησης για τα 1000 πρώτα γονίδια από τα 33096 γονίδια για τις κυτταρικές σειρές των τεσσάρων τύπων καρκίνου (του μαστού, του τραχήλου της μήτρας, των ωοθηκών και του ενδομητρίου).

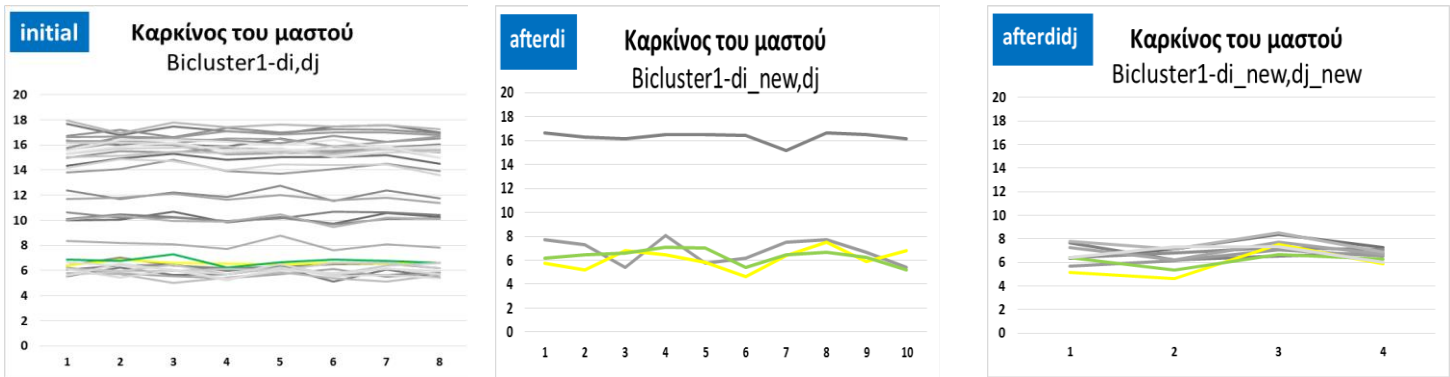
Αφού εκτελέστηκε ο αλγόριθμος Cheng και Church στην αρχική του μορφή (1η εκτέλεση) για κάθε καρκίνο, στη συνέχεια εκτελέστηκε ο αλγόριθμος C&C με αλλαγή στον τρόπο υπολογισμού του di (2η εκτέλεση). Αυτή η εκτέλεση δίνει βαρύτητα στην ομοιομορφία μεταξύ των γονιδίων με βάση την τάση τους, όπως καθορίζεται από το μέτρο που αλλάζουμε. Τέλος, αλλάζοντας και το dj (3η εκτέλεση) για κάθε τύπο καρκίνου επιδιώκουμε την μείωση της απόστασης μεταξύ των γονιδίων που βρίσκονται στο ίδιο Biclusters.

Στον παρακάτω συγκεντρωτικό Πίνακα 6.1, όπως έχει παρουσιαστεί και στο κεφάλαιο 5, παρουσιάζεται το συνολικό πλήθος των γονιδίων (probeid) που ομαδοποιείται σε καθένα από τα τρία Biclusters στον καρκίνο του μαστού για καθεμιά από τις τρεις εκτελέσεις του αλγορίθμου για τα 1000 πρώτα γονίδια. Με χρωματιστό

πλαίσιο παρουσιάζονται τα γονίδια (probeid) τα οποία μεταπηδούν από το ένα Bicluster σε ένα από τα άλλα δύο Biclusters κατά την εφαρμογή των τριών εκτελέσεων του αλγορίθμου. Επίσης, με έντονη πλάγια γραφή και χρωματιστό πλαίσιο παρουσιάζονται τα γονίδια τα οποία εμφανίζονται στο ίδιο Bicluster κατά τη διάρκεια δύο τουλάχιστον διαφορετικών εκτελέσεων του αλγορίθμου.

Πίνακας 6.1. - Παρουσίαση τριών πρώτων Biclusters του καρκίνου του μαστού - για καθεμιά από τις τρεις εκτελέσεις του αλγορίθμου στα 1000 πρώτα γονίδια. Τα γονίδια δίνονται με τη μορφή κωδικών (probeid)

Πρώτη εκτέλεση του αλγορίθμου(di,dj)			Δεύτερη εκτέλεση του αλγορίθμου(di_new,dj)			Τρίτη εκτέλεση του αλγορίθμου(di_new,dj_new)		
Bicluster1	Bicluster2	Bicluster3	Bicluster1	Bicluster2	Bicluster3	Bicluster1	Bicluster2	Bicluster3
226368	164432	217298	139301	162744	231349	219524	142719	151787
213726	151787	165735	214589	181532	140490	231349	181532	162744
224284	223314	223080	180501	201489		111270	101235	125603
224286	161079	225083	214717	236216		140490	201489	132754
199782	140336	136093		140456		214589	146732	115325
173890	127051	156066		170969		190994	109288	
101695	137543	180730				180501	234158	
147346	160701	154584				214717		
138543	217593	140229				163565		
186673	136075	173299						
174944	230085	162744						
155326	123356	188198						
236908	151430	109012						
155784	178114	200313						
149734	182060	166147						
228186	149009	120825						
139301	216613	183840						
180707	194406	174851						
128561	121364	182708						
181532	136278	123369						
182504		109258						
139732		140456						
116808		202153						
180501		189928						
214717		209910						
236216		163565						
215964		104468						
128192		201366						
170969								
229775								
161284								
139151								
108780								
109210								
104969								
118956								



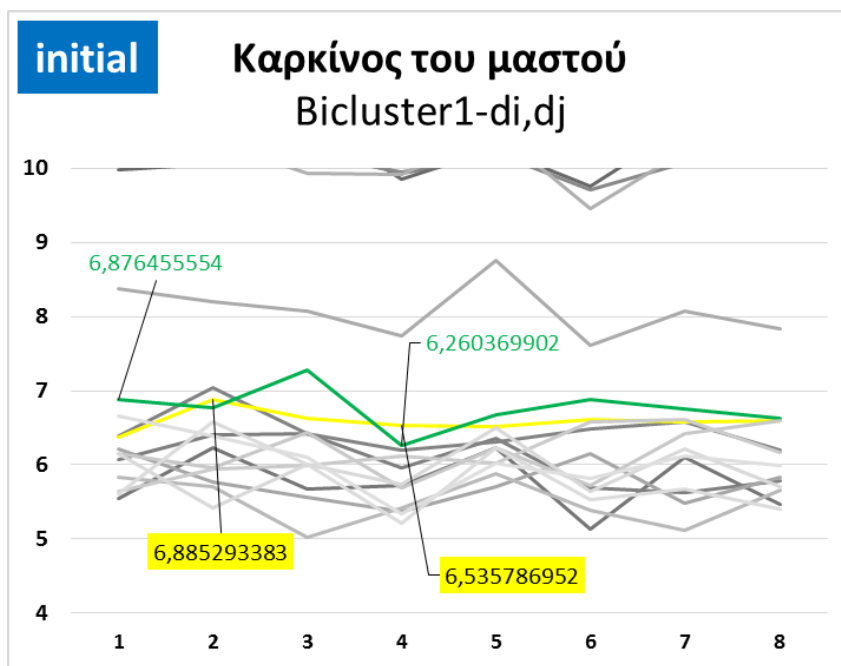
Εικόνα 6.1. - Συγκεντρωτική Παρουσίαση των Τριών εκτελέσεων του Αλγορίθμου Cheng και Church για 1000 Γονίδια από τις Κυτταρικές Σειρές Καρκίνου του Μαστού (Bicluster1), εστιάζοντας στη συμπεριφορά των γονιδίων με probeid 214717 και 180501.

Για παράδειγμα, στον Πίνακα 6.1 απεικονίζεται η συμπεριφορά του γονιδίου με probeid 214717 και του γονιδίου με probeid 180501. Τα γονίδια 214717 και 180501 τυγχάνει να παραμένουν στο Bicluster1 και στις τρεις εκτελέσεις. Δηλαδή, δεν χρειάστηκε να μετακινηθούν σε άλλο Bicluster, ενώ ταυτόχρονα άλλα γονίδια που «ταίριαζαν» με τα δύο γονίδια που εξετάζουμε, στην τάση και στην τιμή, μετακινήθηκαν στο Bicluster 1.

Στα παρακάτω διαγράμματα (Εικόνες 6.1.α, 6.1.β, 6.2.α, 6.2.β, 6.3.α, 6.3.β, 6.4.α, 6.4.β, 6.5.α, 6.5.β, 6.6.α, 6.6.β), η κλίμακα του άξονα y έχει προσαρμοστεί από την τιμή 4 μέχρι την τιμή 10, ώστε να είναι πιο ευδιάκριτη η συμπεριφορά των γονιδίων όσον αφορά το εύρος της μεταξύ τους απόστασης και τη διακύμανσή τους (κορυφές). Η αλλαγή της κλίμακας έχει ως αποτέλεσμα την παρουσίαση ενός υποσυνόλου γονιδίων.

Επίσης, στις Εικόνες 6.1.α, 6.2.α, 6.3.α, 6.4.α, 6.4.β, 6.5.α, και 6.6.α απεικονίζονται οι κυτταρικές σειρές που συμμετέχουν σε κάθε Bicluster, με σκοπό την ανάδειξη του φαινομένου όπου κυτταρικές σειρές προστίθενται ή αφαιρούνται κατά την εφαρμογή των τριών διαφορετικών εκτελέσεων.

Παρατήρηση της συμπεριφοράς των γονιδίων με κωδικό «probeid» 241717 και 180501

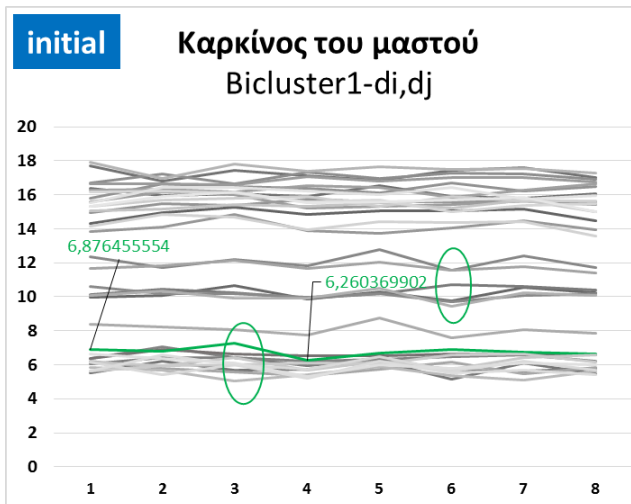


Εικόνα 6.1.α. - Πρώτη Εκτέλεση του Αλγορίθμου (di,dj) για 1000 Γονίδια του Καρκίνου του Μαστού
Ανάδειξη των μεγαλύτερων και μικρότερων τιμών των Καρκινικών Κυτταρικών Σειρών που Συμμετέχουν

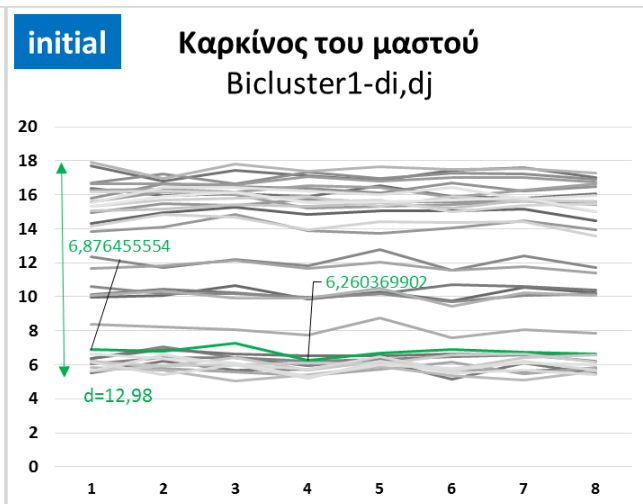
Σε πρώτη φάση, στην Εικόνα 6.1.α παρατηρείται η απόσταση των γονιδίων και η διακύμανση των γονιδίων στο Bicluster 1, κατά την πρώτη εκτέλεση του αλγορίθμου C&C για τα 1000 γονίδια του καρκίνου του μαστού.

Επίσης, στην παραπάνω Εικόνα 6.1.α απεικονίζονται οι κυτταρικές σειρές που περιλαμβάνονται στο Bicluster1 μετά την πρώτη εκτέλεση του αλγορίθμου και οι τιμές των γονιδίων (π.χ. τιμές γονιδίου με probeid 214717) σε κάθε καρκινική κυτταρική σειρά. Οι 10 κυτταρικές σειρές που απαρτίζουν τον πίνακα (των δεδομένων) στον καρκίνο του μαστού, μετά την πρώτη εκτέλεση μειώνονται στις 8.

Απομονώνοντας το γονίδιο με κωδικό 214717 παρατηρούμε: 1) την ανομοιομορφία μεταξύ των γονιδίων που συγκεντρώνονται στο Bicluster1, και 2) την μεγάλη απόσταση που παρουσιάζουν τα γονίδια μεταξύ τους. Συγκεκριμένα, τα γονίδια διαφέρουν στην τάση και παρουσιάζουν μεγάλες διακυμάνσεις (Εικόνα 6.1.α), ενώ το εύρος της απόστασης των γονιδίων στο Bicluster1 μετά την πρώτη εκτέλεση είναι $d=12,98$ όπως φαίνεται στην Εικόνα 6.1.γ, στην οποία παρουσιάζεται το συνολικό πλήθος των εξαγόμενων γονιδίων που ομαδοποιούνται στο Bicluster1.

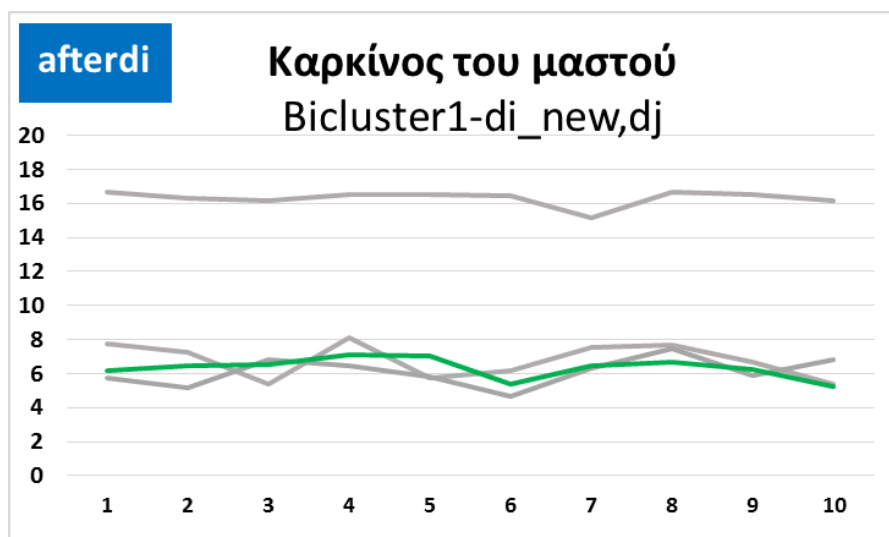


Εικόνα 6.1.β. - Πρώτη Εκτέλεση του Αλγορίθμου (di,dj) για 1000 Γονίδια του Καρκίνου του Μαστού Αναδεικνύοντας την Ανομοιομορφία μεταξύ των γονιδίων



Εικόνα 6.1.γ. - Πρώτη Εκτέλεση του Αλγορίθμου (di,dj) για 1000 Γονίδια του Καρκίνου του Μαστού Αναδεικνύοντας το Εύρος της Απόστασης μεταξύ των γονιδίων

Στη συνέχεια, κατά τη δεύτερη εκτέλεση, δηλαδή μετά την τροποποίηση του di (di_{new},dj) το γονίδιο το οποίο “παρακολουθείται”, παραμένει στο Bicluster1, όπου συσπειρώνεται με άλλα γονίδια με παρόμοια τάση. Αυτό επιτυγχάνεται λόγω των αλλαγών που έχουμε κάνει στο μέτρο.

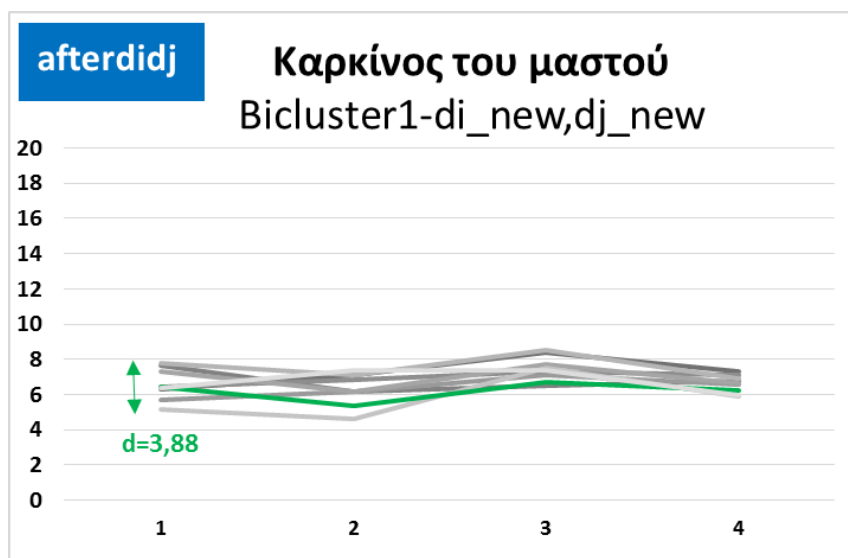


Εικόνα 6.2. - Δεύτερη Εκτέλεση του Αλγορίθμου (di_{new},dj) για 1000 Γονίδια του Καρκίνου του Μαστού Ανάδειξη των Καρκινικών Κυτταρικών Σειρών που Συμμετέχουν

Επίσης, στην παραπάνω Εικόνα 6.2.α απεικονίζονται οι κυτταρικές σειρές που περιλαμβάνονται στο Biclust_{er1} μετά τη δεύτερη εκτέλεση του αλγορίθμου και οι τιμές των γονιδίων (π.χ. τιμές γονιδίου με probeid 214717) σε κάθε καρκινική κυτταρική σειρά. Από τις 6 κυτταρικές σειρές που παρέμειναν μετά την πρώτη εκτέλεση του αλγορίθμου, μετά το πέρας της δεύτερης εκτέλεσης του αλγορίθμου παρατηρούμε ότι συμμετέχουν και οι 10 κυτταρικές σειρές .

Τέλος, μπορούμε να διακρίνουμε τη μείωση των γονιδίων από τα 30 που είχαμε στην πρώτη εκτέλεση, στα 4 γονίδια.

Ακολουθώς, η παρακάτω Εικόνα 6.3.α απεικονίζει τις κυτταρικές σειρές που περιλαμβάνονται στο Biclust_{er1} μετά την τρίτη εκτέλεση του αλγορίθμου και οι τιμές των γονιδίων (π.χ. τιμές γονιδίου με probeid 214717) σε κάθε καρκινική κυτταρική σειρά. Ωστόσο, παρατηρείται το γεγονός ότι έχουν παραμείνει μόνο τέσσερις κυτταρικές σειρές.

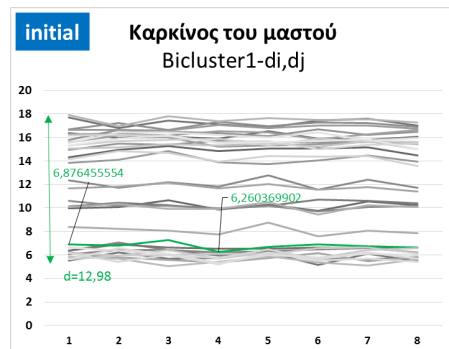


Εικόνα 6.3.α. - Τρίτη Εκτέλεση του Αλγορίθμου (di_new, dj_new) για 1000 Γονίδια του Καρκίνου του Μαστού
Ανάδειξη των Καρκινικών Κυτταρικών Σειρών που Συμμετέχουν

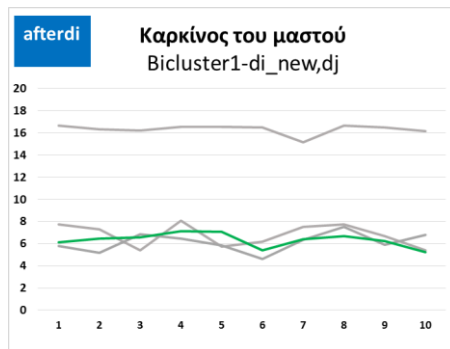
Επομένως, απομονώνοντας το γονίδιο με κωδικό 214717 παρατηρούμε την ομοιομορφία πλέον μεταξύ των γονιδίων που συγκεντρώνονται στο Biclust_{er1}, δηλαδή τα γονίδια παρουσιάζουν παρόμοια τάση και μικρές διακυμάνσεις (Εικόνα 6.3.α), αλλά και η απόσταση που παρουσιάζουν τα γονίδια μεταξύ τους είναι πλέον μικρή. Συγκεκριμένα, το εύρος των γονιδίων στο Biclust_{er1} μετά την πρώτη εκτέλεση ήταν $d=12,98$ όπως φαίνεται στην Εικόνα 6.3.α, ενώ μετά τις τρεις εκτελέσεις, το εύρος των γονιδίων είναι $d'=3,88$.

Στα διαγράμματα (Εικόνες 6.3.γ, 6.3.δ, 6.3.ε), παρατηρούμε τη συμπεριφορά των γονιδίων του Biclust_{er1} και στις τρεις εκτελέσεις, με κλίμακα του άξονα y από την τιμή 0 μέχρι την τιμή 20, ώστε να έχουμε την πλήρη εικόνα των γονιδίων.

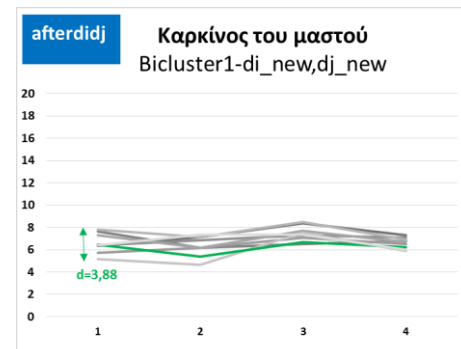
Τα διαγράμματα (Εικόνες 6.3.γ, 6.3.δ, 6.3.ε) απεικονίζουν τη συμπεριφορά των γονιδίων κατά τις τρεις εκτελέσεις του αλγορίθμου για τα 1000 πρώτα Γονίδια του Καρκίνου του Μαστού για το Bicluster1 .



Εικόνα 6.3.β. - Πρώτη Εκτέλεση του Αλγορίθμου (di,dj) για 1000 Γονίδια του Καρκίνου του Μαστού

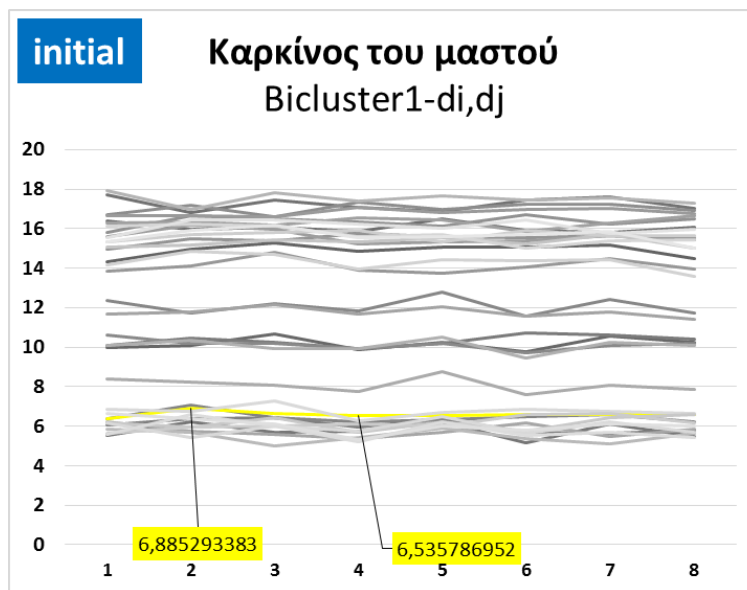


Εικόνα 6.3.γ. - Δεύτερη Εκτέλεση του Αλγορίθμου (di_new,dj) για 1000 Γονίδια του Καρκίνου του Μαστού

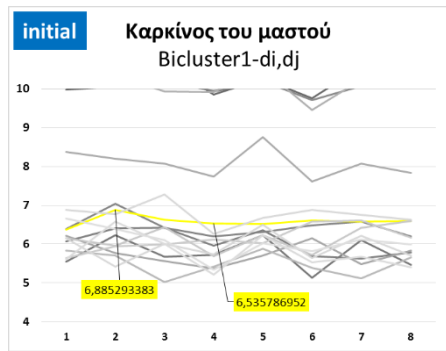


Εικόνα 6.3.δ. - Τρίτη Εκτέλεση του Αλγορίθμου (di_new,dj_new) για 1000 Γονίδια του Καρκίνου του Μαστού

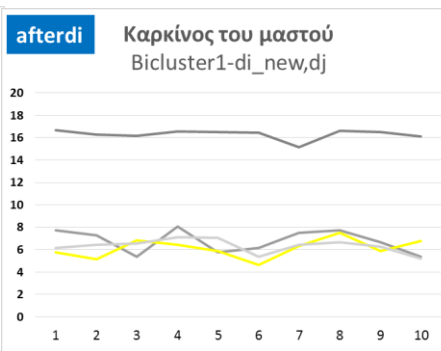
Παρατήρηση της συμπεριφοράς του γονιδίου με αριθμό probeid 180501



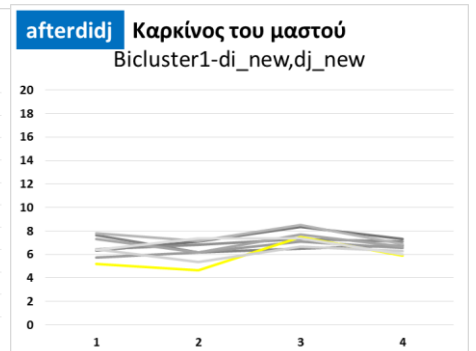
Εικόνα 6.4.α. - Πρώτη Εκτέλεση του Αλγορίθμου (didj) για 1000 Γονίδια του Καρκίνου του Μαστού Ανάδειξη των Καρκινικών Κυτταρικών Σειρών που Συμμετέχουν



Εικόνα 6.4.β. - Πρώτη Εκτέλεση του Αλγορίθμου (didj) για 1000 Γονίδια του Καρκίνου του Μαστού Αναδεικνύοντας το γονίδιο 180501



Εικόνα 6.4.γ. - Δεύτερη Εκτέλεση του Αλγορίθμου (di_new,dj) για 1000 Γονίδια του Καρκίνου του Μαστού Αναδεικνύοντας το γονίδιο 180501



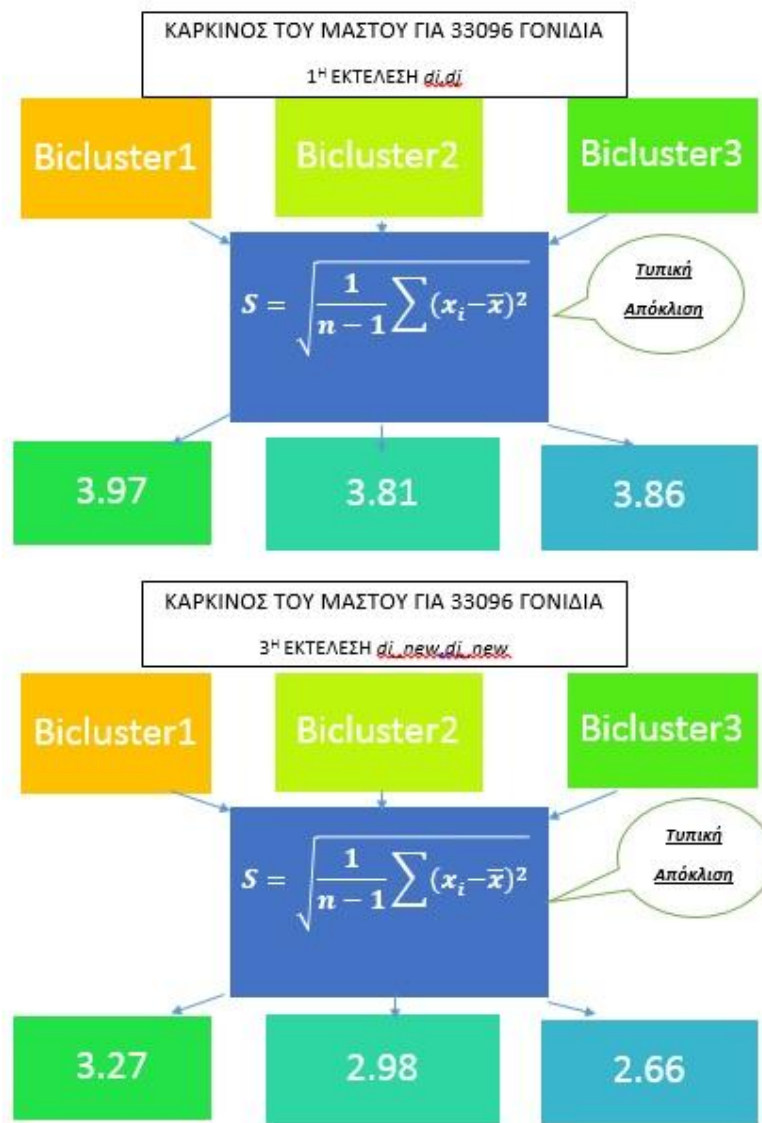
Εικόνα 6.4.δ. - Τρίτη Εκτέλεση του Αλγορίθμου (di_new,dj_new) για 1000 Γονίδια του Καρκίνου του Μαστού Αναδεικνύοντας το γονίδιο 18050

Σε δεύτερη φάση, στην Εικόνα 6.4.α παρατηρείται η απόσταση των γονιδίων και η διακύμανση των γονιδίων στο Bicluster1, κατά την πρώτη εκτέλεση του αλγορίθμου C&C για τα 1000 γονίδια του καρκίνου του μαστού. Επειδή το γονίδιο με probeid 180501 ακολουθεί την ίδια πορεία κατά τη διάρκεια των τριών εκτελέσεων με το γονίδιο 214717 που παρουσιάστηκε παραπάνω, δηλαδή παραμένει και στις τρεις εκτελέσεις στο Bicluster 1, παρουσιάστηκε συνοπτικά η συμπεριφορά του κατά τη διάρκεια των τριών εκτελέσεων.

Στατιστική Επικύρωση Αποτελεσμάτων

Η επικύρωση των αποτελεσμάτων Biclusters είναι ένα σημαντικό πρόβλημα και ταυτόχρονα μια πρόκληση, λόγω των διαφορετικών κριτηρίων που χρησιμοποιούνται και των διαφορετικών στόχων που τίθενται στο πλαίσιο της διπλής κατηγοριοποίησης. Υπάρχουν αρκετές κοινές στρατηγικές επικύρωσης για τα αποτελέσματα διπλής κατηγοριοποίησης, συμπεριλαμβανομένης της επικύρωσης βάσει δεικτών, χρησιμοποιώντας επικύρωση με τη γνώση του πεδίου (domain knowledge) και στατιστικές δοκιμές.

Στην εν λόγω εργασία, για να επικυρώσουμε τα εξαγόμενα αποτελέσματα, χρησιμοποιήσαμε την Τυπική Απόκλιση. Υπολογίσαμε την Τυπική Απόκλιση στην πρώτη και στην τρίτη εκτέλεση του αλγορίθμου. Στην παρακάτω εικόνα παρουσιάζονται οι εξαγόμενες Τυπικές Αποκλίσεις. Παρατηρείται το γεγονός ότι μετά το πέρας και των τριών εκτελέσεων οι Τυπικές Αποκλίσεις έχουν μειωθεί, γεγονός που μας δείχνει ότι πετύχαμε το στόχο μας.



Εικόνα 6.5. - Στατιστική Επικύρωση Αποτελεσμάτων

Οπτικοποίηση Αποτελεσμάτων Διπλής Κατηγοριοποίησης

Η οπτικοποίηση των αποτελεσμάτων διπλής κατηγοριοποίησης είναι δύσκολη λόγω των μη συνεχών και αλληλεπικαλυπτόμενων κατατάξεων τόσο στη σειρά όσο και στη στήλη. Η πιο δημοφιλής τεχνική οπτικοποίησης που αντιπροσωπεύει μια μονή ομάδα διπλής κατηγοριοποίησης είναι η τεχνική θερμικού χάρτη. Ένας θερμικός χάρτης είναι ένα ορθογώνιο πλέγμα που αποτελείται από εικονοστοιχεία, κάθε ένα από τα οποία αντιστοιχεί σε μια τιμή δεδομένων. Οι διαφορετικές κλίμακες γκρι/χρώματος αντιστοιχούν σε διαφορετικές τιμές δεδομένων. Συνήθως, όσο πιο φωτεινό είναι το χρώμα, τόσο μεγαλύτερη είναι η τιμή. Με αυτόν τον τρόπο, αν οι σειρές ή/και οι στήλες του συνόλου δεδομένων αναδιατάσσονται κατάλληλα, το μοτίβο διπλής κατηγοριοποίησης γίνεται αντιληπτό οπτικά. Οι θερμικοί χάρτες αρκούν συνήθως για τον έλεγχο ενός μονού bicluster. Στην παρούσα εργασία, η οπτικοποίηση έγινε με τη βοήθεια της πλατφόρμας

απεικόνισης και ανάλυσης πινάκων GENE-E[61] που σχεδιάστηκε για να υποστηρίξει την εξερεύνηση οπτικών δεδομένων. Η πλατφόρμα GENE-E περιέχει επίσης εργαλεία ειδικά σχεδιασμένα για γονιδιωματικά δεδομένα. Για την παρακολούθηση των αλλαγών που επέρχονται κατά την εφαρμογή της προτεινόμενης μεθοδολογίας, έγινε προσαρμογή (adjustment) των τιμών στην κλίμακα που εμφανίζεται κατά την αρχική εκτέλεση του αλγορίθμου Cheng και Church. Στα σχήματα 6.5. – 6.9., οι εικόνες παρουσιάζονται σε σφαιρική (global) μορφή και όχι σχετική (relative)¹. [62]

Στα σχήματα 6.5.- 6.8 γίνεται οπτικά αντιληπτή η μεγάλη βελτίωση που έχει επιτευχθεί κατά την εφαρμογή της προτεινόμενης μεθοδολογίας στον **υποχώρο δεδομένων** για τον καρκίνο του μαστού, του τραχήλου της μήτρας, του ενδομητρίου, και των ωοθηκών. Αντίθετα, στο σχήμα 6.9. που αφορά στην οπτικοποίηση των αποτελεσμάτων διπλής κατηγοριοποίησης του συνόλου των δεδομένων για τον καρκίνο του μαστού δεν φαίνεται να είναι ορατή η αποτελεσματικότητα της εφαρμογής της προτεινόμενης μεθοδολογίας, γεγονός που μπορεί να οφείλεται στο μεγάλο όγκο των δεδομένων.

Πέρα από την οπτικοποίηση, η οποία επέτρεψε την επικύρωση των αποτελεσμάτων και ανέδειξε την προτεινόμενη μεθοδολογία ως κατάλληλη για εφαρμογή σ' έναν υποχώρο δεδομένων, επιχειρήθηκε παράλληλα και η βιολογική επικύρωση των αποτελεσμάτων μέσω της λειτουργικής ανάλυσης εμπλουτισμού για α) όρους γονιδιακής οντολογίας και β) βιολογικών μονοπατιών με τη βοήθεια της πλατφόρμας WebGestalt.[63] Όπως φαίνεται στους Πίνακες 6.2. και 6.3. μπόρεσαν να εξαχθούν ενδεικτικά συμπεράσματα βάσει των αποσπασματικών αποτελεσμάτων εμπλουτισμού ($p \leq 0.05$) βιολογικών διεργασιών και μοριακών μονοπατιών για τις ομάδες διπλής κατηγοριοποίησης, για τις τρεις εκτελέσεις. Αυτό οφείλεται στους παρακάτω κυρίως λόγους:

Για τα Αποτελέσματα Υποχώρου Δεδομένων για όλους τους τύπους καρκίνου

- Στην έλλειψη σχολιασμού (annotation) των γονιδίων που συμμετέχουν στα Biclusters κατά τη δεύτερη και τρίτη εκτέλεση του αλγορίθμου.

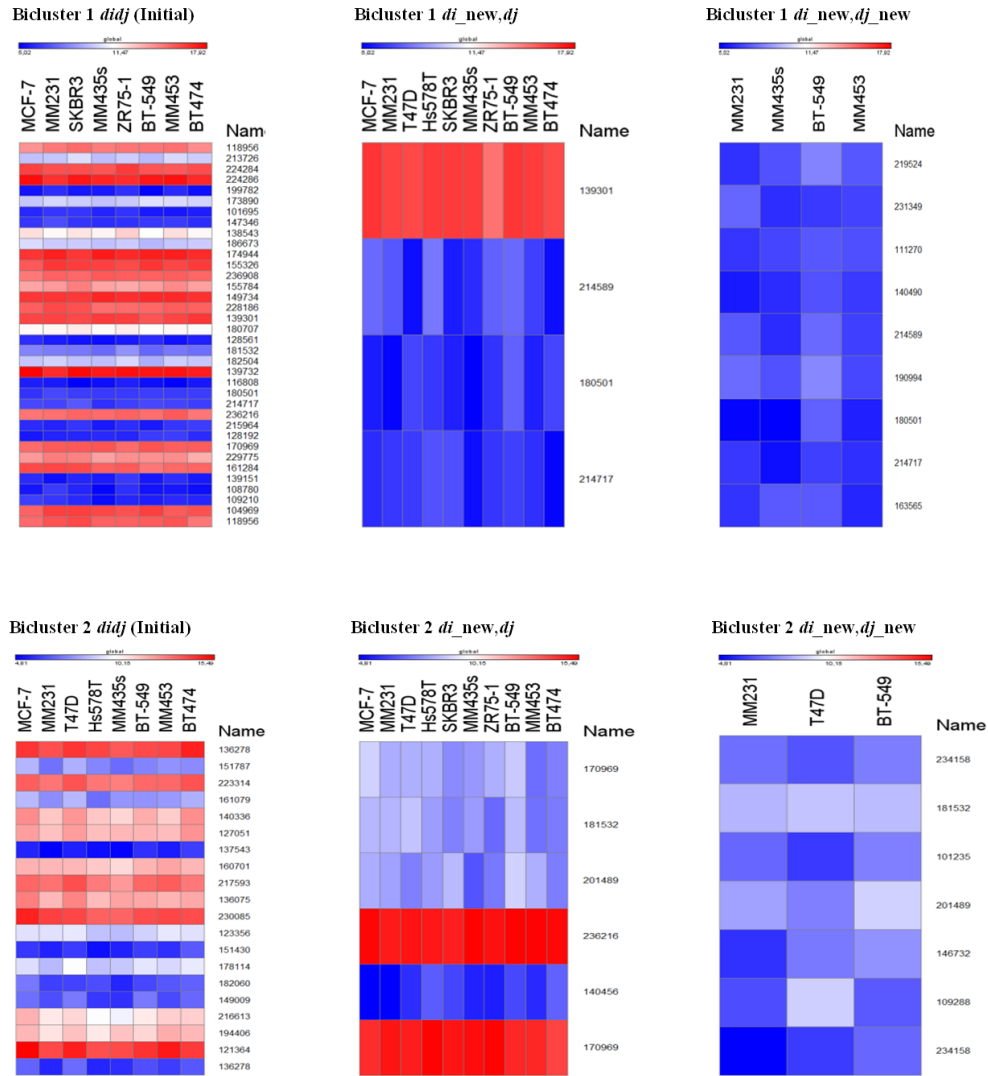
Για τα Αποτελέσματα Συνόλου Δεδομένων για τον καρκίνο του μαστού

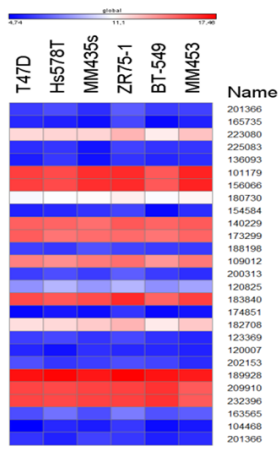
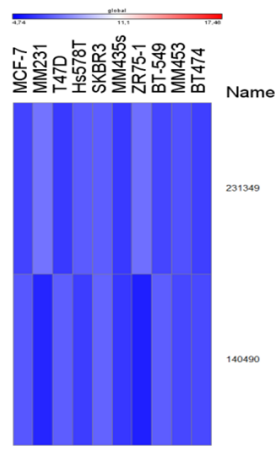
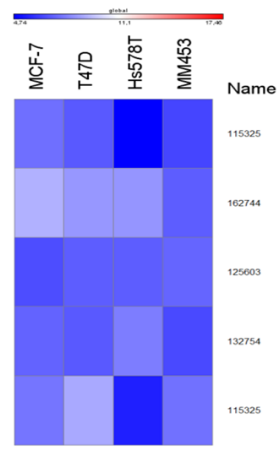
- Στην παρατήρηση ενός πολύ μικρού υποσυνόλου Biclusters (3 Biclusters) σε σχέση με τον πραγματικό αριθμό των Biclusters (π.χ. 2411 Biclusters στον καρκίνο του μαστού).

Για παράδειγμα στον πίνακα 6.2. παρουσιάζονται ικανοποιητικά αποτελέσματα για τα εμπλουτισμένα βιολογικά μονοπάτια στο Bicluster 1 και εν μέρει στο Bicluster 2, που δεν φαίνεται να συνάδουν ωστόσο με τα αποτελέσματα της οπτικοποίησης, ενώ δυστυχώς δεν μπορούμε να αξιολογήσουμε τα αποτελέσματα για το Bicluster 3 γιατί δεν βρέθηκαν εμπλουτισμένα βιολογικά μονοπάτια κατά τη 2^η και 3^η εκτέλεση του αλγορίθμου. Επίσης, στον πίνακα 6.3. παρουσιάζονται ικανοποιητικά αποτελέσματα για τις εμπλουτισμένες βιολογικές διεργασίες στο Bicluster 2 που συνάδουν εν μέρει με τα αποτελέσματα της οπτικοποίησης, αλλά δεν μπορούμε να αξιολογήσουμε τα αποτελέσματα για τα Biclusters 1 και 3 γιατί δεν βρέθηκαν εμπλουτισμένες βιολογικές διεργασίες κατά τη 2^η ή/και 3^η εκτέλεση του αλγορίθμου.

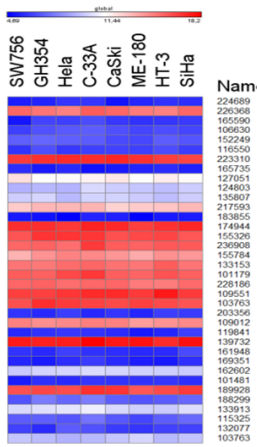
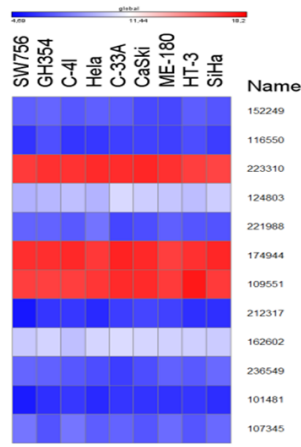
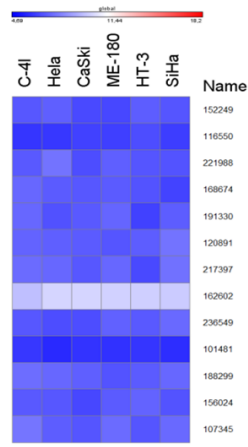
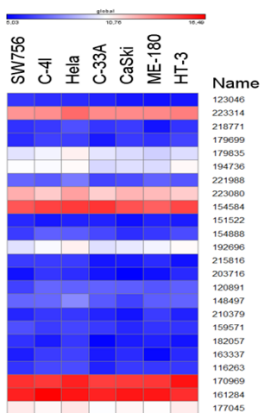
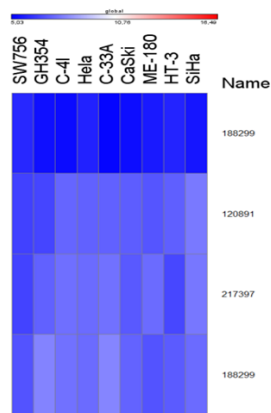
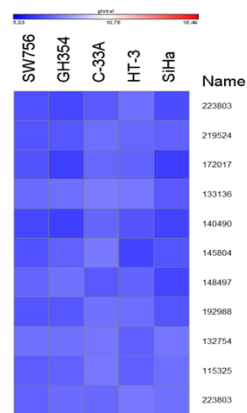
¹ **Σφαιρική (Global):** Η πλατφόρμα GENE-E μετατρέπει τις τιμές σε χρώματα θερμικού χάρτη με τη χρήση των ελάχιστων και μέγιστων τιμών σε ολόκληρο το σύνολο δεδομένων. **Σχετική (Relative):** Η πλατφόρμα GENE-E μετατρέπει τις τιμές σε χρώματα θερμικού χάρτη χρησιμοποιώντας τις μέσες και μέγιστες τιμές για κάθε γονίδιο ή τις τυπικές αποκλίσεις από τον μέσο όρο για κάθε γονίδιο.

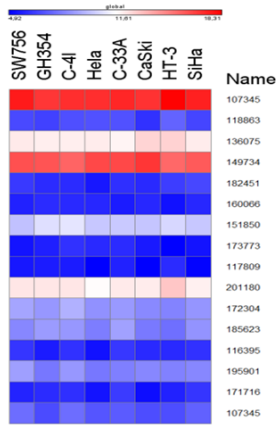
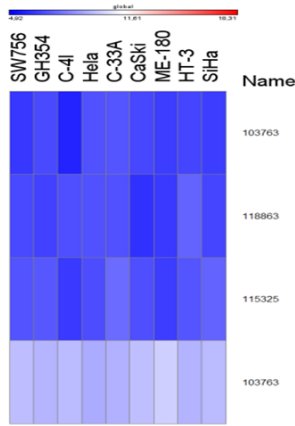
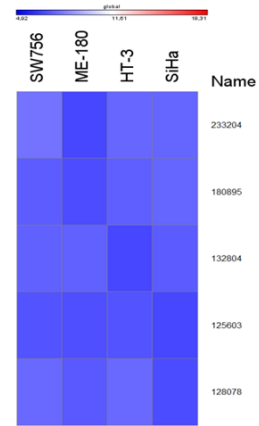
Σχήμα 6.5. – Οπτικοποίηση Αποτελεσμάτων Διπλής Κατηγοριοποίησης Υποχώρου Δεδομένων
(1000 γονίδια - 10 κυτταρικές σειρές)
Καρκίνος του Μαστού



Bicluster 3 *didj* (Initial)Bicluster 3 *di_new,dj*Bicluster 3 *di_new,dj_new*

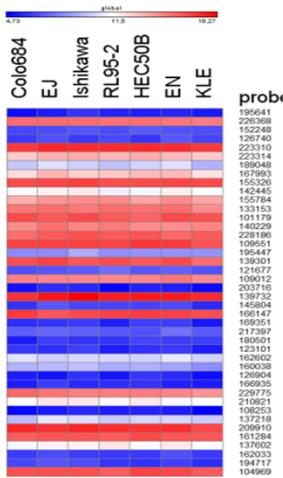
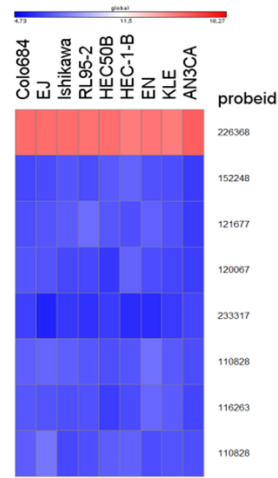
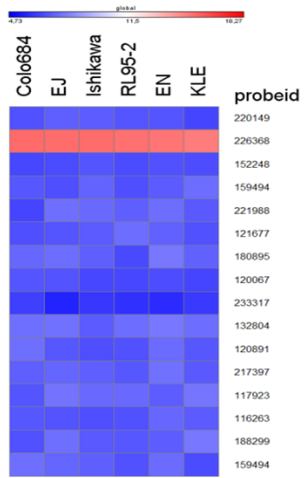
Σχήμα 6.6. - Οπτικοποίηση Αποτελεσμάτων Διπλής Κατηγοριοποίησης Υποχώρου Δεδομένων (1000 γονίδια - 9 κυτταρικές σειρές)
Καρκίνος του Τραχήλου της Μήτρας

Bicluster 1 *didj* (Initial)Bicluster 1 *di_new,dj*Bicluster 1 *di_new,dj_new*Bicluster 2 *didj* (Initial)Bicluster 2 *di_new,dj*Bicluster 2 *di_new,dj_new*

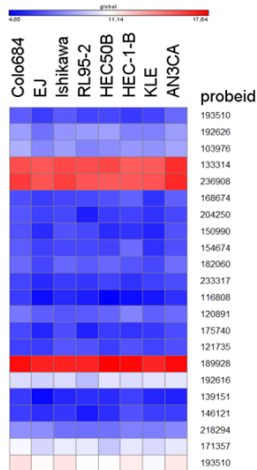
Bicluster 3 *didj* (Initial)Bicluster 3 *di_new,dj*Bicluster 3 *di_new,dj_new*

Σχήμα 6.7. – Οπτικοποίηση Αποτελεσμάτων Διπλής Κατηγοριοποίησης Υποχώρου Δεδομένων (1000 γονίδια - 9 κυτταρικές σειρές)

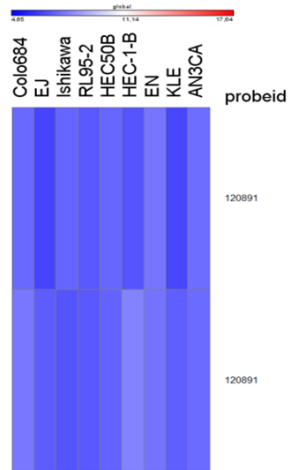
Καρκίνος του Ενδομητρίου

Bicluster 1 *didj* (Initial)Bicluster 1 *di_new,dj*Bicluster 1 *di_new,dj_new*

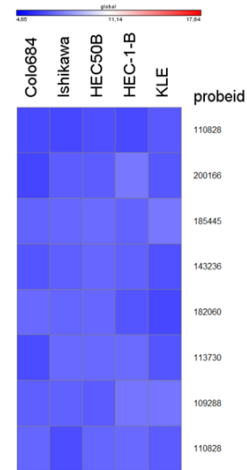
Bicluster 2 *didj* (Initial)



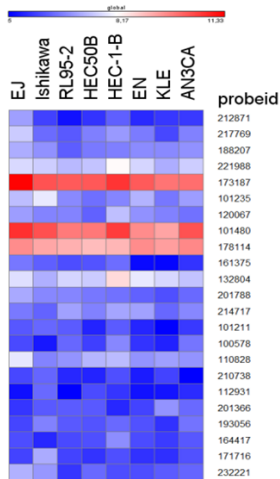
Bicluster 2 *di_new, dj*



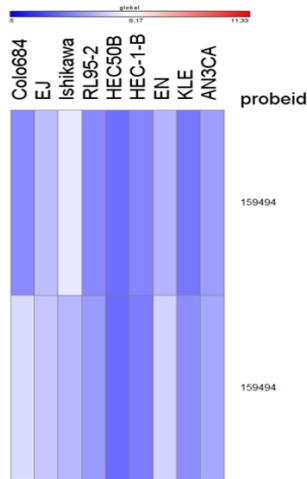
Bicluster 2 *di_new, dj_new*



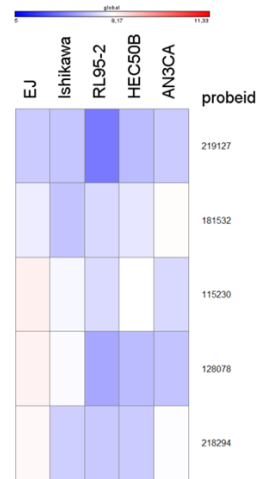
Bicluster 3 *didj* (Initial)



Bicluster 3 *di_new, dj*



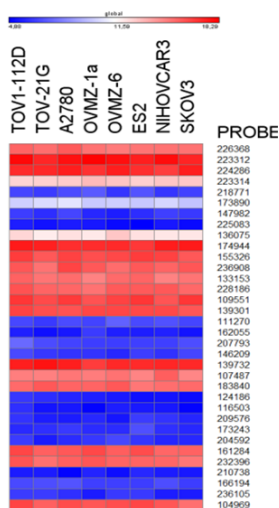
Bicluster 3 *di_new, dj_new*



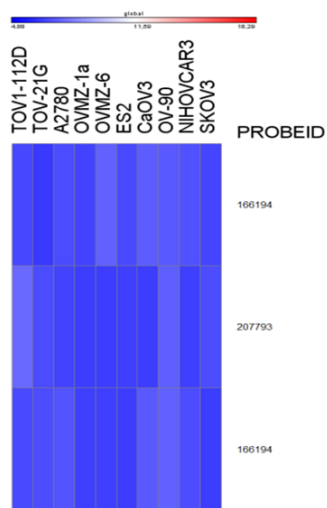
Σχήμα 6.8. - Οπτικοποίηση Αποτελεσμάτων Διπλής Κατηγοριοποίησης Υποχώρου Δεδομένων (1000 γονίδια - 10 κυτταρικές σειρές)

Καρκίνος των Ωοθηκών

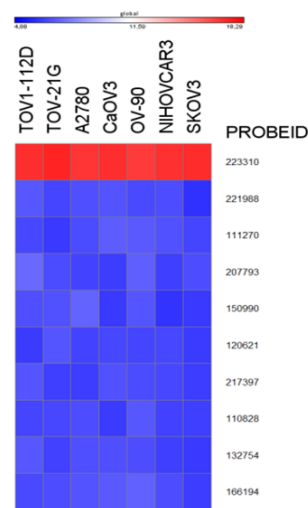
Bicluster 1 *didj* (Initial)



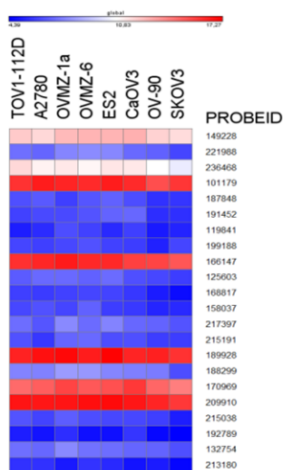
Bicluster 1 *di_new,dj*



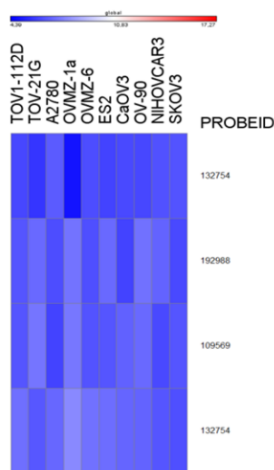
Bicluster 1 *di_new,dj_new*



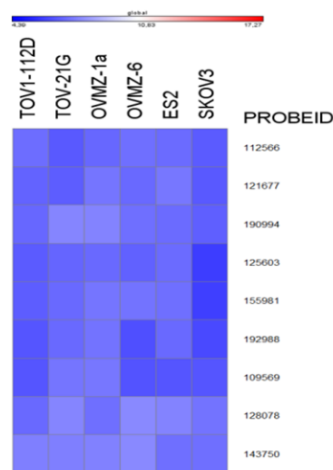
Bicluster 2 *didj* (Initial)



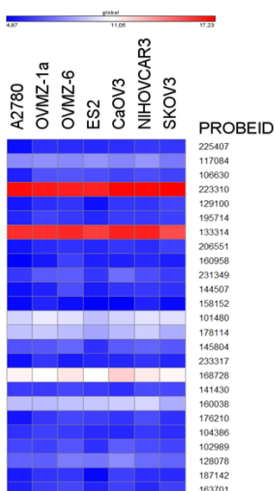
Bicluster 2 *di_new,dj*



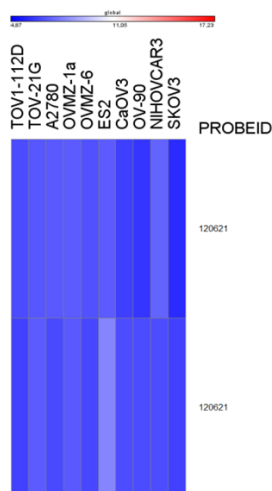
Bicluster 2 *di_new,dj_new*



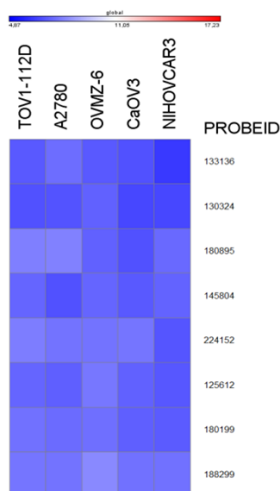
Bicluster 3 *didj* (Initial)



Bicluster 3 *di_new,dj*



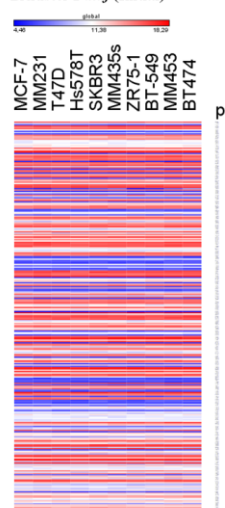
Bicluster 3 *di_new,dj_new*



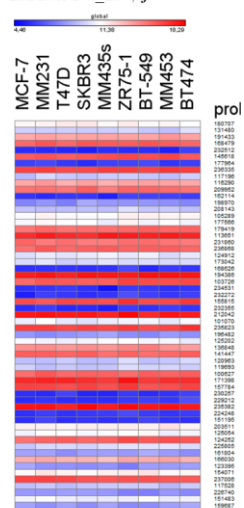
Σχήμα 6.9. - Διπλή Κατηγοριοποίηση Συνόλου Δεδομένων (33096 γονίδια - 10 κυτταρικές σειρές)

Καρκίνος του Μαστού

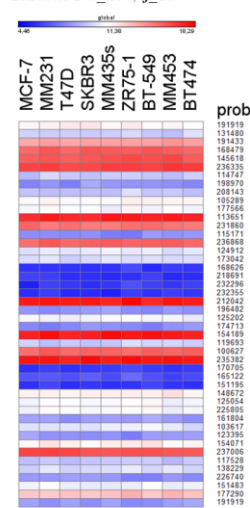
Bidcluster 1 *didj* (Initial)



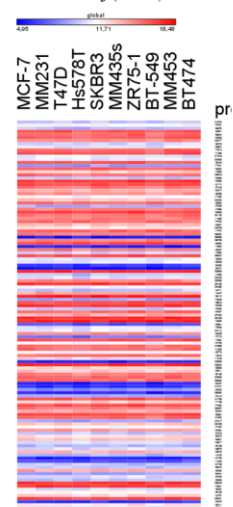
Bidcluster 1 *di_new,dj*



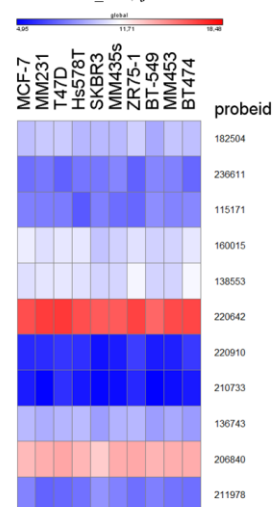
Bidcluster 1 *di_new,dj_new*



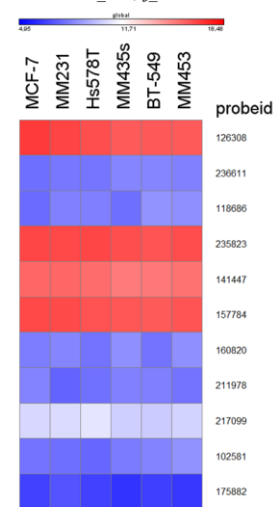
Bidcluster 2 *didj* (Initial)



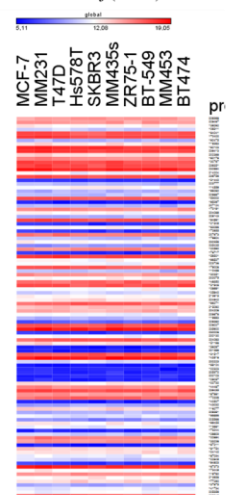
Bidcluster 2 *di_new,dj*



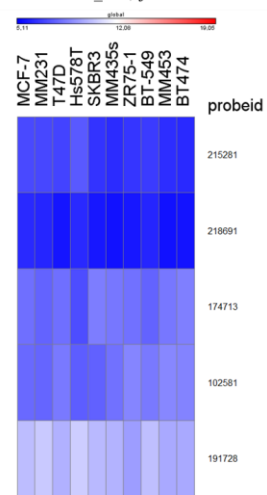
Bidcluster 2 *di_new,dj_new*



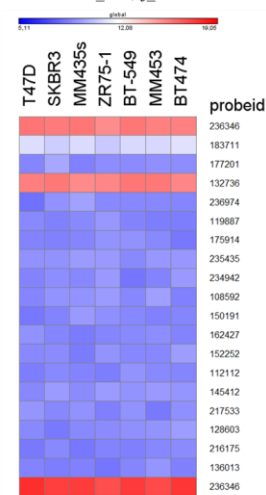
Bidcluster 3 *didj* (Initial)



Bidcluster 3 *di_new,dj*



Bidcluster 3 *di_new,dj_new*



Πίνακας 6.2. - Παρουσίαση των εμπλουτισμένων βιολογικών μονοπατιών ($p \leq 0.05$) για τα τρία Biclusters κατά την διπλή κατηγοριοποίηση του συνόλου των δεδομένων για τον καρκίνο του μαστού

ΒΙΟΛΟΓΙΚΑ ΜΟΝΟΠΑΤΙΑ Καρκίνος του Μαστού - Διπλή Κατηγοριοποίηση Συνόλου Δεδομένων (33096 γονίδια – 10 κυτταρικές σειρές)	
Bicluster 1 <i>didj</i>	<u>Ribosome</u> , <u>Spliceosome</u> , RNA transport, Ribosome biogenesis in eukaryotes, One carbon pool by folate, Ubiquitin mediated proteolysis, Purine metabolism, Pyrimidine metabolism, ABC transporters, Oxidative phosphorylation, Metabolic pathways
Bicluster 1 <i>di_new,dj</i>	<u>Ribosome</u> , Endocytosis
Bicluster 1 <i>di_new,dj_new</i>	<u>Ribosome</u> , <u>Spliceosome</u> , Endocytosis
Bicluster 2 <i>didj</i>	<u>Ribosome</u> , <u>Spliceosome</u> , Oxidative phosphorylation, RNA transport, Protein processing in endoplasmic reticulum, Metabolic pathways, Ubiquitin mediated proteolysis
Bicluster 2 <i>di_new,dj</i>	-
Bicluster 2 <i>di_new,dj_new</i>	<u>Ribosome</u>
Bicluster 3 <i>didj</i>	Ribosome, <u>Spliceosome</u> , Proteasome, Peroxisome
Bicluster 3 <i>di_new,dj</i>	-
Bicluster 3 <i>di_new,dj_new</i>	-

Πίνακας 6.3. - Παρουσίαση των εμπλουτισμένων βιολογικών διεργασιών ($p \leq 0.05$) για το δεύτερο Bicluster κατά την διπλή κατηγοριοποίηση του συνόλου των δεδομένων για τον καρκίνο του μαστού

ΒΙΟΛΟΓΙΚΕΣ ΔΙΕΡΓΑΣΙΕΣ

Καρκίνος του Μαστού - Διπλή Κατηγοριοποίηση Συνόλου Δεδομένων (33096 γονίδια - 10 κυτταρικές σειρές) - Bicluster 2

Βιολογικές Διεργασίες Bcl2-BC <i>didj</i>		Κοινές Βιολογικές Διεργασίες μεταξύ BC-Bcl2 <i>didj</i> και BC-Bcl2 <i>di_new,dj_new</i>	
intracellular protein transport	GO:0006886	translational elongation	GO:0006414
mRNA metabolic process	GO:0016071	translational initiation	GO:0006413
viral reproduction	GO:0016032	translational termination	GO:0006415
intracellular transport	GO:0046907	cotranslational protein targeting to membrane	GO:0006613
cellular protein localization	GO:0034613	SRP-dependent cotranslational protein targeting to membrane	GO:0006614
cellular macromolecule localization	GO:0070727	protein targeting to ER	GO:0045047
multi-organism reproductive process	GO:0044703	establishment of protein localization to endoplasmic reticulum	GO:0072599
macromolecule localization	GO:0033036	translation	GO:0006412
Βιολογικές Διεργασίες BC-Bcl2 <i>di_new,dj</i>		viral transcription	GO:0019083
nucleobase-containing compound biosynthetic process	GO:0034654	viral genome expression	GO:0019080
RNA biosynthetic process	GO:0032774	nuclear-transcribed mRNA catabolic process, nonsense-mediated decay	GO:0000184
nucleobase-containing compound metabolic process	GO:0006139	protein localization to endoplasmic reticulum	GO:0070972
RNA metabolic process	GO:0016070	establishment of protein localization to organelle	GO:0072594
organic cyclic compound metabolic process	GO:1901360	viral infectious cycle	GO:0019058
aromatic compound biosynthetic process	GO:0019438	nuclear-transcribed mRNA catabolic process	GO:0000956
cellular aromatic compound metabolic process	GO:0006725	mRNA catabolic process	GO:0006402
cellular nitrogen compound metabolic process	GO:0034641	protein targeting to membrane	GO:0006612
organic cyclic compound biosynthetic process	GO:1901362	cellular protein complex disassembly	GO:0043624
heterocycle biosynthetic process	GO:0018130	viral reproductive process	GO:0022415
gene expression	GO:0010467	protein targeting	GO:0006605
heterocycle metabolic process	GO:0046483	RNA catabolic process	GO:0006401
cellular nitrogen compound biosynthetic process	GO:0044271	protein complex disassembly	GO:0043241
nitrogen compound metabolic process	GO:0006807	protein localization to organelle	GO:0033365
nucleic acid metabolic process	GO:0090304	cellular macromolecular complex disassembly	GO:0034623
virus-host interaction	GO:0019048	macromolecular complex disassembly	GO:0032984
cellular macromolecule metabolic process	GO:0044260	cellular component disassembly at cellular level	GO:0071845
cellular macromolecule biosynthetic process	GO:0034645	cellular component disassembly	GO:0022411
macromolecule biosynthetic process	GO:0009059	protein transport	GO:0015031
interaction with host	GO:0051701	establishment of protein localization	GO:0045184
interspecies interaction between organisms	GO:0044419	protein localization	GO:0008104
symbiosis, encompassing mutualism through parasitism	GO:0044403	cellular protein metabolic process	GO:0044267
Βιολογικές Διεργασίες BC-Bcl2 <i>di_new,dj_new</i>		establishment of localization in cell	GO:0051649
cellular macromolecular complex subunit organization	GO:0034621		
regulation of translation	GO:0006417		
protein complex subunit organization	GO:0071822		
macromolecular complex subunit organization	GO:0043933		
regulation of cellular protein metabolic process	GO:0032268		
posttranscriptional regulation of gene expression	GO:0010608		
regulation of protein metabolic process	GO:0051246		
cellular process involved in reproduction	GO:0048610		

ΣΥΜΠΕΡΑΣΜΑΤΑ ΚΑΙ ΜΕΛΛΟΝΤΙΚΕΣ ΕΠΕΚΤΑΣΕΙΣ

7.1 Συμπεράσματα

Στην παρούσα εργασία μελετήσαμε τις μεθόδους διπλής κατηγοριοποίησης σε πίνακες γονιδιακής έκφρασης για την εξαγωγή υποομάδων με στατιστική και βιολογική σημασία.

Η πειραματική διαδικασία εφαρμόστηκε σε αρχεία που περιλάμβαναν πίνακες γονιδίων-κυτταρικών σειρών προερχόμενα από μικροσυστοιχίες DNA από καρκινικές κυτταρικές σειρές του μαστού, του τραχήλου της μήτρας, των ωοθηκών και του ενδομητρίου. Τα αποτελέσματα που λάβαμε ήταν σημαντικά και επικυρώθηκαν μέσω της οπτικοποίησής τους.

Συγκεκριμένα:

Ο βελτιωμένος biclustering αλγόριθμος που χρησιμοποιήθηκε απέδειξε τους λόγους για τους οποίους, ο αρχικός αλγόριθμος έχριζε βελτίωσης και ότι ο βελτιωμένος αλγόριθμος διπλής κατηγοριοποίησης υπερτερεί, δίνοντας έμφαση στην πληροφορία, μειώνοντας το θόρυβο και παρέχοντας μας μια εικόνα για το σύνολο δεδομένων πιο συγκεκριμένη και λεπτομερής.

Σε αυτό το σημείο είναι σημαντικό να αναφερθεί το γεγονός ότι ο βελτιωμένος αλγόριθμος που προτείναμε, φαίνεται να είναι πιο αποδοτικός και να λειτουργεί καλύτερα στον υποχώρο γονιδίων. Η εξαγωγή ενός μεγάλου αριθμού γονιδίων σε μία ομάδα διπλής κατηγοριοποίησης (Biclusters) δυσκολεύει το έργο του αλγορίθμου και δεν έχουμε τόσο ικανοποιητικά αποτελέσματα.

Όσον αφορά τον υποχώρο γονιδίων, οι ομάδες που δημιουργήθηκαν, εμφάνιζαν γονίδια με όμοια συμπεριφορά κατά μήκος των κυτταρικών σειρών και παράλληλα το εύρος των τιμών τους ήταν μικρό, κάνοντας έτσι τις ομάδες πιο συνεκτικές και ικανές να εξάγουν σημαντικές βιολογικές πληροφορίες, γενικές αλλά και πιο εξειδικευμένες.

7.2 Μελλοντικές επεκτάσεις

Όπως συμβαίνει με τις περισσότερες μελέτες, σε περιπτώσεις όπως η δική μας, όπου ο όγκος των δεδομένων είναι αρκετά μεγάλος, επιλέξαμε να μελετήσουμε έναν μικρό αριθμό Biclusters για την καλύτερη αξιολόγηση των αποτελεσμάτων. Έτσι μελλοντικά θα μπορούσαμε να επεξεργαστούμε επιπλέον ομαδοποιήσεις που είναι πιθανό να μας παρέχουν χρήσιμη πληροφορία κι έτσι να κρίνουμε ακόμα πιο αντικειμενικά τα τωρινά μας αποτελέσματα.

Επιπλέον, ένα επόμενο βήμα θα αποτελούσε η πιο αναλυτική επεξεργασία των μεγάλων Biclusters που μας παρείχαν αρκετή και σημαντική πληροφορία. Αυτά δηλαδή, που μελετώντας τα πιο μεθοδικά θα μπορούσαμε να καταλήξουμε σε βιολογικές διεργασίες με μεγαλύτερη λεπτομέρεια και εξειδίκευση.

Τέλος, πέραν του αλγορίθμου Cheng και Church, θα ήταν πολύ σημαντική η εφαρμογή και άλλων αλγορίθμων biclustering, για τη σύγκριση των αποτελεσμάτων με τη δική μας εφαρμογή και έτσι θα μπορούσαμε να έχουμε μία πλήρη και πιο σφαιρική εικόνα για την τεχνική biclustering.

Βιβλιογραφία

- [1] Tanay A., Sharan R., Shamir R., “Discovering statistically significant biclusters in gene expression data”, *Bioinformatics*, 18 Suppl 1:S136-44. 2002.
- [2] Saber H.B., and Elloumi M., “A comparative study of clustering and biclustering of microarray data”, *International Journal of Computer Science & Information Technology (IJCSIT)*, vol. 6, no. 6, pp. 93-111, 2014.
- [3] Graphical Abstract: <http://dx.doi.org/10.1016/j.jbi.2015.06.028>.
- [4] Cheng Y., and Church G.M., “Biclustering of Expression Data” *Proc Int Conf Intell Syst Mol Biol.*, vol. 8, pp. 93-103, 2000.
- [5] Alberts B., “Molecular biology of the cell”, 5th ed. New York: Garland Science, 2002.
- [6] Hartigan J.A., “Clustering Algorithms”, John Wiley & Sons, New York, 1975.
- [7] Mirkin B., “Mathematical Classification and Clustering”, Kluwer Academic Publishers, 1996. ISBN 0-7923-4159-7.
- [8] Tanay A., Sharon R., and Shamir R., “Biclustering gene expression data”, In *Proceedings of ISMB 2002*, pp. S136–S144, 2002.
- [9] Ihmels J. Bergmann S. Barkai N. Defining transcriptional modules using large-scale gene expression data. *Bioinformatics.*, vol. 20, no. 13, pp. 1993–2003, 2004.
- [10] Madeira S.C., and Oliveira A.L., “Biclustering algorithms for biological data analysis: a survey”, *IEEE/ACM Trans. Comput. Biol. Bioinform.*, vol. 1, no. 1, pp. 24–45, 2004.
- [11] Pontes B., Giráldez R., Aguilar-Ruiz J.S., “Biclustering on expression data: A review”, *Journal of Biomedical Informatics*, vol. 57, pp. 163-180, 2015.
- [12] Αλεβυζάκη Ανδρονίκη, «Μέθοδοι διπλής κατηγοριοποίησης “biclustering” για επιλογή γονιδιακών δεικτών από κυτταρικές σειρές». Διπλωματική Εργασία, Πολυτεχνείο Κρήτης, 2010.
- [13] Tang A.H., Neufeld T.P., Rubin G.M., Müller H.-A.J., “Transcriptional regulation of cytoskeletal functions and segmentation by a novel maternal pair-rule gene, lilliputian”, *Development* vol. 128, pp. 801-813, 2001.
- [14] Ben-Dor A., Chor B., Karp R., and Yakhini Z., “Discovering local structure in gene expression data: the order-preserving submatrix problem”, *J. Comput. Biol.* vol. 10, no. 3-4, pp. 373–384, 2003.
- [15] Segal E., Battle A., and Koller D., “Decomposing gene expression into cellular processes”, *Pac. Symp. Biocomput.* vol. 8, pp. 89–100, 2003.
- [16] Murali T.M., and Kasif S., “Extracting conserved gene expression motifs from gene expression data”, *Pac. Symp. Biocomput.* vol. 8, pp. 77–88, 2003.
- [17] Sheng Q., Moreau Y., and De Moor B., “Biclustering micrarray data by gibbs sampling”, *Bioinformatics*, vol. 19, no. Suppl.2, p. ii196–ii205, 2003.
- [18] Cano C., Adarve L., López J., and Blanco A., “Possibilistic approach for biclustering microarray data”, *Comput Biol Med.*, vol. 37, no. 10, pp. 1426–1436, 2007.
- [19] Maulik U., Mukhopadhyay A., and Bandyopadhyay S., “Finding multiple coherent biclusters in microarray data using variable string length multiobjective genetic algorithm”, *IEEE Trans. Inf. Technol. Biomed.*, vol. 13, no. 6, pp. 969–975, 2009.
- [20] Chen C.-P., Fushing H., Atwill R., and Koehl P., “biDCG: a new method for discovering global features of DNA microarray data via an iterative re-clustering procedure”, *PLoS One*, vol. 9, no. 7, p. e102445, 2014.
- [21] [https://gr.dreamstime.com/\[image24298020\]](https://gr.dreamstime.com/[image24298020])
- [22] Lodish H., Berk A., Matsudaira P., Kaiser C.A., Krieger M., Scott M.P., Zipursky L., and Darnell J., “Molecular Cell Biology”, 5th edition. W.H. Freeman & Co (Sd), 2003.
- [23] <http://ebooks.edu.gr/>
- [24] <https://www.slideshare.net/annpyl/dna-rna>
- [25] http://mathimatabiologias.weebly.com/uploads/1/0/9/7/10979738/313406_orig.png
- [26] Βιοπληροφορική-αναζήτηση πληροφορίας σε βιολογικές βάσεις δεδομένων: http://ecourses.dbnet.ntua.gr/fsr/17066/ERGASTHRIAKH%20ASKHSH%202_sumplhrwmatiko.pdf
- [27] <http://www.chem.uoa.gr/?cat=39>
- [28] <http://www.sciencedirect.com/topics/page/Transcriptome>
- [29] https://repository.kallipos.gr/bitstream/11419/1585/1/Chapter07_geneexpression_R.pdf
- [30] Μαλατράς Α., «Εφαρμογές Υπολογιστικής Νοημοσύνης σε Μικροσυστοιχίες DNA», Πτυχιακή εργασία. Πανεπιστήμιο Στερεάς Ελλάδας, 2010.
- [31] *microarrays dna_images*

- [32] Jiang D., Tang C., and Zhang A., “Cluster Analysis for Gene Expression Data: A Survey” IEEE Transactions on Knowledge and Data Engineering, vol. 16, no. 11, pp. 1370-1386, 2004.
- [33] Brazma A., Parkinson H., Sarkans U., Shojatalab M., Vilo J., Abeygunawardena N., Holloway E., Kapushesky M., Kemmeren P., Lara G.G., Oezcimen A., Rocca-Serra P., Sansone S.-A., ArrayExpress—a public repository for microarray gene expression data at the EBI. Nucleic Acids Res., vol. 31, no. 1, pp. 68-71, 2003.
- [34] Lee M.L.T., “Analysis of Microarray Gene Expression Data”, Springer, 2004.
- [35] Zhao H., Liew A.W.-C., Wang D.Z., and Yan, H. Biclustering analysis for pattern discovery: Current techniques, comparative studies and applications. Current Bioinformatics, vol. 7, no. 1, pp. 43-55, 2012.
- [36] <http://www.eng.ucy.ac.cy/cpitris/courses/ECE370/presentations/14.%20Biomaterials.pdf>
- [37] <http://eclass.uth.gr/eclass/modules/document/file.php/DIB200/chapter2final.pdf>
- [38] <https://www.ncbi.nlm.nih.gov/genbank/>
- [39] <https://www.nlm.nih.gov/bsd/pmresources.html>
- [40] <https://www.ncbi.nlm.nih.gov/geo/>
- [41] GenBank Statistics
- [42] http://students.ceid.upatras.gr/~maurakis/Gene_Expression_Analysis.pdf
- [43] Agrawal R., Gehrke J., Gunopulos D., and Raghavan P., “Automatic Subspace Clustering of High Dimensional Data for Data Mining Applications”, SIGMOD '98 Proceedings of the 1998 ACM SIGMOD international conference on Management of data, vol. 27, no. 2, pp. 94-105, 1998.
- [44] Eisen M.B., Spellman P.T., Brown P.O., and Botstein D., “Cluster analysis and display of genome-wide expression patterns”, Proc. Natl. Acad. Sci. USA, vol. 95, no. 25, pp. 14863– 14868, 1998.
- [45] <https://www.slideshare.net/raedald/automatic-subspace-clustering-of-high-dimensional-data-for-data-mining-application-32900971>
- [46] Tchagang A.B., Pan Y., Famili F., Tewfik A.H., and Benos P.V., “Biclustering of DNA Microarray Data: Theory, Evaluation, and Applications”, Handbook of Research on Computational and Systems Biology: Interdisciplinary Applications, IGI Global, 2011.
- [47] Sharan R., Porat U.B., and Bleiberg O., “Analysis of Biological Networks: Network Modules– Clustering and Biclustering ”, Lecture 5, November 23 , pp. 1–21, 2006.
- [48] [http://www.microarrays.ca/services/“What Is Hierarchical Clustering ?”](http://www.microarrays.ca/services/“What%20Is%20Hierarchical%20Clustering%20?”)
- [49] <http://homes.di.unimi.it/~valenti/SlideCorsi/MB0910/HierarchicalClustering.pdf>
- [50] Maimon O., and Rokach L., “Clustering methods”, in Data mining and knowledge discovery handbook. Springer US, pp. 321-352, 2005.
- [51] Manning C.D., Raghavan P., and Schütze H., “Introduction to Information Retrieval”, Cambridge University Press, 2008.
- [52] https://www.google.gr/url?sa=i&rct=j&q=&esrc=s&source=images&cd=&ved=0ahUKEwiJj9a99szUAhVFNxQKHUy4DDYQjxwIAw&url=https%3A%2F%2Fwww.slideshare.net%2Fgveress%2Fcluster-training-2013&psig=AFQjCNGxX6_dkRmkENDcCdFObjYojWyCsA&ust=1498065313474950
- [53] https://www.google.gr/url?sa=i&rct=j&q=&esrc=s&source=images&cd=&ved=0ahUKEwiJj9a99szUAhVFNxQKHUy4DDYQjxwIAw&url=https%3A%2F%2Fwww.slideshare.net%2Fgveress%2Fcluster-training-2013&psig=AFQjCNGxX6_dkRmkENDcCdFObjYojWyCsA&ust=1498065313474950
- [54] <https://en.wikipedia.org/wiki/Biclustering>
- [55] Blackshaw S., Harpavat S., Trimarchi J., Cai L., Huang H., Kuo W. P., Weber G., Brazma, and Vilo J., “Minireview: Gene Expression Data Analysis”, FEBS Letters 480, pp. 17-24, 2003.
- [56] Pontes B., Girdes R., and Aguilar-Ruiz J.S., “Quality Measures for Gene Expression Biclusters”, Margis R, ed. PLoS ONE, vol. 10, no. 3, p. e0115497, 2015.
- [57] Obermayr, Sanchez-Cabo E. F., Tea M. K. M., Singer C. F., Krainer M., Fischer M. B., Sehouli J., Reinthaller A., Horvat R., Heinze G., Tong D., and Zeillinger R., “Assessment of a six gene panel for the molecular detection of circulating tumor cells in the blood of female cancer patients.” BMC Cancer, vol. 10, no. 1, p. 666, 2010.
- [58] NCI Dictionary of Cancer Terms
- [59] Cancer Cell Line Encyclopedia (CCLE)
- [60] Jacob F., Nixdorf S., Hacker N.F., and Heinzelmann-Schwarz V.A., “Reliable in vitro studies require appropriate ovarian cancer cell lines.” J. Ovarian Res., vol. 7, p. 60, 2014.
- [61] <https://software.broadinstitute.org/GENE-E/>
- [62] <https://software.broadinstitute.org/GENE-E/doc.html>
- [63] http://www.webgestalt.org/webgestalt_2013/option.php

ΠΑΡΑΡΤΗΜΑ

Στο Παράρτημα παρατίθενται τα γονίδια που ομαδοποιούνται στα τρία πρώτα Biclusters στην πρώτη, δεύτερη και τρίτη εκτέλεση του αλγορίθμου. Δίνονται οι αναλυτικοί Πίνακες για τον καρκίνο του μαστού, του τραχήλου της μήτρας, του ενδομητρίου και των ωοθηκών για τα 1000 γονίδια του “dataset”.

Καρκίνος του Μαστού

di,dj ->Πρώτη εκτέλεση αλγορίθμου

initial									
di_dj									
bicluster1	PROBEID	MCF-7	MM231	SKBR3	MM435s	ZR75-1	BT-549	MM453	BT474
18	226368	14,31719	14,95013	15,29338	14,83902	15,05665	15,05267	15,16847	14,47759
42	213726	9,975586	10,05667	10,67213	9,856768	10,26455	9,752279	10,56243	10,2598
105	224284	16,35968	15,97745	16,04696	15,87873	16,51371	15,88931	15,78106	16,05529
107	224286	17,70126	16,79674	17,45814	17,09883	16,8679	17,45312	17,58593	17,00214
120	199782	5,545611	6,226255	5,671115	5,716348	6,237167	5,131927	6,104767	5,46478
150	173890	10,10765	10,45769	10,24529	9,928914	10,17369	10,70491	10,61749	10,41654
153	101695	6,072557	6,397267	6,426852	5,952134	6,356477	5,687495	5,634138	5,790055
177	147346	6,395415	7,038766	6,417404	6,198541	6,31326	6,483754	6,572028	6,190864
190	138543	12,37473	11,7103	12,21802	11,83362	12,77035	11,55955	12,40636	11,74627
259	186673	10,62225	10,19384	10,20132	9,946289	10,20132	9,702818	10,07868	10,18312
273	174944	16,70533	17,20289	16,62122	17,35966	16,98538	17,25555	17,22564	16,93257
289	155326	15,78106	16,62461	16,47307	16,3555	16,12322	16,70292	16,22788	16,50926
303	236908	14,95278	15,50219	15,37988	15,76231	15,59826	15,56187	15,64221	15,42449
310	155784	13,82736	14,08613	14,82038	13,9083	13,71764	14,05261	14,49957	13,93928
320	149734	16,62461	16,66969	16,57633	17,08047	16,81259	16,99287	17,00442	16,77015
336	228186	15,59826	16,19141	15,98615	15,22927	15,32926	15,40191	15,65025	15,56187
343	139301	16,29054	16,31057	16,15504	16,53159	16,44803	15,86532	16,26508	16,65611
387	180707	11,68976	11,80061	12,13062	11,66029	12,03644	11,58405	11,77812	11,39865
390	128561	6,221719	5,762957	5,561125	5,375277	5,711893	6,151954	5,480831	5,835731
438	181532	8,373335	8,207122	8,078394	7,736715	8,757705	7,60793	8,076276	7,838595
499	182504	10,07569	10,3546	9,93524	9,913831	10,50159	9,463651	10,23471	10,09346
563	139732	17,923	16,95017	17,80707	17,40093	17,64809	17,4688	17,55437	17,27832
606	116808	5,829279	5,701691	5,020673	5,47176	5,874409	5,389694	5,113216	5,655472
678	180501	6,377281	6,885293	6,627092	6,535787	6,513497	6,610454	6,575171	6,590028
718	214717	6,876456	6,768432	7,275703	6,26037	6,673344	6,873969	6,747133	6,620442
730	236216	15,04358	15,14456	15,44834	15,33159	15,53159	15,2265	15,69976	14,99796
792	215964	6,148683	5,965177	6,433183	5,690075	6,229993	5,722257	6,421423	6,589592
795	128192	5,636823	5,939517	5,998356	6,12362	6,028395	6,584104	6,61524	6,173348
801	170969	15,31567	15,7042	15,79074	15,8865	15,31334	15,91738	15,63375	15,64662
803	229775	14,17414	14,85701	14,69149	13,97004	14,42462	14,37224	14,43587	13,5761
818	161284	16,14322	16,2218	16,10162	15,56673	15,69749	15,01808	15,39034	15,4559
845	139151	6,155465	5,426909	6,012651	5,732012	6,499128	5,648874	6,208971	5,703918
909	108780	5,605265	6,572742	6,000083	5,215229	6,237474	5,52839	5,673374	5,401041
915	109210	6,659545	6,387754	6,096934	5,33436	6,018268	5,809462	6,126813	5,99687
978	104969	15,53605	16,42501	16,4331	16,22309	15,99071	16,41701	15,73791	15,95591
994	118956	15,32926	15,90111	16,19141	15,55201	15,58661	15,69231	15,87653	14,98848
[1,2,5,6,7,8,9,10]									

Bicluster2	Probeid	MCF-7	MM231	SKBR3	MM435s	ZR75-1	BT-549	MM453	BT474
47	164432	14,46505	13,87184	14,46997	14,11255	13,7036	14,00062	14,07032	14,88423
69	151787	8,635632	7,239536	8,52605	7,634946	7,077521	7,59948	7,89729	7,779308
108	223314	13,5969	13,12667	13,84732	13,09266	12,91258	13,36127	13,387	13,83196
146	161079	8,747381	7,80924	8,677046	7,145436	7,992976	7,929998	8,043446	8,517164
149	140336	12,61235	11,4145	12,40398	11,42642	11,12099	11,96932	11,29658	12,58113
187	127051	12,24356	11,5182	12,28037	11,56633	11,72818	11,77746	11,61646	12,47905
193	137543	5,568	4,812809	5,485596	5,130313	4,920732	5,859754	5,432801	6,11472
196	160701	11,81033	11,7453	11,7571	11,38507	11,0268	11,76741	12,11732	11,76157
244	217593	13,41063	13,18044	13,9013	12,99592	12,38513	13,42249	13,64528	13,03966
256	136075	12,90104	11,69179	12,77364	11,63942	12,1696	12,07241	11,80657	12,04721
296	230085	14,89965	14,21666	14,07188	13,54719	13,45154	13,97161	14,07841	13,96269
503	123356	9,555903	9,573616	9,699554	8,893515	8,768327	9,23255	9,70536	9,550708
538	151430	6,002187	5,542083	6,075485	5,17647	5,541288	6,146669	5,819265	6,723677
540	178114	9,462366	8,730112	10,14004	8,890063	9,030732	9,643751	9,347486	9,746303
587	182060	7,347799	6,167629	6,248132	6,476599	5,644152	6,329511	6,667814	6,868262
655	149009	7,197602	6,459343	7,373792	6,023935	5,870365	7,385542	6,382085	7,023492
770	216613	11,88129	10,7117	11,05489	10,17339	9,921691	10,63467	11,23681	11,24038
814	194406	11,72418	10,8154	11,52313	10,70896	11,48223	11,09369	11,97138	11,69343
825	121364	15,49235	14,03838	14,94737	13,93928	14,01721	14,50678	15,0384	14,35543
941	136278	6,931848	5,631319	7,093789	5,628536	5,816121	6,51293	6,126555	6,671591
[1,2,3,4,6,8,9,10]									

Bicluster3	probeid	T47D	Hs578T	MM435s	ZR75-1	BT-549	MM453
6	217298	6,253944	6,656327	5,898956	7,070186	6,277799	6,412153
159	165735	5,735893	5,567824	5,136867	6,532075	5,005613	5,584856
238	223080	12,15183	12,22873	12,22808	13,02781	11,65189	12,70626
240	225083	5,909558	6,096884	5,248097	6,421529	6,059	6,090798
263	136093	5,510129	6,165769	5,198784	6,044713	5,731562	5,914829
313	101179	15,97745	15,69485	16,3966	16,47307	15,34586	16,70081
315	156066	15,97745	15,69485	16,3966	16,47307	15,34586	16,70081
321	180730	10,96456	11,09423	11,26086	11,61886	11,05733	11,40091
324	154584	5,611654	5,744894	6,053406	6,471021	5,034153	5,896183
326	140229	15,13648	15,05267	14,78292	15,33382	15,02597	15,29581
354	173299	15,23826	14,64175	14,84127	14,87739	14,64782	15,0492
404	188198	6,365204	6,147532	6,823609	6,452447	6,273496	6,216863
411	109012	14,35402	13,98778	14,50318	14,81061	13,8092	14,54862
550	200313	6,288208	6,681701	6,234391	7,102128	6,236535	5,903964
617	120825	8,398501	9,320165	8,55609	9,322817	8,375298	8,854036
632	183840	15,67562	15,32926	15,63375	16,48492	15,02964	15,95401
635	174851	4,944075	5,343539	5,132999	6,102083	5,02972	5,248587
665	182708	11,9843	12,16154	12,47628	13,03966	11,7103	12,52424
711	123369	6,325246	6,571876	6,047078	6,575121	6,012033	5,698496
716	120007	5,732923	5,229325	5,968539	5,75379	5,780379	5,631384
743	202153	6,766566	6,102506	6,314251	6,713627	6,446853	5,849025
774	189928	17,13934	17,43537	16,99287	17,45814	17,09883	16,86047
813	209910	15,85037	15,71945	15,89665	15,91533	16,26755	15,30499
839	232396	15,85037	15,71945	15,89665	15,91533	16,26755	15,30499
853	163565	6,68222	7,602164	6,686453	7,84799	6,809898	7,113893
863	104468	4,73585	5,759724	5,468223	6,154048	5,047554	5,377096
888	201366	6,079344	5,622237	6,424458	6,177023	5,565188	5,843548
[3,4,6,7,8,9]							

Καρκίνος του Μαστού

di_new,dj ->Δεύτερη εκτέλεση αλγορίθμου

dinew_dj										
Bicluster1	probeid	MCF-7	MM231	T47D	Hs578T	SKBR3	MM435s	ZR75-1	BT-549	MM453
343	139301	16,66128	16,30631	16,19141	16,53715	16,51115	16,45829	15,13869	16,64168	16,49421
462	214589	7,734273	7,289839	5,399297	8,090775	5,751375	6,179302	7,512713	7,718731	6,661609
678	180501	5,762265	5,173406	6,836401	6,451068	5,851414	4,637098	6,353068	7,501148	5,874923
718	214717	6,145841	6,44112	6,568804	7,100972	7,053609	5,376219	6,428902	6,670218	6,255565
[1,2,3,4,5,6,7,8,9,10]										

Bicluster2	probeid	MCF-7	MM231	T47D	Hs578T	SKBR3	MM435s	ZR75-1	BT-549	MM453
395	162744	9,219379	8,467646	8,538543	8,513284	7,785826	8,020587	8,526212	9,064433	7,122914
438	181532	8,680829	8,653343	8,929641	8,125524	7,542553	8,101936	7,047804	8,760486	7,231413
448	201489	8,487513	8,214824	7,499573	8,090161	8,74074	6,582001	7,408209	9,212658	8,513943
730	236216	15,44289	15,00473	15,16472	15,10782	14,93144	15,50422	15,20032	15,79866	16,25412
740	140456	4,842105	4,749718	5,81241	6,768801	6,30261	5,561439	5,471868	6,341218	5,695855
801	170969	14,65284	15,13456	14,95476	15,61031	15,13456	15,30316	15,91908	14,77876	14,77669
[1,2,3,4,5,6,7,8,9,10]										

Bicluster3	probeid	MCF-7	MM231	T47D	Hs578T	SKBR3	MM435s	ZR75-1	BT-549	MM453
329	231349	6,45168	7,646295	6,229041	7,080557	7,118636	6,186273	7,554411	6,510242	6,787215
371	140490	6,96837	5,711136	7,125772	6,299203	7,263096	6,178563	5,543545	7,096224	6,536871
[1,2,3,4,5,6,7,8,9,10]										

Καρκίνος του Μαστού

di_new,dj_new ->Τρίτη εκτέλεση αλγορίθμου

di_new_dj_new						
Bicluster1 probeid			MM231	MM435s	BT-549	MM453
245	219524		6,313778	7,104665	8,36813	7,289293
329	231349		7,646295	6,186273	6,510242	6,787215
363	111270		6,418241	6,835804	7,2975	7,075434
371	140490		5,711136	6,178563	7,096224	6,536871
462	214589		7,289839	6,179302	7,718731	6,661609
553	190994		7,780951	7,126229	8,510808	6,882617
678	180501		5,173406	4,637098	7,501148	5,874923
718	214717		6,44112	5,376219	6,670218	6,255565
853	163565		6,378258	7,346787	7,344067	5,988559
[2,6,8,9]						

Bicluster2	probeid		MM231	T47D	BT-549
168	142719		7,175448	6,612224	7,457644
438	181532		8,653343	8,929641	8,760486
440	101235		7,004152	6,055474	7,531619
448	201489		8,214824	7,499573	9,212658
533	146732		5,869345	7,373126	7,896661
780	109288		6,449557	9,199736	6,710841
955	234158		4,497453	6,105809	7,016225
[2,3,8]					

Bicluster3	probeid		MCF-7	T47D	Hs578T	MM453
69	151787		7,571363	6,983927	4,753833	6,488392
395	162744		9,219379	8,538543	8,513284	7,122914
602	125603		6,730458	7,083339	7,112675	7,299382
904	132754		7,238138	7,022943	7,891657	6,606438
945	115325		7,712005	8,998439	5,550856	7,611143
[1,3,4,9]						

$d_i, d_j \rightarrow$ Πρώτη εκτέλεση αλγορίθμου

[1,2,4,5,6,7,8,9]

[1,3,4,5,6,7,8]

[1,2,3,4,5,6,8,9]

Καρκίνος του τραχήλου της μήτρας
di_new,dj ->Δεύτερη εκτέλεση αλγορίθμου

di_new,dj										
Biclusteri	probeid	SW756	GH354	C-4I	Hela	C-33A	CaSki	ME-180	HT-3	SiHa
49	152249	7,2556324	7,4255726	7,0260866	7,2391865	7,1776018	6,675348	6,605961	7,1831391	7,0102056
92	116550	6,0934023	6,8134234	6,1339317	6,1966267	6,4284483	6,4643673	6,4093507	6,7781048	6,3648074
106	223310	16,761318	16,992873	16,932573	17,150353	17,092611	17,283723	17,004424	16,615897	16,4656
210	124803	9,433284	9,5914623	9,8105594	9,5093646	10,464566	10,147551	9,9777943	9,7560875	10,149089
234	221988	7,3640563	7,3476925	7,1074731	7,7761992	6,5824233	6,7726323	7,2718647	6,9926969	6,92741
273	174944	17,077839	17,06098	17,239012	16,820318	17,364681	17,239012	16,627222	17,024784	17,294955
338	109551	16,615897	16,608639	16,648786	16,862559	17,161044	17,153337	16,739178	17,637396	16,739178
551	212317	5,3460891	6,139169	6,2208492	5,8530042	6,3726268	6,5747708	6,4603028	5,9973073	5,99219
719	162602	10,050049	10,312374	9,8098455	10,361391	10,523458	10,374621	10,25464	10,202138	10,129886
720	236549	7,4425682	7,3287883	7,0360753	6,7838702	6,3687385	6,7847416	7,3366123	7,0931715	7,5441673
772	101481	5,4471093	5,9642345	6,1967716	5,8457607	6,1359576	6,0876415	6,0641783	6,1406199	5,9407602
1000	107345	7,844192	6,927098	7,8328134	7,2221772	7,3571121	6,9895693	7,5530775	6,9954964	7,5336043
[1,2,3,4,5,6,7,8,9]										

Bicluster2	probeid	SW756	GH354	C-4I	Hela	C-33A	CaSki	ME-180	HT-3	SiHa
159	165735	5,9391925	5,3989462	5,34104	5,8140107	5,1945288	5,3873598	5,6457532	5,8606492	5,5600311
623	120891	6,5295374	6,6117393	7,3545927	7,2493306	7,2700248	7,1550659	6,9694486	7,2185135	7,7782864
661	217397	6,5714892	7,2341783	7,5866759	7,486672	7,6314021	7,0574531	7,4964565	6,692654	7,7362114
789	188299	6,9468478	8,0099766	7,6111433	7,4654998	8,0572237	7,2810735	6,9269982	7,1301289	7,4889247
[1,2,3,4,5,6,7,8,9]										

Bicluster3	probeid	SW756	GH354	C-4I	Hela	C-33A	CaSki	ME-180	HT-3	SiHa
44	106630	6,4156662	6,8901907	5,8842316	7,1092478	7,1890325	6,821831	6,5629594	6,8059	6,5970332
172	118863	6,8654793	6,6613918	7,0667775	7,1037585	7,0758282	6,2684029	6,4555176	7,5674798	6,7716146
945	115325	7,1664459	7,2642105	6,4911273	6,9813359	7,7679318	7,0662762	6,5887087	7,2179296	7,5481648
988	103763	9,9362311	9,7216546	9,8782942	9,4493952	9,6243857	9,9805756	10,405236	9,6597138	9,8957957
[1,2,3,4,5,6,7,8,9]										

Καρκίνος του τραχήλου της μήτρας
di_new,dj_new ->Τρίτη εκτέλεση αλγορίθμου

di_new,dj_new							
Bicluster1	probeid	C-4I	Hela	CaSki	ME-180	HT-3	SiHa
49	152249	7,0260866	7,2391865	6,675348	6,605961	7,1831391	7,0102056
92	116550	6,1339317	6,1966267	6,4643673	6,4093507	6,7781048	6,3648074
234	221988	7,1074731	7,7761992	6,7726323	7,2718647	6,9926969	6,92741
342	168674	7,4809475	7,1709002	7,037937	7,1603067	6,9564847	6,5262645
509	191330	7,4928272	6,9077857	7,1172795	7,3937337	6,4906461	7,2274741
623	120891	7,3545927	7,2493306	7,1550659	6,9694486	7,2185135	7,7782864
661	217397	7,5866759	7,486672	7,0574531	7,4964565	6,692654	7,7362114
719	162602	9,8098455	10,361391	10,374621	10,25464	10,202138	10,129886
720	236549	7,0360753	6,7838702	6,7847416	7,3366123	7,0931715	7,5441673
772	101481	6,1967716	5,8457607	6,0876415	6,0641783	6,1406199	5,9407602
789	188299	7,6111433	7,4654998	7,2810735	6,9269982	7,1301289	7,4889247
847	156024	7,0959633	7,042771	6,6325935	7,1431925	7,3662256	6,9547596
1000	107345	7,8328134	7,2221772	6,9895693	7,5530775	6,9954964	7,5336043
[3,4,6,7,8,9]							

Bicluster2	probeid	SW756	GH354	C-33A	HT-3	SiHa
172	118863	6,8654793	6,6613918	7,0758282	7,5674798	6,7716146
245	219524	6,9933466	7,0079331	7,5889966	7,3853697	7,2401637
261	172017	6,926658	6,486689	7,3822668	7,2808587	6,4238303
272	133136	7,4919227	7,5574718	7,8574376	7,7323969	7,0680994
371	140490	6,6485941	6,481245	7,3411922	7,0022584	6,5775083
569	145804	6,9062355	7,2519258	7,7897533	6,598699	6,9973845
666	148497	7,3514019	7,6312422	7,0943583	7,2770573	6,6092993
768	192988	6,9915397	7,0320979	7,6214427	7,5304628	6,9806237
904	132754	7,6118188	7,6046634	7,6840971	7,2112607	7,709287
945	115325	7,1664459	7,2642105	7,7679318	7,2179296	7,5481648
982	223803	7,2359009	7,4773407	7,4454309	7,7696152	7,6125504
[1,2,5,8,9]						

Bicluster3					
	probeid	SW756	ME-180	HT-3	SiHa
295	233204	8,0055847	6,8137238	7,6743726	7,6191485
419	180895	7,4157829	6,9701876	7,4893019	7,6257569
596	132804	7,4968769	7,5022686	6,820095	7,3395559
602	125603	7,1482264	7,0567261	7,1587824	6,8531651
886	128078	7,7477107	7,3307188	7,7053705	6,9397308
[1,7,8,9]					

$d_i, d_j \rightarrow$ Πρώτη εκτέλεση αλγορίθμου

[1,2,3,4,5,6,9,10]

[1,3,4,5,6,7,8,10]

[3,4,5,6,7,9,10]

di_new,dj -> Δεύτερη εκτέλεση αλγορίθμου

[illegible]

di_new,dj_new -> Τρίτη εκτέλεση αλγορίθμου

Biclust3	PROBEID	TOV1-112D	A2780	OVIMZ-6	CaOV3	NIHOVCAR3
272	133136	7,0674754	7,5522646	7,071667	6,8624547	6,316743597
288	130324	6,90423372	6,9275144	7,275436	6,6212345	6,614169313
419	180895	7,99619397	8,0429035	7,2522509	6,8778556	7,467860777
569	145804	7,38660904	6,8780671	7,3857148	7,0772121	7,306752463
578	224152	7,96894867	7,7239716	7,7246916	7,7413579	6,984970036
654	125612	7,38884662	7,2248598	7,7998048	7,2645463	7,019960828
691	180199	7,6332486	7,4964565	7,5844839	7,204759	7,145835637
789	188299	7,74884783	7,6950231	8,2474016	7,6556814	7,652811016
[1,3,5,7,9]						

$d_i, d_j \rightarrow$ Πρώτη εκτέλεση αλγορίθμου

[1,2,3,4,5,7,8]

[1,2,3,4,5,6,8,9]

97

[2,3,4,5,6,7,8,9]

di_new,dj -> Δεύτερη εκτέλεση αλγορίθμου

[illegible]

$di_new, dj_new \rightarrow$ Τρίτη εκτέλεση αλγορίθμου

Bicluster3	probeid	EJ	Ishikawa	RL95-2	HEC50B	AN3CA
116	219127	7,5580813	7,4747797	6,5354105	7,352694	7,5579332
438	181532	7,9707598	7,4392251	7,7187311	7,881333	8,2323641
608	115230	8,3862721	8,0889528	7,7353895	8,1837973	7,7163123
886	128078	8,3740923	8,1413744	7,1051917	7,3554995	7,4534147
983	218294	8,2868653	7,5959366	7,5316195	7,5560139	8,1741964
[2,3,4,5,9]						