

Exploring Clustering Techniques for Microarray Gene Expression Data

*Purnendu Mishra¹, Nilamani Bhoi²

^{1,2}Dept. of Electronics and Telecommunication Engineering

^{1,2}Veer Surendra Sai University of Technology, Burla-768018, Odisha, India

¹erpurnendumishra@gmail.com, ²nilamani2003@gmail.com

Abstract

Analysis of DNA microarray has become the most commonly used valuable genomics approach in the field of bioinformatics. Clustering of genes information is a standard exploratory procedure used to distinguish firmly related genes. Clustering is fundamental in the process of data mining to uncover regular structures and distinguish intriguing patterns with regards to the basic large number of genes and the complex of biological data. Progressing research in this domain show that microarray managing will be supportive for the characterization disease genes. Various Artificial Intelligent strategies are in like way used to perceive the tumors and malady cells. Here, in this paper, distinctive proficient clustering techniques are discussed and analyzed about for the grouping of genes from the microarray gene expression dataset.

Keywords: Genomic Signal Processing, Clustering, Microarray gene expression data, Cancer Cells.

1. Introduction

Microarray is an investigation tool used to recognize the expression of thousands of genes at the same time. DNA microarrays are microscopic slides that are printed with an enormous number of little spots in described positions, with each spot containing an acknowledged DNA course of action or gene. With the advancement of DNA microarray innovation, in any case, researchers would now be able to look at how dynamic a huge number of genes are at some random time. Image analysis is a significant perspective on resulting examination, for example, image improvement, segmentation and clustering [1]. Microarray innovation will assist scientists with studying a wide range of maladies. With the assistance of microarray innovation drug improvement and creation, disease characterization, analytic turn of events, classification of cancer on the examples of gene action in the tumor cells can be conceivable. [2]

DNA microarrays are made by robotic machines that organize tiny measures of hundreds or thousands of gene successions on a solitary microscope slide as appeared in figure 1.

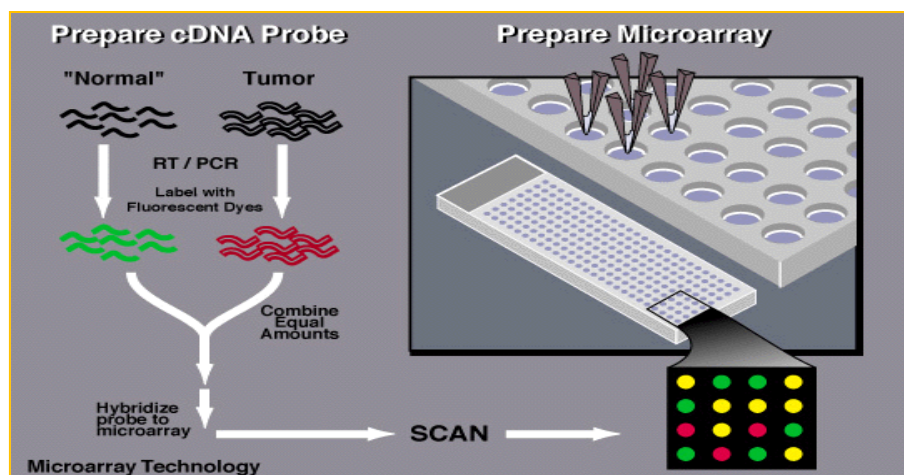


Figure 1. Microarray Image Preparation

Scientists have a database of more than a great many gene sequences that they can use for this reason. At the point when a gene is actuated, cell apparatus starts to duplicate certain fragments of that gene. The subsequent item is known as messenger RNA (mRNA), which is the body's format for making proteins. The mRNA created by the cell is integral, and thusly will tie to the first segment of the DNA strand from which it was replicated. To figure out the genes turned on and turned off in a given cell, an analyst should initially gather the messenger RNA particles present in that cell. The analyst at that point marks every mRNA particle by connecting a fluorescent color. Next, the scientist puts the named mRNA to the microarray slide. The messenger RNA which was available in the cell will at that point hybridize - or tie - to its correlative DNA on the microarray, leaving its fluorescent tag. A scientist should then utilize an uncommon scanner to quantify the fluorescent regions on the microarray. On the off chance that a specific gene is extremely dynamic, it produces numerous particles of messenger RNA, which hybridize to the DNA on the microarray and generate a bright fluorescent territory as shown in Fig 2 & Fig 3. Genes that are to some degree dynamic produce less mRNAs, which brings about dimmer fluorescent spots. In the event that there is no fluorescence demonstrating that the gene is inactive. Analysts oftentimes utilize this method to inspect the movement of different genes at various times. [2, 3]

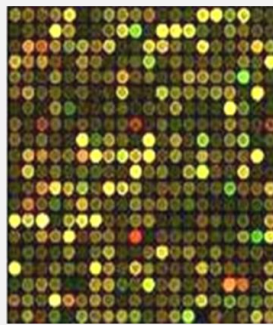


Figure 2. Microarray Image

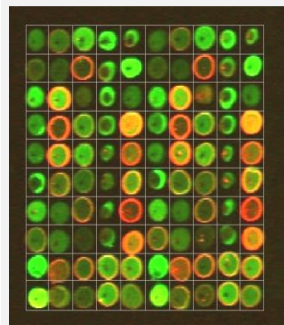
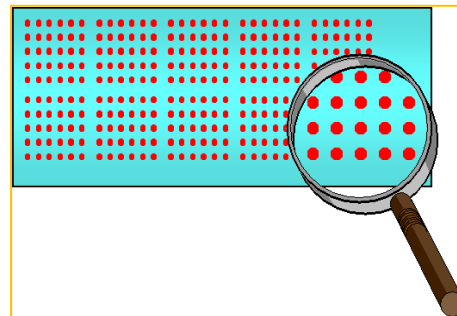


Figure 3. Meta-array



Array to place in an orderly

arrangement

Due to Microarrays simultaneously screen the gene expression of number of genes under various test conditions. Conspicuous verification of co-conveyed genes and understandable models is the prime goal in microarray or gene expression data on examination and is a noteworthy task. A microarray analyze commonly evaluates countless DNA sequences. These conditions might be a period arrangement during a biological or an assortment of various tissue tests e.g., normal versus cancerous tissues.

In microarray tests, hybridized arrangements are imaged in a microarray scanner to deliver red and green fluorescence power estimations at every one of a huge assortment of pixels which together spread the exhibit. These fluorescence powers relate to the degrees of hybridisation of the two examples to the DNA sequences spotted on the slide. Fluorescence intensities are typically put away as 16-bit pictures which we see as 'raw' data. The preparing of filtered pictures fundamentally go through three significant zones as shown in the Fig. 4. Those are denoising, segmentation and clustering [1].

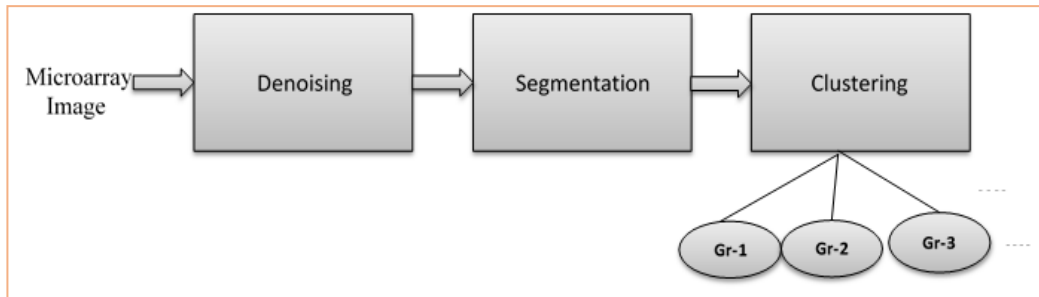


Figure 4. Block Diagram of Microarray Image Analysis

Through this paper, we will focus on the clustering investigation of gene expression data without making a differentiation among DNA arrangements, which will consistently be classified "genes." Clustering of gene expression data can be partitioned into two fundamental classifications. Gene-based clustering and sample based clustering [4]. In gene based clustering, genes are treated as objects and tests are highlights or traits for clustering. The objective of gene-based clustering is to sort differentially expressed genes and sets of genes or conditions with comparative expression patterns or profiles, and to generate a list of expression estimations.

Clustering procedures demonstrated are to be useful to comprehend gene work, gene guideline, cellular processes, and subtypes of cells. Genes with comparable expression patterns can be clustered along with comparable cell capacities. This methodology may additionally comprehension of the elements of numerous genes for which data has not been beforehand accessible. [5].

2. Effective Clustering of Microarray Gene Expression

In this segment, we will discuss about the diverse unmistakably utilized clustering approaches for microarray gene expression data like k-means clustering, Particle Swarm Optimization, Genetic Algorithm, Cuckoo search Algorithm and so forth.

2.1 K-means clustering

Nidheesh et al. [6] introduced an improved deterministic K-Means clustering calculation for forecast of disease subtype, which gives stable outcomes and which has a novel technique for choosing initial centroids. The calculation abuses the way that clusters exist at thick locales in highlight space thus, it is more reasonable to pick data centers which have a spot with thick regions and which are sufficient detached in incorporate space, as starting centroids. Being deterministic, not normal for the old style K-Means, there was no compelling reason to execute the calculation ordinarily to get a satisfactory outcome.

Jothi et al. [7] have presented a deterministic introduction algorithm for K-means (DK-means) by exploring a great deal of likely concentrations through a constrained bi- partitioning technique. This algorithm was contrasted and traditional K-implys with arbitrary initialization and improved K-means variations, for instance, K-means ++ and Min Max calculations. It was also contrasted and three deterministic introduction techniques. Exploratory assessment on gene expression datasets displays that DK-means achieves improved results with respect to faster and stable intermingling with higher cluster quality when contrasted and various algorithms.

Gasch et al. [8] used a heuristically modified version of fuzzy k-means clustering to recognize covering clusters of yeast genes dependent on distributed gene-articulation information following the reaction of yeast cells to natural changes. They have approved the strategy by distinguishing gatherings of practically related and co- regulated genes, and in the process we have revealed new connections between yeast genes and between the tests conditions dependent on likenesses in gene- expression patterns. To examine the guideline of gene expression, they associated the clusters with realized transcription factor restricting locales present in the genes' advertisers. These outcomes give experiences into the system of the guideline of gene expression in yeast cells reacting to environmental c changes. : Fuzzy k-means clustering was a helpful

explanatory tool for extricating biological bits of knowledge from gene-expression data. Their analysis introduced here recommends that a pervasive subject in the guideline of yeast gene expression was the condition-explicit co-regulation of covering sets of genes.

The Model Based Modified k-means technique was presented by Suresh et al. [9] to locate the specific number of clusters and conquer the issues in the current k-implies clustering strategy. Their test results shows that the proficiency of their technique by figuring and contrasting the entirety of squares and diverse k values.

2.2 Fuzzy C-Means Clustering

Doulaye et al. [10] proposed the Fuzzy C-means (FCM) clustering technique for microarray data analysis. Clustering analysis is basic requirement for recognizing naturally important gene groups. Partitional clustering strategies, for example, K-means or SOMs, appoint every gene to a solitary cluster. Notwithstanding, these strategies don't give data about the impact of a given gene for the general state of clusters. They have applied the strategy of fuzzy partitioning, Fuzzy C-means, to credit the participation of cluster esteems to genes.

Dhiraj K at al. [11] applied the Fuzzy C-means clustering algorithm to the microarray data. Two pattern recognition data (IRIS and WBCD information) and thirteen microarray data is utilized to assess execution of K-means and Fuzzy C-means. Broad simulation results shows that the FCM clustering algorithm had the option to give the most elevated accuracy and generalization results contrasted with K-means clustering algorithm.

2.3 Particle Swarm Optimization (PSO)

A hybrid clustering approach that depended on Self Organizing Maps and Particle Swarm Optimization was proposed by Xiao et al. [12]. In their technique, the rate of convergence was improved by adding a conscience factor to the Self-Organizing Maps algorithm. The strength of the outcome was estimated by utilizing a resampling method. The algorithm was actualized on a cluster of workstations.

Sun et al. [13] proposed another plan for clustering gene expression dependent on an altered rendition of Quantum-acted PSO (QPSO) algorithm, known as the Multi-Elitist QPSO (MEQPSO) model. The clustering strategy likewise utilizes a one-step K-means operator to adequately quicken the convergence speed of the algorithm. The MEQPSO algorithm was tried and contrasted and some other as of late proposed PSO and QPSO variations on a set-up of benchmark capacities. Their outcomes demonstrate that MEQPSO was a leading strategy for gene clustering.

Lam et al. [14] executed an upgraded cluster matching (CM2) for improvement the nature of original cluster matching (CM) with streamlined structure. Thusly, it will likewise improve PSO-based K-Means clustering algorithm (PSO-KM). The CM2 was utilized to recommend that the arrangement of cluster centroids in every particle position can be very much coordinated with the relating ones in the worldwide best molecule, with the nearest separation. It brings about ideal sequencing of groups in every particle so as to well get ready for next PSO updating cycle, which was correspondence measure between particles to discover next better position. Three famous genuine datasets were chosen to assess the proposed strategy with other related algorithms. Trial results show that PSO-KM with the proposed scheme CM2 can beat the first form, PSO-KM with CM, and different strategies regarding compactness within clusters. Particularly, their plan makes PSO-KM to converge faster with great compactness with the fresher strategy PK-Means.

In light of the recently demonstrated Quantum-carried on Particle Swarm Optimization (QPSO) algorithm, Chen et al. [15] focused in on the utilization of QPSO in gene expression data clustering which can be decreased to an optimization issue. Their clustering algorithm parcels the N examples of the gene expression dataset into user-defined K categories to limit the fitness function of Total Within-Cluster Variation. Along these lines a partition with superior was acquired. The experiment results on four gene expression data sets show that their QPSO-based clustering algorithm have been compelling and promising tool for gene expression data investigation.

Dey et al. [16] proposed an improved clustering strategy utilizing PSO based K-means clustering which gives preferable accuracy over other clustering algorithms in microarray data clustering. There are two kinds of gene expression data clustering measure (a) gene-based clustering where genes were clustered taking samples as highlights for example sample size was steady and (b) sample based clustering where samples were grouped accepting genes as highlights. The sample based clustering technique was explained as single-objective PSO based K-means algorithm was thought of. It is conceivable to generate multi objective PSO based K-means clustering algorithm which can cluster the two genes and samples all the while for gene expression data.

Deng et al. [17] presented a K-means clustering algorithm subject to particle swarm optimization (PSO K-means) for microarray data clustering. Their algorithm discovers clusters in microarray data without prior information on achievable cluster numbers or complex parameter settings, which were required by other clustering procedures. PSO K-means had the alternative to draw out the latent structure of microarray data sets. A huge decrease in inside cluster blunder was procured using the technique, without noteworthy diminishing in the middle of between cluster distinctions.

2.4 Genetic Algorithm

Chakraborty et al. [18] presented two genetic algorithms which utilize greedy algorithm as the nearby pursuit methodology. The outcomes for yeast and lymphoma datasets show that these algorithms have performed superior to different methodologies. In spite of the fact that the reason behind choosing the fitness functions of the two genetic algorithms have been distinctive both the algorithms have performed similarly well regarding quality of biclusters and their arrangement with known biological outcomes. The main advantage of the method was that both the algorithms don't need threshold score as input parameter, subsequently disposing of the trouble of figuring the edge for each information data. The tables for the comparison of various algorithms plainly show that utilizing the genetic algorithms they can find bigger bi-clusters i.e., containing more genes and conditions with smaller residue score. This infers they can find the best bi-clusters utilizing the genetic algorithm on the grounds that the smaller the residue and/or as the bigger the volume, the better was the nature of the bi-cluster.

A hybrid GA (genetic algorithm)- based clustering (HGACCLUS) schema, joining benefits of the Simulated Annealing, was depicted by Pan et al. [19] for finding an optimal or near-optimal set of medoids. Their schema boosted the clustering accomplishment by accomplishing inside cluster attachment and outside cluster segregation. The exhibition of HGACCLUS and different techniques was analyzed by utilizing reenacted information and open microarray gene-expression datasets. HGACCLUS was generally discovered to be more precise and strong than different strategies talked about by the specific approval methodology and the unequivocal cluster number.

2.5 Cuckoo search with crossover (CSC)

Balamurugan et al. [20] found the essential biclusters in enormous expression data using rearranged cuckoo search with Nelder–Mead (SCS-NM). The widening and acceleration of the pursuit space were gained through improving and simplex NM individually. Their model was taken a stab at four benchmark datasets, and the results are contrasted and the swarm intelligence techniques and the distinctive biclustering algorithms. The results show that there was basic improvement in the health assessment of SCS-NM. Also, chooses the organic significance of the biclusters with Gene Ontology to the extent function, process and component.

Sampathkumar et al. [21] displayed a enhanced bio- inspired algorithm explicitly cuckoo search with crossover (CSC) for picking genes from development of microarray that had the alternative to aggregate different disease sub-types with high precision. The assessment results were finished with five benchmark cancer gene expression datasets. The results portray that CSC was beats than CS and other eminent procedures.

Table 1 shows the effective clustering methods for microarray gene expression.

Table 1: Effective Clustering of Microarray Gene Expression

Author	Objective	Clustering Techniques	Significance
Nidheesh et al [6]	An modified K-Means clustering algorithm for cancer subtype prediction	K-Means clustering algorithm	There was no need to run the algorithm number of times for finding the suitable result
Jothi et al. [7]	A deterministic k-means clustering algorithm for gene expression analysis	K-Means clustering algorithm	Results improved in terms of faster and stable convergence, and better cluster quality
Suresh et al. [9]	Model based modified k-means clustering for microarray data	K-Means clustering algorithm	Efficiency of their method by calculating and comparing the sum of squares with different k values
Xiao et al. [12]	Gene clustering using self-organizing maps and particle swarm optimization	PSO	The rate of convergence was improved by adding a conscience factor
Sun et al. [13]	Gene expression data analysis with the clustering method based on an improved quantum-behaved PSO	PSO	Their results shows that MEQPSO clustering algorithm was an efficient technique and can be widely used for gene clustering.
Lam et al. [14]	PSO-based K-Means clustering with enhanced cluster matching for gene expression data	PSO with K-means	Their scheme makes PSO-KM to converge faster with good compactness
Chakraborty et al. [18]	Biclustering of gene expression data using genetic algorithm.	Genetic Algorithm	Discovered the best biclusters using the genetic algorithm
Balamurugan et al. [20]	A hybrid cuckoo search algorithm for biclustering of microarray gene-expression data	Cuckoo search	There was significant improvement in the fitness value

3. Cluster Validation for Microarray Gene Expression

Datta et al. [22] thought about six clustering algorithms (of different flavors) and assess their exhibitions on a notable openly accessible microarray data set on sporulation of sprouting yeast and on two recreated data set. In addition to other things, they define three sensible approval methodologies that can be utilized with any clustering algorithm when transient perceptions or replications were available. They assess every one of these six clustering strategies with these validation measures. While the 'best' technique was reliant on the specific approval methodology

and the quantity of clusters to be utilized, generally speaking Diana seems, by all accounts, to be a strong performer. Strikingly, the performance of correlation-based hierarchical clustering and model-based clustering (another strategy that has been upheld by various analysts) seem, by all accounts, to be on inverse boundaries, contingent upon what validation measure one utilizes. Next it was demonstrated that the group implies delivered by Diana were the nearest and those created by UPGMA were the farthest from a model profile dependent on a lot of hand-picked genes.

Bolshakova et al. [23] presented a cluster validation method for gene expression information. Machaon Clustering and Validation Environment (CVE) system hopes to partition genes into bunches portrayed by comparable expression patterns, and to locate the nature of the groups obtained. The Machaon Clustering and Validation Environment (Machaon CVE) was a cross-stage Java-based apparatus, which offers various clustering and validity strategies for DNA microarray data assessment. It centers: (a) to segment genes into bunches arranged by comparative expression models and (b) to assess the validity of the bunch found.

Bolshakova et al. [24] arranged a connection of the data based and Gene Ontology (GO)-based approaches to manage cluster approval methods for gene microarray examination. They apply a homogeneous method to manage getting measurements from different GO-based similarity measures and a standardization of validation index esteems that grants to compare them with each other similarly as to databased approval lists. The results show robust connection between both GO-based and data based validation indices. The results recommend that may address a practical device to help biomedical data exposure tasks subject to gene expression data.

Jaskowiak et al. [25] explored the decision of nearness measures for the clustering of microarray data by assessing the exhibition of 16 vicinity measures in 52 datasets from time-course and disease tests. Their outcomes uphold that measures once in a while utilized in the gene expression writing can give preferred outcomes over generally utilized ones, for example, Pearson, Spearman and, Euclidean separation. Given that various apportionments represented time-course and disease data assessments, their decision ought to be explicit to every situation. To assess gauges on time-course data they pre-handled and ordered 17 datasets from the microarray writing in a benchmark alongside another strategy, called Intrinsic Biological Separation Ability (IBSA). Both can be utilized in future exploration to survey the viability of new measures for gene time-course data.

A Hybrid Microarray Clustering Algorithm was proposed by Ghalib et al. [26] which was consolidated with Hubert's Statistic Technique, Jaccard's coefficient and Dunn's Index utilized for cluster Validation. Its fundamental point was to improve the proficiency level of the nature of clusters, with optimized validation and diminish the memory prerequisites lower than practically all the current clustering algorithms. It likewise controls in accomplishing quality clusters.

Costa et al. [27] played out a similar investigation of clustering strategies which is utilized for analysis of gene expression time courses. Five different clustering techniques found in the writing of gene expression analysis were thought about: agglomerative various leveled clustering, CLICK, dynamical clustering, k-means and SOM. To assess the strategies, a k-fold cross-validation techniques was applied. The accuracy of the outcomes was surveyed with the correlation of the segments got in this investigations.

4. Genomic Signal Processing of Microarrays for Cancer Gene Expression

To reveal definite genetic mechanisms in hepatocellular carcinoma (HCC) with a view toward improvement of novel restorative targets, Okabe et al. [28] dissected expression profiles of 20 essential HCCs and their comparing noncancerous tissues by methods for cDNA microarrays comprising of 23,040 genes. Up-guideline of mitosis- promoting genes was seen in most of the tumors analyzed. A few genes indicated expression designs in hepatitis B infection positive HCCs that were not the same as those in hepatitis C infection positive HCCs; a large portion of them encoded compounds that utilize cancer-causing agents or potentially anticancer specialists. Besides, they distinguished various genes related with dangerous histological sort or obtrusive phenotype. Collection of such data will make it conceivable to characterize the idea of

individual tumors, to give pieces of information to recognizing new restorative targets, and eventually to streamline treatment of every patient.

Zembutsu et al. [29] used a cDNA microarray speaking to 23,040 genes to investigate expression profiles in a board of 85 cancer xenografts got from nine human organs. The xenografts, embedded into naked mice, were inspected for affectability to nine anticancer medications (5-fluorouracil, 3-[(4-amino-2-methyl-5-pyrimidinyl) methyl]-1-(2-chloroethyl)- 1-nitrosourea hydrochloride, Adriamycin, cyclophosphamide, cisplatin, mitomycin C, methotrexate, vincristine, and vinblastine). Examination of the gene expression profiles of the tumors with sensitivities to each medication recognized 1,578 genes whose expression levels associated fundamentally with chemo affectability; 333 of those genes indicated noteworthy relationship with at least two medications, and 32 corresponded with six or seven medications. These data should contribute helpful data for recognizing prescient markers for drug affectability that may in the end give "customized chemotherapy" for singular patients, just as for advancement of novel medications to beat procured obstruction of tumor cells to chemical specialists.

Subramanian et al. [30] presented a unimaginable scientific method called Gene Set Enrichment Analysis (GSEA) for understanding gene expression data. The procedure decides its ability by focusing in on gene sets, that was, social affairs of genes that share fundamental biological function, chromosomal region, or guideline. They show how GSEA yields encounters into a couple of malignant growth related data sets like leukemia, cellular breakdown in the lungs. Strikingly, where single-gene analysis finds little similarity between two autonomous examinations in endurance of patient in cellular breakdown in the lungs, GSEA reveals various natural pathways in like way. The GSEA procedure was exemplified in a straightforwardly open programming bundle, alongside a starter database of 1,325 organic gene sets.

The cDNA microarray procedure was an as of late created apparatus that abuses abundance of data for the analysis of gene expression. In cDNA microarray strategy, DNA tests speaking to cDNA probes were displayed onto a glass slide and examined with fluorescently marked cDNA targets. The intensity of the innovation was the capacity to play out a genome-wide expression profile of thousands of genes in a single examination. Khan et al. [31] depicted the standards of the microarray innovation as applied to disease research, sum up the writing on its utilization up until now, and guess on the future use of this incredible strategy.

Khan et al. [32] used cDNA microarrays containing 1238 cDNAs to explore the gene expression profile of a gathering of seven alveolar rhabdomyosarcoma (ARMS) cell lines portrayed by the presence of the PAX3-FKHR combination gene. Utilizing the technique for multidimensional scaling to speak to the connections among the cell lines in two-dimensional Euclidean space, they established that ARMS cells show a steady pattern of gene expression, which permits the cells to be clustered together. These outcomes in ARMS show the capability of cDNA microarray innovation to clarify tumor-explicit gene expression profiles in human malignant growths.

Gardina et al. [33] examined 20 matched tumor-typical colon malignancy tests utilizing a microarray intended to distinguish more than 1,000,000 putative exons that can be practically collected into potential gene-level transcripts as indicated by different degrees of earlier supporting proof. Analysis of high certainty (observationally upheld) transcripts recognized 160 differentially communicated genes, with 42 genes possessing an organization affecting cell expansion and another 29 genes with obscure capacities. A more theoretical analysis, including transcripts dependent on computational forecast, created another 160 differentially communicated genes, three-fourths of which have no past comment. They likewise present an examination of gene signal assessments from the Exon 1.0 ST and the U133 in addition to 2.0 exhibits. Novel splicing occasions were anticipated by trial algorithms that look at the general commitment of every exon to the related transcript force in each tissue. Top scoring up-and-comers from our analysis were additionally found to generously cover with EST-based bioinformatics forecasts of elective splicing in cancer.

Takata et al. [34] investigated the gene expression profiles of biopsy materials from 27 intrusive bladder tumors utilizing a cDNA microarray comprising of 27,648 genes, after populaces of malignant growth cells had been cleansed by laser laser micro beam micro

dissection. They recognized many genes that were expressed distinctively between nine 'responder' and nine "non-responder" tumors; from that rundown we chose the 14 "prescient" genes that demonstrated the most critical contrasts and formulated a mathematical forecast scoring framework that obviously isolated the responder group from the non-responder gathering. That framework precisely anticipated the medication reactions of 8 of 9 experiments that were held from the first 27 cases. Since continuous reverse transcription ^ PCR data were profoundly concordant with the cDNA microarray information for those 14 genes, we built up a quantitative reverse transcription ^ PCR ^ based prediction framework that could be possible for clinical use.

Clarke et al. [35] investigated advancements in gene expression microarray and show the advancement and capability of the procedure in malignancy science, pharmacology, and medication improvement. Significant applications include: (a) advancement of a more worldwide comprehension of the gene expression irregularities that add to cancer; (b) disclosure of new symptomatic and prognostic indicators and biomarkers of remedial reaction; (c) recognizable proof and validation of new molecular focuses for drug improvement; (d) arrangement of an improved comprehension of the sub-atomic method of activity during lead ID and enhancement, including structure-action connections for on track versus off-target impacts; (e) forecast of expected results during preclinical turn of events and toxicology contemplates; (f) affirmation of a sub-atomic method of activity during theory testing clinical preliminaries; (g) Identification of genes engaged with giving medication affectability and obstruction; and (h) expectation of patients destined to profit by the medication and use in general pharmacogenomics examines. Because of additional mechanical enhancements and diminishing costs, the utilization of microarrays have become a basic and conceivably routine instrument for malignant growth and biomedical research.

Hisaminato et al. [36] played out a complete analysis of the expression profiles in 25 grown-up and 4 fetal human tissues by methods for a cDNA microarray comprising of 23,040 human genes. Their investigation uncovered various genes that were communicated explicitly in every one of those tissues. Among the 29 tissues analyzed, 4,080 genes were profoundly communicated (in any event a five-overlap expression proportion) in one or just a couple of tissues and 1,163 of those were communicated only (in excess of a ten times higher expression proportion) in a specific tissue. Expression of a portion of the genes in the last classification was affirmed by northern analysis. A hierarchical clustering analysis of gene-expression profiles in nerve tissues (grown-up cerebrum, fetal mind, and spinal rope), lymphoid tissues (bone marrow, thymus, spleen, and lymph hub), muscle tissues (heart and skeletal muscle), or fat tissues (mesenteric fat and mammary organ) recognized a lot of genes that were normally communicated among related tissues. These information should give valuable data to clinical exploration, particularly for endeavors to distinguish tissue-specific molecules as expected focuses of novel medications to treat human infections.

Cho et al. [37] attempted to research various highlights and classifiers using three benchmark datasets to evaluate deliberately the exhibitions of the component choice procedures and AI classifiers. Three benchmark datasets were Leukemia dataset, Colon dataset and Lymphoma disease dataset. Pearson's and Spearman's correlation coefficients, cosine coefficient, sign to commotion proportion Euclidean separation, shared data, data gain and so forth have been utilized for highlight choice. k-nearest neighbour, Multi-layer perceptron, structure versatile self-sorting out guide, uphold vector machine and so on have been used for characterization. Also, they have joined the classifiers to better the introduction of characterization. Exploratory results implies that the gathering with a couple of reason classifiers conveys the best pace of acknowledgment on the benchmark dataset.

Segal et al. [38] demonstrated that dynamic imaging qualities in non- computed tomography (CT) intentionally correspond with the worldwide gene expression projects of primary human liver malignancy. Blends of 28 imaging credits can reproduce 78% of the worldwide gene expression profiles, revealing cell extension, liver engineered capacity, and patient perception. Henceforth, genomic action of human liver tumors can be decoded by non-intrusive imaging, therefore enabling non-obtrusive, successive and continuous molecular profiling for modified prescription.

MicroRNAs (miRNAs) were a class of small noncoding RNAs that control gene expression by focusing in on mRNAs and setting off either translation limitation or RNA corruption. Their twisted expression may be related with human infections, including malignancy. Indeed, miRNA mutilated expression has been recently found in human constant lymphocytic leukemias, where miRNA marks were connected with unequivocal clinicobiological features. Here, Iorio et al. [39] differentiated and ordinary breast cancer, miRNAs were furthermore bizarrely imparted in human breast disease. The general miRNA expression could obviously separate ordinary versus disease tissues, with the most out and out freed miRNAs being mir-125b, mir-145, mir-21, and mir-155. Results were confined by microarray and Northern smudge examinations. They perceive miRNAs whose expression was correspond with bosom disease biopathologic highlights, for instance, estrogen and progesterone receptor expression, tumor stage, vascular attack, or proliferation index.

Watanabe et al. [40] recognized a lot of separating genes that can be utilized for portrayal and prediction of reaction to radiotherapy in rectal disease. 52 rectal malignant growth patients who went through preoperative radiotherapy were contemplated. Biopsy specimens were acquired from rectal malignancy before preoperative radiotherapy. Reaction to radiotherapy was controlled by histopathologic assessment of surgically resected specimens and delegated responders or non-responders. By deciding gene expression profiles utilizing human U95Av2 Gene Chip, they recognized 33 novel separating genes of which the expression varied essentially among responders and non-responders. Utilizing gene set, they had the option to build up another model to foresee reaction to radiotherapy in rectal malignancy with an accuracy of 82.4%. The rundown of separating genes included growth factor, apoptosis, cell expansion, signal transduction, or cell adhesion-related genes. They recommended the likelihood that gene expression profiling might be valuable in foreseeing reaction to radiotherapy to build up an individualized custom-made treatment for rectal malignancy. Worldwide expression profiles of responders and non-responders may give bits of knowledge into the improvement of novel helpful targets.

Sgroi et al. [41] uncovered the joined usage of laser catch microdissection and high-throughput cDNA microarrays to screen in vivo gene expression levels in purified common, prominent, and metastatic breast cell populaces from a singular patient. These in vivo gene expression profiles were affirmed by ongoing quantitative PCR and immunohistochemistry. The merged usage of laser catch microdissection and cDNA microarray assessment gives a better approach to manage explain the in vivo sub-atomic occasions sub-atomic occasions encompassing the turn of events and movement of breast cancer disease and is ordinarily relevant to the examination of malignancy.

5. Micro Array Data Clustering and Classification

Asyali et al. [42] applied fluffy c-implies (FCM) and the normal mixture modeling (NMM) based classification procedures to separate microarray data into strong and deceitful signal intensity populaces. They contrasted the eventual outcomes of FCM classification and the NMM classification wards. The two strategies were approved against the given organic data sets which comprises of true positives and true negatives. They found that the two procedures performed correspondingly well with respect to sensitivity and specificity. In spite of the fact that the correlation of the calculation times shows that the fuzzy methodology was computationally more proficient, various contemplations uphold the usage of NMM in the respect of dependability proportions of microarray data.

Two assessment measurements of attribution accuracy were used by Ouyang et al. [43]. First and foremost, the RMS botch evaluates the differentiation between the true values and the attributed qualities. Furthermore, the amount of mis-clustered genes measures the variance with the clustering true values and that with credited qualities; it takes a gander at the tendency familiar by attribution with clustering. The Gaussian blend clustering with model averaging ascription was superior to all other attribution methods, as demonstrated by both assessment measurements, on both time-series(correlated) and non-time arrangement (uncorrelated) datasets.

Dembale et al. [44] applied a fuzzy partitioning strategy, Fuzzy C- means, to ascribe membership values of cluster to genes. A significant issue for placing the FCM strategy for

clustering a the decision of the fuzziness parameter. They propose a practical procedure, in view of the circulation of separations between genes in a given data set, to decide a satisfactory incentive for m . Utilizing a yeast cell cycle data set for instance, they show that the determination builds the general biological significance of the genes inside the cluster.

Maulik et al. [45] introduced a real coded Simulated Annealing (VSA) based fuzzy clustering technique with variable length design was created and joined with mainstream Artificial Neural Network (ANN) based classifier. The thought was to refine the clustering created by VSA utilizing ANN classifier to get improved clustering execution. The proposed strategy was utilized to cluster three freely accessible genuine microarray data sets. The predominant presentation of the proposed method has been shown by contrasting and some broadly utilized existing algorithms. Likewise measurable centrality test has been led to set up the factual hugeness of the prevalent presentation of the proposed clustering calculation. At last organic pertinence of the clustering arrangements were built up.

Istepanian et al. [46] introduced a similar examination of two genomic signal preparing strategies to be specific Linear Predictive Coding, and Discrete Wavelet Decomposition coefficients for vigorous microarray data clustering. Vector quantization was applied to the resultant coefficients to give the clustering of the data tests. The two methods were validated for a standard data set. Relative investigations of the outcomes show that these strategies give improved clustering exactness contrasted with some customary clustering methods. Besides, there classifiers don't need any earlier preparing techniques.

The proposed algorithm, CKNN ascribe, was an augmentation of k closest neighbor attribution with nearby data clustering being joined by Keerin et al. [47] for improved quality and proficiency. Gene expression data was normally spoken to as a matrix whose lines and segments relate to genes and examinations, individually. CKNN commences by finding a total dataset by means of the evacuation of columns with missing value(s). At that point, k clusters and their comparing centroids were gotten by applying a clustering method on the total dataset. A lot of comparable genes of the objective gene (with missing qualities) were those having a place with the cluster, whose centroid was the nearest the objective. Having known this, the objective gene was credited by applying Clustering k nearest neighbor (CKNN) technique with comparative genes recently decided.

Liu et al. [48] diverged from NMF and PCA for reduction of dimension. The diminished data was used on 11 genuine gene expression datasets for representation, and clustering examination through k -means. They apply NMF and PCA for decrease in perception before the clustering examination. The results on one leukemia dataset suggests that NMF may discover common clusters and unquestionably recognize one mislabeled test while PCA can't. NMF most normally beats PCA for clustering examination through k -means. Their results display the pervasiveness of NMF than PCA in reducing the microarray data.

Qu et al. [49] indicated the usage of the model-based calculation in supervised clustering of microarray data. They have applied the proposed procedures to reenacted data just as to real gene expression data. They exhibited that the administered model-based algorithm was much better over the unsupervised method just as over the support vector machines (SVM) procedure.

Tibshirani et al. [50] surveyed different techniques for clustering, and delineate how they can be utilized to organize both the genes and cell lines from a lot of DNA microarray tests. The strategies examined were global clustering methods including hierarchical, K -means. At long last, they propose another technique for recognizing structure in subsets of the two genes and cell lines that were conceivably darkened by the global clustering draws near.

Another clustering algorithm, named as fuzzy-rough supervised attribute clustering (FRSAC) was proposed by Maji et al. [51] to discover such groups of genes. The proposed algorithm depended on the hypothesis of fuzzy-rough sets, which straightforwardly fuses the data of test classifications into the gene clustering measure. Another quantitative measure was presented dependent on fuzzy-rough sets that consolidates the data of test classes to quantify the closeness among genes. The proposed algorithm depended on estimating the similitude between genes utilizing the new quantitative measure, whereby excess among the genes was eliminated. The clusters were refined gradually dependent on test classes. The adequacy of the proposed FRSAC algorithm, alongside an examination with existing directed and solo gene determination

and clustering algorithms, was exhibited on six disease and two joint pain data sets dependent on the class distinctness list and prescient precision of the guileless Bayes' classifier, the K-nearest neighbor rule, and the support vector machine.

Hengpraprom et al. [52] introduced a technique for choosing instructive highlights utilizing K-Means clustering and SNR positioning. The exhibition of the proposed technique was tried on malignant growth classification issues. Genetic Programming was utilized as a classifier. The trial results demonstrate that the proposed technique yields higher accuracy than utilizing the SNR positioning alone and higher than utilizing the entirety of the genes in classification. The clustering step guarantees that the chose genes have low excess, subsequently the classifier can misuse these highlights to acquire better execution.

Díaz-Uriarte et al. [53] examined the utilization of random forest for classification of microarray data (counting multi-class issues) and propose another strategy for gene determination in classification issues dependent on random forest. Utilizing reproduced and nine microarray data sets they show that random forest has equivalent execution to other classification techniques, including DLDA, KNN, and SVM, and that the new gene determination system yields exceptionally little arrangements of genes (frequently littler than elective strategies) while saving prescient accuracy.

P.Mishra et al.[54] proposed Microarray Filtering-Based Fuzzy C-Means Clustering and Classification in Genomic Signal Processing. Proposed Kalman filter based fuzzy c-means cluster and artificial neural network (KF-FANN) update the genomic signal processing to the ideal level. The dissected result shows that the proposed technique is a capable strategy for the classification of microarray data with regards to the existing approaches.

P.Mishra et al. [55] proposed Genomic signal getting ready of microarrays for cancer gene expression and acknowledgment using cluster- fuzzy adaptive networking. The outcome demonstrated that the proposed strategy can give suitable and ideal plan for identification of microarray disease gene than the ordinary techniques, individually.

Table 2 shows the micro array data clustering and classification techniques.

Table 2: Micro Array Data Clustering and Classification Techniques

Author	Objective	Clustering	Classification	Significance
Asyali et al. [42]	Reliability investigation of microarray data utilizing fuzzy c-means and normal mixture modeling based classification strategies	Fuzzy c-means clustering	Normal mixture modeling classification	Fuzzy approach was more efficient
Ouyang et al. [43]	Gaussian mixture clustering & imputation of microarray data	Gaussian mixture clustering	Classification Expectation–Maximization algorithm	Less complexity and less time
Dembele et al. [44]	Fuzzy C-means technique for clustering microarray data. bioinformatics	Fuzzy C-means	-	Rises the biological essentialness of the genes inside the cluster

Maulik et al. [45]	VSA based fuzzy clustering method with variable length	Fuzzy clustering method	Artificial Neural Network (ANN)	Found Improvement in Clustering performance
Istepanian et al. [46]	Linear predictive coding and wavelet decomposition for robust microarray data clustering	Linear Predictive Coding and Discrete Wavelet Decomposition coefficients	-	The classifiers donot require any past training processes
Keerin et al. [47]	Cluster-based KNN missing value imputation for DNA microarray data	Clustering k nearest neighbour (CKNN)	-	Centroid was closest to the cluster
Liu et al. [48]	Deduction microarray data through nonnegative matrix factorization for visualization & clustering analysis	k-means	-	Microarray data was found reduced
Qu et al. [49]	Supervised cluster data Investigation for microarray data basing upon multivariate Gaussian mixture	-	support vector machines	High accuracy
Tibshirani et al. [50]	Different Clustering techniques for the analysis of DNA microarray data	Hierarchical, K-means, and block clustering, and tree-structured vector quantization	-	Low complexity and high efficiency found
Maji et al. [51]	Fuzzy-rough supervised attribute clustering algorithm & classification of microarray data	Fuzzy-rough supervised attribute clustering (FRSAC)	Naive Bayes' classifier, the K-nearest neighbor rule, and the support vector machine.	Predictive accuracy was high
Hengpraprom et al. [52]	Selection of Informative Genes from Microarray Data for Cancer Classification with the classifier Genetic Programming by Using the K-Means Clustering & SNR Ranking	K-Means clustering and SNR ranking	Genetic Programming	Low redundancy and better performance

Mishra,P., et al . [54]	Filtering-Based Fuzzy C-Means Clustering & Classification in Genomic Signal Processing	Kalman filter-based fuzzy c-means Clustering	Artificial Neural Network (ANN)	KF-FANN is found as an efficient method for the purpose of classification of microarray data as compared to the existing various approaches in the genomic signal processing.
Mishra,P., et al [55]	Genomic signal processing of microarrays for cancer gene expression and identification using cluster-fuzzy adaptive networking	Kalman filter-based grid density-based clustering	Adaptive neuro fuzzy interference system (ANFIS)	This technique can get effective and optimal classification and identification of microarray cancer genes than the conventional techniques.

6. Filtering with Genomic Data Clustering and Classification

Wang et al. [56] have observationally assess its presentation on three distributed malignancy classification data sets. They utilize the straight SVM and the k-NN as classifiers in the trials, and contrast the exhibition of Relief-F and other element sifting strategies, including Information Gain, Gain Ratio, and χ^2 - statistic Utilizing the leave-one-out cross validation, exploratory outcomes show that the exhibition of Relief-F was tantamount with different strategies

Huerta et al. [57] planned a hybrid framework for gene determination and classification of DNA microarray data. Right off the bat, five conventional measurable techniques were consolidated for primer gene choice (Multiple Fusion Filter). At that point diverse important gene subsets were chosen by utilizing an inserted technique that utilizes a Genetic Algorithm (GA), with a Tabu Search (TS) and Support Vector Machine (SVM). A gene subset, comprising of the most pertinent genes was gotten from this process, by dissecting the recurrence of every gene in the diverse gene subsets. At last, the most continuous genes were assessed by the inserted way to deal with get a last applicable small gene subset with superior. The proposed technique was tried in four DNA microarray datasets. From reproduction study, it was seen that their methodology works in a way that is better than different techniques.

Yang et al. [58] portrayed an improved hybrid framework for gene choice. It depended on an as of late proposed genetic ensemble (GE) framework. To upgrade the generalization property of the chose genes or gene subsets and to conquer the over fitting issue of the GE framework, they formulated a mapping technique to combine the decency data of every gene gave by numerous filtering algorithms. That data was then utilized for initialization and transformation activity of the ensemble framework.

Canedo et al. [59] proposed another system for highlight choice including a group of filters and classifiers. Five channels, taking into account different measurements, were used. Each filter picks a substitute subset of highlights which was used to prepare and to test a specific classifier. The yields of these five classifiers were combined by straightforward voting. Three prominent classifiers were used for the portrayal task: C4.5, naïve Bayes and IB1. The legitimization of the

grouping was to diminish the fluctuation of the highlights chose by channels in different classification regions. Its ampleness was displayed by using 10 microarray data sets.

Leung et al. [60] introduced a multiple-filter multiple-wrapper (MFMW) technique that utilizes numerous channels and various coverings to improve the accuracy and robustness of the classification, and to recognize potential biomarker genes. Investigations dependent on six benchmark data sets show that the MFMW approach outflanks SFSW models (generated by all blends of channels and coverings utilized in the relating MFMW model) in all cases and for every one of the six data sets. Some of MFMW-chose genes have been affirmed to be biomarkers or add to the advancement of specific malignancies by different examinations.

Tritchler et al. [61] introduced a secluded models for demonstrating the network structure in order to examine the overall effects of different filtering methods. They show that cluster examination and principal parts were insistently impacted by filtering. filtering methods proposed expressly for cluster and network examination were introduced and differentiated by reenacting secluded networks and known factual properties. To focus more pragmatic conditions, they dismember recreated "real" data dependent on particularly depicted *E. coli* and *S. cerevisiae* regulatory networks.

To quantitatively evaluate the quality of microarray tests, Miller et al. [62] straightforwardly contrast RNA-Seq with Agilent microarrays by preparing 231 exceptional samples from the Allen Human Brain Atlas utilizing RNA-Seq. The two procedures give profoundly steady, exceptionally reproducible gene expression estimations in grown-up human mind, with RNA-Seq marginally beating microarray results generally. They show that RNA-Seq can be utilized as ground truth to survey the dependability of most microarray tests, eliminate tests with askew impacts, and scale test forces to coordinate the expression levels distinguished by RNA-Seq. These sequencing scaled microarray intensities (SSMIs) give more reliable, quantitative evaluations of outright expression levels for some genes when contrasted and unscaled intensities. At long last, they approve that bring about two human cell lines, indicating that linear scaling factors can be applied across tests utilizing the equivalent microarray platform.

FiGS was an online workbench that thusly breaks down about different gene decision systems and gives the ideal gene determination result by Hwang et al [63] for a data microarray dataset. FiGS makes orchestrated gene assurance approach by changing indisputable segment decision strategies and classifiers. Close by the exceptionally accepted techniques, FiGS broadens the gene assurance methodologies by joining gene clustering choices in the component decision advance and data pre- handling decisions of choices in classifier getting ready advance.

Ke et al. [64] a score-based criteria fusion (SCF) highlight determination method was proposed for malignant growth expectation, and this procedure targets improving the forecast presentation of the classification model. The SCF procedure was evaluated on five open gene microarray datasets and three low dimensional datasets, and it shows better execution over some striking component choice methodologies while using two classifiers SVM and KNN for estimating the nature of chose highlights. Examinations watch that SCF had the choice to find more discriminative highlights than the fighting procedures and can be used as a pre-handling algorithm to join with various different strategies effectively.

A filter algorithm utilizing F-measure has been utilized with include excess expulsion dependent on the Kolmogorov-Smirnov (KS) test for rough equality of factual circulations by Biesiada et al. [65]. Therefore computationally productive K-S Correlation Based Selection algorithm has been created and tried on three high-dimensional microarray datasets utilizing four sorts of classifiers. Results were very promising and a few enhancements were proposed.

Fleury et al [66] presented a technique for identifying emphatically monotone evolutionary patterns of gene expression from a temporal sequence of microarray data. They perform gene filtering through multi-target optimization to uncover genes which have the properties of: solid monotonic increment, high end-to-end slope and low slope deviation. Both a worldwide Pareto optimization and a pair-wise local Pareto optimization are researched. Gene filtering strategy was explained on mouse retinal genes gained at various focuses over the lifetimes of a populace of mice.

7. Conclusion

Here, we have included the effective clustering of microarray gene expression data, microarray filtering base clustering and classification techniques, filtering based clustering and classification for microarray data and genomic signal processing of microarray for cancer gene expression was analysed. The significance of the microarray filtering base clustering and classification techniques were added the significance of effective clustering of microarray gene expression. Here we have given a comprehensive study of different existing techniques, algorithm, tools used in the field of clustering. The basic objective of this paper to motivate the researchers for developing potentially more effective new algorithms as till date many problems in clustering process of microarray gene expression data are still unsolved.

References

- [1] Jie Hwa Yang, M.J.Buckley,T.P.Spped, “Analysis of cDNA microarray images”, Henry Stewart Publications,Briefings in bioinformatics,Vol 2,No.4,Page 341-349,December (2001)
- [2] P.Jaluria,K.Konstantopaolos,M.Betenbaugh, “A perspective on microarrays:current applications,pitfalls,and potential uses”, Microbial cell factories,Vol 6, No.4, Page 1-14 (2007)
- [3] Jie Liang, Semen Kachalo, “Computatinal analysis of microarray gene expression profiles: clustering, classification and beyond”, Chemometrics and intelligent laboratory systems,Elsevier,Vol 62,Page 199-216 (2002)
- [4] Daxin Jiang, Chun Tang, Aidong Zhang,”Cluster analysis for gene expression data: a survey, knowledge and Data Engineering”,IEEE Transactions on , (2004);16(11),1370-1386.
- [5] M.B. Eisen, P.T. Spellman, P.O. Brown, and D. Botstein, “Cluster Analysis and Display of Genome-Wide Expression Patterns,” Proc. Nat’l Academy of Science, vol. 95, no. 25, pp. 14863-14868, Dec. (1998).
- [6] Nidheesh, N., Nazeer, K.A. and Ameer, P.M.. An enhanced deterministic K-Means clustering algorithm for cancer subtype prediction from gene expression data. Computers in biology and medicine, 91, pp.213-221. (2017)
- [7] Jothi, R., Mohanty, S.K. and Ojha, A.. DK-means: a deterministic k-means clustering algorithm for gene expression analysis. Pattern Analysis and Applications, 22(2), pp.649-667. (2019)
- [8] Gasch, A.P. and Eisen, M.B., Exploring the conditional coregulation of yeast gene expression through fuzzy k-means clustering. Genome biology, 3(11), pp.research0059-1. (2002)
- [9] Suresh, R.M., Dinakaran, K. and Valarmathie, P., April. Model based modified k-means clustering for microarray data. In 2009 International Conference on Information Management and Engineering (pp. 271-273). IEEE.(2009)
- [10] <https://www.researchgate.net/publication/10748487> : Fuzzy C-Means Method for Clustering Microarray Data, June (2003).
- [11] Dhiraj K., Rath S.K., Babu K.S. , FCM for Gene Expression Bioinformatics Data. In: Ranka S. et al. (eds) Contemporary Computing. IC3 2009. Communications in Computer and Information Science, vol 40. Springer, Berlin, Heidelberg. https://doi.org/10.1007/978-3-642-03547-0_50 (2009)
- [12] Xiao, X., Dow, E.R., Eberhart, R., Miled, Z.B. and Oppelt, R.J., April. Gene clustering using self-organizing maps and particle swarm optimization. In Proceedings International Parallel and Distributed Processing Symposium (pp. 10-pp). IEEE. (2003)
- [13] Sun, J., Chen, W., Fang, W., Wun, X. and Xu, W., Gene expression data analysis with the clustering method based on an improved quantum-behaved Particle Swarm Optimization. Engineering Applications of Artificial Intelligence, 25(2), pp.376-391. (2012)

- [14] Lam, Y.K., Tsang, P.W.M. and Leung, C.S., PSO-based K-Means clustering with enhanced cluster matching for gene expression data. *Neural Computing and Applications*, 22(7-8), pp.1349-1355.(2013)
- [15] Chen, W., Sun, J., Ding, Y., Fang, W. and Xu, W. Clustering of gene expression data with quantum-behaved particle swarm optimization. In *International Conference on Industrial, Engineering and Other Applications of Applied Intelligent Systems* (pp. 388-396). (2008), June Springer, Berlin, Heidelberg.
- [16] Dey, L. and Mukhopadhyay, A. Microarray gene expression data clustering using PSO based K-means algorithm. In *Proceedings of the International Conference on Advanced Computing, Communication and Networks* (Vol. 1, pp. 587-591).(2011)
- [17] Deng, Y., Kayarath, D., Elasm, M.O. and Brown, S.J. Microarray data clustering using particle swarm optimization K-means algorithm. *Proc. 8th JCIS*, pp.1730-1734.(2005)
- [18] Chakraborty, A. and Maki, H., November. Biclustering of gene expression data using genetic algorithm. In *2005 IEEE Symposium on Computational Intelligence in Bioinformatics and Computational Biology* (pp. 1-8). IEEE. (2005)
- [19] Pan, H., Zhu, J. and Han, D. Genetic algorithms applied to multi-class clustering for gene expression data. *Genomics, Proteomics & Bioinformatics*, 1(4), pp.279-287. (2003)
- [20] Balamurugan, R., Natarajan, A.M. and Premalatha, K. A new hybrid cuckoo search algorithm for biclustering of microarray gene-expression data. *Applied Artificial Intelligence*, 32(7-8), pp.644-659 (2018).
- [21] Sampathkumar, A., Rastogi, R., Arukonda, S., Shankar, A., Kautish, S. and Sivaram, M. An efficient hybrid methodology for detection of cancer-causing gene using CSC for micro array data. *Journal of Ambient Intelligence and Humanized Computing*, pp.1-9 (2020).
- [22] Datta, S. and Datta, S., Comparisons and validation of statistical clustering techniques for microarray gene expression data. *Bioinformatics*, 19(4), pp.459-466 (2003).
- [23] Bolshakova, N. and Azuaje, F. Machaon CVE: cluster validation for gene expression data. *Bioinformatics*, 19(18), pp.2494-2495 (2003).
- [24] Bolshakova, N., Zamolotskikh, A. and Cunningham, P., Comparison of the data-based and gene ontology-based approaches to cluster validation methods for gene microarrays. In *19th IEEE Symposium on Computer-Based Medical Systems (CBMS'06)* (pp. 539-543) June (2006). IEEE.
- [25] Jaskowiak, P.A., Campello, R.J. and Costa, I.G., Proximity measures for clustering gene expression microarray data: a validation methodology and a comparative analysis. *IEEE/ACM transactions on computational biology and bioinformatics*, 10(4), pp.845-857 (2013).
- [26] Ghalib, M.R., Nandeibam, B. and Ghosh, D.K., Microarray Gene Expression Data Analysis through a Hybrid Clustering Algorithm incorporated with Validation Techniques.
- [27] Costa, I.G., de Carvalho, F.D.A. and de Souto, M.C. Comparative analysis of clustering methods for gene expression time course data. *Genetics and Molecular Biology*, 27(4), pp.623-631 (2004).
- [28] Okabe, H., Satoh, S., Kato, T., Kitahara, O., Yanagawa, R., Yamaoka, Y., Tsunoda, T., Furukawa, Y. and Nakamura, Y. Genome-wide analysis of gene expression in human hepatocellular carcinomas using cDNA microarray: identification of genes involved in viral carcinogenesis and tumor progression. *Cancer research*, 61(5), pp.2129-2137 (2001).
- [29] Zembutsu, H., Ohnishi, Y., Tsunoda, T., Furukawa, Y., Katagiri, T., Ueyama, Y., Tamaoki, N., Nomura, T., Kitahara, O., Yanagawa, R. and Hirata, K., Genome-wide cDNA microarray screening to correlate gene expression profiles with sensitivity of 85 human cancer xenografts to anticancer drugs. *Cancer research*, 62(2), pp.518-527 (2002).
- [30] Subramanian, A., Tamayo, P., Mootha, V.K., Mukherjee, S., Ebert, B.L., Gillette, M.A., Paulovich, A., Pomeroy, S.L., Golub, T.R., Lander, E.S. and Mesirov, J.P., Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proceedings of the National Academy of Sciences*, 102(43), pp.15545-15550 (2005).

- [31] Khan, J., Saal, L.H., Bittner, M.L., Chen, Y., Trent, J.M. and Meltzer, P.S. Expression profiling in cancer using cDNA microarrays. *ELECTROPHORESIS: An International Journal*, 20(2), pp.223-229 (1999).
- [32] Khan, J., Simon, R., Bittner, M., Chen, Y., Leighton, S.B., Pohida, T., Smith, P.D., Jiang, Y., Gooden, G.C., Trent, J.M. and Meltzer, P.S., Gene expression profiling of alveolar rhabdomyosarcoma with cDNA microarrays. *Cancer Research*, 58(22), pp.5009-5013 (1998).
- [33] Gardina, P.J., Clark, T.A., Shimada, B., Staples, M.K., Yang, Q., Veitch, J., Schweitzer, A., Awad, T., Sugnet, C., Dee, S. and Davies, C., Alternative splicing and differential gene expression in colon cancer detected by a whole genome exon array. *BMC genomics*, 7(1), p.325 (2006).
- [34] Takata, R., Katagiri, T., Kanehira, M., Tsunoda, T., Shuin, T., Miki, T., Namiki, M., Kohri, K., Matsushita, Y., Fujioka, T. and Nakamura, Y., Predicting response to methotrexate, vinblastine, doxorubicin, and cisplatin neoadjuvant chemotherapy for bladder cancers through genome-wide gene expression profiling. *Clinical cancer research*, 11(7), pp.2625-2636 (2005).
- [35] Clarke, P.A., te Poele, R., Wooster, R. and Workman, P., Gene expression microarray analysis in cancer biology, pharmacology, and drug development: progress and potential. *Biochemical pharmacology*, 62(10), pp.1311-1336 (2001).
- [36] Saito-Hisaminato, A., Katagiri, T., Kakiuchi, S., Nakamura, T., Tsunoda, T. and Nakamura, Y., Genome-wide profiling of gene expression in 29 normal human tissues with a cDNA microarray. *DNA research*, 9(2), pp.35-45 (2002).
- [37] Cho, S.B. and Won, H.H., 2003, January. Machine learning in DNA microarray analysis for cancer classification. In *Proceedings of the First Asia-Pacific bioinformatics conference on Bioinformatics Volume 19* (pp. 189-198) (2003).
- [38] Segal, E., Sirlin, C.B., Ooi, C., Adler, A.S., Gollub, J., Chen, X., Chan, B.K., Matcuk, G.R., Barry, C.T., Chang, H.Y. and Kuo, M.D., Decoding global gene expression programs in liver cancer by noninvasive imaging. *Nature biotechnology*, 25(6), pp.675-680 (2007).
- [39] Iorio, M.V., Ferracin, M., Liu, C.G., Veronese, A., Spizzo, R., Sabbioni, S., Magri, E., Pedriali, M., Fabbri, M., Campiglio, M. and Ménard, S., MicroRNA gene expression deregulation in human breast cancer. *Cancer research*, 65(16), pp.7065-7070 (2005).
- [40] Watanabe, T., Komuro, Y., Kiyomatsu, T., Kanazawa, T., Kazama, Y., Tanaka, J., Tanaka, T., Yamamoto, Y., Shirane, M., Muto, T. and Nagawa, H., Prediction of sensitivity of rectal cancer cells in response to preoperative radiotherapy by DNA microarray analysis of gene expression profiles. *Cancer research*, 66(7), pp.3370-3374 (2006).
- [41] Sgroi, D.C., Teng, S., Robinson, G., LeVangie, R., Hudson, J.R. and Elkahoul, A.G., In vivo gene expression profile analysis of human breast cancer progression. *Cancer research*, 59(22), pp.5656-5661 (1999).
- [42] Asyali, M.H. and Alci, M., Reliability analysis of microarray data using fuzzy c-means and normal mixture modeling based classification methods. *Bioinformatics*, 21(5), pp.644-649 (2005).
- [43] Ouyang, M., Welsh, W.J. and Georgopoulos, P., Gaussian mixture clustering and imputation of microarray data. *Bioinformatics*, 20(6), pp.917-923 (2004).
- [44] Dembele, D. and Kastner, P., Fuzzy C-means method for clustering microarray data. *bioinformatics*, 19(8), pp.973-980 (2003).
- [45] Maulik, U. and Mukhopadhyay, A., Simulated annealing based automatic fuzzy clustering combined with ANN classification for analyzing microarray data. *Computers & operations research*, 37(8), pp.1369-1380 (2010).
- [46] stepanian, R.S., Sungoor, A. and Nebel, J.C., August. Linear predictive coding and wavelet decomposition for robust microarray data clustering. In *2007 29th Annual International Conference of the IEEE Engineering in Medicine and Biology Society* (pp. 4629-4632) (2007). IEEE.

- [47] Keerin, P., Kurutach, W. and Boongoen, T., October. Cluster-based KNN missing value imputation for DNA microarray data. In 2012 IEEE International Conference on Systems, Man, and Cybernetics (SMC) (pp. 445-450) **(2012)**. IEEE.
- [48] Liu, W., Yuan, K. and Ye, D., Reducing microarray data via nonnegative matrix factorization for visualization and clustering analysis. *Journal of biomedical informatics*, 41(4), pp.602-606 **(2008)**.
- [49] Qu, Y. and Xu, S., Supervised cluster analysis for microarray data based on multivariate Gaussian mixture. *Bioinformatics*, 20(12), pp.1905-1913 **(2004)**.
- [50] Tibshirani, R., Hastie, T., Eisen, M., Ross, D., Botstein, D. and Brown, P., Clustering methods for the analysis of DNA microarray data. *Dept. Statist., Stanford Univ., Stanford, CA, Tech. Rep* **(1999)**.
- [51] Maji, P., Fuzzy-rough supervised attribute clustering algorithm and classification of microarray data. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, 41(1), pp.222-233 **(2010)**.
- [52] Hengpraprom, S. and Chongstitvatana, P., October. Selecting Informative Genes from Microarray Data for Cancer Classification with Genetic Programming Classifier Using K-Means Clustering and SNR Ranking. In 2007 Frontiers in the Convergence of Bioscience and Information Technologies (pp. 211-218) **(2007)**. IEEE.
- [53] Díaz-Uriarte, R. and De Andres, S.A., Gene selection and classification of microarray data using random forest. *BMC bioinformatics*, 7(1), p.3 **(2006)**.
- [54] Mishra, P., Bhoi, N. Microarray Filtering-Based Fuzzy C-Means Clustering and Classification in Genomic Signal Processing. *Arab J Sci Eng* 44, 9381–9395 **(2019)**. <https://doi.org/10.1007/s13369-019-03945-0>
- [55] Mishra, P., Bhoi, N. Genomic signal processing of microarrays for cancer gene expression and identification using cluster-fuzzy adaptive networking. *Soft Comput* **(2020)**. <https://doi.org/10.1007/s00500-020-05068-3>
- [56] Wang, Y. and Makedon, F., August. Application of Relief-F feature filtering algorithm to selecting informative genes for cancer classification using microarray data. In *Proceedings. 2004 IEEE Computational Systems Bioinformatics Conference, 2004. CSB 2004.* (pp. 497-498) **(2004)**. IEEE.
- [57] Bonilla-Huerta, E., Hernandez-Montiel, A., Morales-Caporal, R. and Arjona-Lopez, M., Hybrid framework using multiple-filters and an embedded approach for an efficient selection and classification of microarray data. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 13(1), pp.12-26 **(2015)**.
- [58] Yang, P., Zhou, B.B., Zhang, Z. and Zomaya, A.Y., A multi-filter enhanced genetic ensemble system for gene selection and sample classification of microarray data. *BMC bioinformatics*, 11(S1), p.S5 **(2010)**.
- [59] Bolón-Canedo, V., Sánchez-Marroño, N. and Alonso-Betanzos, A., An ensemble of filters and classifiers for microarray data classification. *Pattern Recognition*, 45(1), pp.531-539 **(2012)**.
- [60] Leung, Y. and Hung, Y., A multiple-filter-multiple-wrapper approach to gene selection and microarray data classification. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 7(1), pp.108-117 **(2008)**.
- [61] Tritchler, D., Parkhomenko, E. and Beyene, J., Filtering genes for cluster and network analysis. *BMC bioinformatics*, 10(1), p.193 **(2009)**.
- [62] Miller, J.A., Menon, V., Goldy, J., Kaykas, A., Lee, C.K., Smith, K.A., Shen, E.H., Phillips, J.W., Lein, E.S. and Hawrylycz, M.J., Improving reliability and absolute quantification of human brain microarray data by filtering and scaling probes using RNA-Seq. *BMC genomics*, 15(1), pp.1-14 **(2014)**.
- [63] Hwang, T., Sun, C.H., Yun, T. and Yi, G.S., FiGS: a filter-based gene selection workbench for microarray data. *BMC bioinformatics*, 11(1), p.50 **(2010)**.

- [64] Ke, W., Wu, C., Wu, Y. and Xiong, N.N., A new filter feature selection based on criteria fusion for gene microarray data. *IEEE Access*, 6, pp.61065-61076 **(2018)**.
- [65] Biesiada, J. and Duch, W., November. A Kolmogorov-Smirnov Correlation-Based Filter for Microarray Data. In *International Conference on Neural Information Processing* (pp. 285-294) **(2007)**. Springer, Berlin, Heidelberg.
- [66] Fleury, G., Hero, A., Yoshida, S., Carter, T., Barlow, C. and Swaroop, A., September. Pareto analysis for gene filtering in microarray experiments. In *2002 11th European Signal Processing Conference* (pp. 1-4) **(2002)**. IEEE.