# ASSOCIATION RULE MINING FROM GENE EXPRESSION USING RANWAR ALGORITHM

**SATHIYAPRIYA J**
*Ph.D Research Scholar*
*Department of Computer Science*
*Bharathiar University*
*Coimbatore – 641046*
*Sathyapriya08@yahoo.com*

**D. RAMYACHITRA**
Assistant Professor
Department of Computer Science
Bharathiar University
Coimbatore – 641046
jaichitra1@yahoo.co.in

*Abstract*— **Ranking of association rules is as of now an intriguing subject in data mining and bioinformatics. The gigantic number of advanced tenets of things (or, qualities) by Association Rule Mining (ARM) calculations settles on perplexity to the leader. In this project, we propose a weighted rule-mining strategy (say, RANWAR or Rank-based Weighted Association Rule-mining) to rank the tenets utilizing two novel lead intriguing quality measures, Rank-based Weighted Condensed Support (WCS) and Weighted Condensed Confidence (WCC) measures to sidestep the issue. These measures are essentially relied on upon the rank of things (qualities). Utilizing the rank, we dole out weight to everything. RANWAR produces a great deal less number of incessant thing sets than the cutting edge affiliation run mining calculations. Along these lines, it spares the time of execution of the calculation. We run RANWAR on quality expression and methylation datasets. The qualities of the top guidelines are organically approved by Gene Ontologies (GOs) and KEGG pathway examinations. Many tops positioned decides removed from RANWAR that hold poor positions in conventional Apriority, are profoundly naturally noteworthy to the related maladies. Finally, the top RANWAR requirements that aren't in Apriority are taken into account.**

**Keywords** — *Rank, Weight, Genes, WCs, Wcc.*

## I. INTRODUCTION

Data Disclosure and Data Mining (KDD) is an interdisciplinary space that generally focuses on the productive techniques for obtaining captivating standards and cases from the data. The basic ordinary illustrations which are surveyed by charming quality measures consolidate alliance rules. Association control mining (ARM), a basic data digging system is utilized for perceiving entrancing associations between things. The large number of regulations makes it difficult to choose the finest among them. As a result, establishing standards from natural data is an important area for research. For this, assorted manage intriguing quality measures (viz., support, assurance, lift, conviction et cetera.) were proposed. In any case, these still make immense number of visit thing sets, and along these lines these make monstrous number of connection standards. Along these lines, bundle of time is taken to run these computations. In this article, we propose a weighted represent mining technique (viz., RANWAR or Rank-based Weighted Alliance Oversee Mining) which has been created using two novel measures rank-based weighted thick support (say, wcs)

and rank-based weighted combined sureness (say, wcc) measures for isolating standards from the data. Sooner or later, it happens that a lot of standards have same support and same conviction. Starting at now, if we require a couple of them, it is difficult to isolate among them.

In this way, on the off chance that we apply the WCs and wcc, we can without a lot of an extend organize them. The noteworthy favorable position of RANWAR is that it creates significantly a smaller number of standard things sets than front line association oversees digging figuring's for same slightest reinforce regard. There is no such ARM methodology which makes lesser number of persistent things sets than RANWAR. In this way, it takes fundamentally less time than alternate figuring's. Another favorable position of RANWAR is that a segment of the standards which hold low rank in traditional represent mining estimations, hold awesome rank in RANWAR. A couple of affirmations of regular criticalness of the characteristics of the propelled measures are furthermore found. As we understand that if the number of characteristics in data is broad, the quantity of thing sets will be also significant, thus, using Limma quantifiable test, we have as of late considered the top differentially imparted (i.e., DE) or differentially methylated (i.e., DM) qualities. Limma is a useful accurate test which performs well for both ordinarily and non-routinely appropriated data for a wide range of test size (i.e., little, medium, tremendous). Our proposed measures are basically rank-based weighted measures. In this manner, situating of characteristics has a significant part here. Limma test gives a rank-wise quality rundown according to their p-values from best to most negative situations. Starting there, we dole out weight to everything/quality as for their p-regard situating, and fuse these into the measures. Subsequently, our measures offer essentialness to everything (quality). Our proposed measures (viz., WCs and wcc) are combined sort of the traditional support and sureness measures. In addition, two quality expression datasets and two methylation datasets are used to test the execution of RANWAR. We have made a relative examination of it with the regular Apriori computation and other bleeding edge oversees mining counts. For endorsement of the tenets, GO terms and KEGG pathways of the qualities in the gauges are perceived.

The characteristics of the propelled benchmarks including most vital number of GOs/pathways are represented naturally kind hearted points. Finally, we report many top situated standards conveyed by RANWAR that

hold poor positions in standard Apriori, however are outstandingly actually imperative to related afflictions.

## II. LITERATURE SURVEY

ARM is an acclaimed framework to gage interesting associations among different things (i.e., qualities). Expect, Itemset = {i1, i2, ..., in} be an itemset (i.e., set of characteristics) and S = {s1, s2, ..., sm} be a game plan of trades (tests) [1] . Thusly, a run the show might be depicted as A) C, where A, C _ Itemset and ATC = _. Here, an is called as herald and C is called as following. In a trade database, a trade may include a game plan of things purchased in it. In a relative sense, in quality expression/methylation [3] dataset, in any tissue test (trans-action), a course of action of characteristics may happen together. Some of them are up managed/hyper-methylated, and some are down-coordinated/hypomethylated, besides, remaining are not differentially conveyed/methylated [4] (i.e., neither up-coordinated/hyper-methylated nor down directed/hypo-methylated). On the off chance that there ought to emerge an event of a characteristic trade, accept, {gene1+, gene2−, gene3nde => gene4+} is an association choose which communicates that if gene1 is up-coordinated (meant by '+'), gene2 is down-coordinated (meant by '- '), and gene3nde is non-differentially conveyed (outlined as 'nde') at the same time, by then it is likely that gene4 will be up-overseen. In this way, those four qualities must occur in some of trades at the same time. The support of an itemset is described as number of trades in which all things of the itemset appear in the meantime. The itemset is visit when its support is more critical than any edge regard (i.e., minimum support). The assurance of the run is described as extent of support of the itemset to the support of antecedent. Apriori is a fundamental estimation for learning alliance guidelines to control on databases that have trades [5]. Apriori utilizes a "base up" approach, where visit subsets are opened up one thing at a chance to make each cheerful and social affairs of the hopefuls are attempted against the data. The count closes if there are no further productive extensions to be perceived. The yield of Apriori is truly the courses of action of rules that make the occasion of things in the dataset. Apriori takes after extensiveness first interest to number the confident item sets. Apriori produces candidate item sets of length k from item sets of length k − 1. Starting there, it wipes out the hopefuls having an uncommon sub-outline. By further examinations, unmistakable limitations have been found in the standard Apriori count, like period of gigantic number of progressive item sets, high snuck past time, different compass issue, stack abnormality issue, acquiring same criticalness to everything et cetera. For diminishing those deficiencies, assorted enhancements have been performed on the main Apriori computation (viz., Orlando et al. in 2001, Pavon et al. in 2006, Yu et al. in 2008 , Oguz et al. in 2012 et cetera.). Other than that, various other ARM methodologies have been proposed (e.g., Tao et al. in 2003). Tao et al. used weighted ARM framework in capable way. In any case, less change on reducing snuck past time has been done through Tao et al. Thusly, progress new procedures are proposed (Yun et al. (WIP) in 2006, Hong et al. in 2008, Sun et al. in 2008, Ahmed et al. in 2008) [6]. Regardless, it has been seen that it is especially difficult to decrease each and every such limitation at the same time. Consequently, we have basically focused how to diminish snuck past time for oversee mining in such way that elite top situated things and their related significantly enormous rules will present in result for considerable trade database.

The present Learning Exposure and Data Mining (KDD) is an interdisciplinary territory that generally focuses on the precise strategies for getting interesting fundamentals and cases from the data [7]. The basic typical cases which are evaluated by charming quality measures fuse connection rules. Association Administer Mining (ARM) is an imperative data burrowing strategy is utilized for recognizing entrancing associations between things. Huge amounts of standards constantly make issue to pick best among them. Thusly, the situating of gauges from the normal data is basic range for research. For this, particular oversee charming quality measures (viz., reinforce sureness, lift, conviction et cetera.)

The Huge number of measures reliably makes issue to pick beat among them. Thusly, the situating of principles from the regular data is basic district for research [7].

The tremendous number of created rules of things (or, qualities) by alliance control mining (ARM) estimations settles on perplexity to the pioneer. Quality of Organization not guaranteed.

It's very difficult to reduce each and every such limitation in the meantime.

## III. PROPOSED METHODOLOGY

The proposed structure is a weighted lead mining framework (viz., RANWAR or Rank-based Weighted Connection Control Mining) which has been created using two novel measures rank-based Weighted Thick Support (say, WCS) and rank-based Weighted Solidified Conviction (say, WCC) measures for expelling rules from the data. Sooner or later, it happens that a lot of standards have same support and same conviction. Starting at now, if we require some of them, it is difficult to isolate among them. Subsequently, if we apply the WCS and WCC, we can without quite a bit of an extend arrange them. The critical favorable position of RANWAR is that it delivers significantly a smaller number of consistent things sets than best in class connection control burrowing figurings for same slightest reinforce regard. There is no such ARM system which makes lesser number of progressive things sets than RANWAR. In this way, it takes extensively less time than exchange figuring. Another preferred standpoint of RANWAR is that a segment of the fundamentals which hold low rank in traditional oversee mining counts hold incredible rank in RANWAR. A couple of affirmations of natural giganticness of the characteristics of the propelled rules are moreover found. The support of a thing set is portrayed as number of trades in which all things of the thing set appear at the same time. There is no such

ARM procedure which produces lesser number of normal things sets than RANWAR and its execution time is less. Guarantee Nature of Organization essential showed by customers
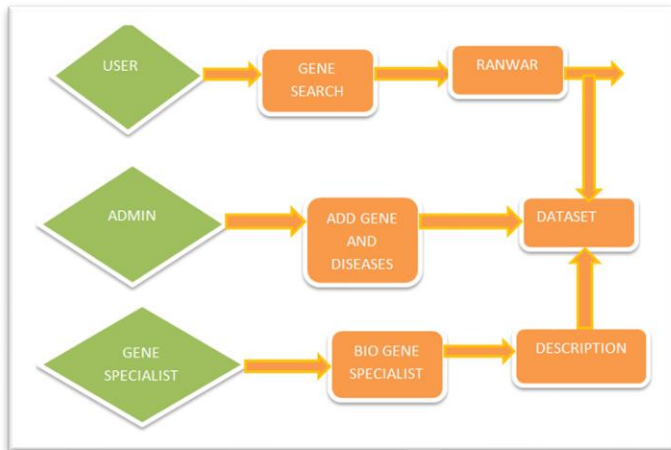at same time.

## IV.  ARCHITECTURE DIAGRAM



Fig. 1.          System Architecture

The Design graph is commonly utilized for a larger amount, less point-by-point depiction pointed more at understanding the general ideas and less at understanding the subtle elements of usage. The design chart obviously says how the client going to the site and looking the item to their loved. The fundamental system behind this engineering work is to explore the client conduct in light of their enjoyed. I have made a site and in that an arrangement of item depictions has been portrayed, through this site just the purchaser will explore the client conduct in which way the explore the item in the item site, while the client track the item in the item site the route of item following will be put away in the database and the fundamental page site will give the related points of interest of the client explored item while they login to the site to prescribe the item.

## V.          ALGORITHM

### A.  ASSOCIATION RULES

An affiliation decides is an example that states when X happens, y happens with certain likelihood. The primary objective is to discover all decides that fulfill the client determined least support (minsup) and least certainty. In this affiliation show distinctive set calculation as been utilized as a part of that Apriori is a standout amongst the most every now and again utilized strategy to examination the issue. There are two procedures to out all the vast thing sets from the database in Apriori calculation.

To ascertion the affiliation administer the arrangement of guidelines been taken after,

(i)   Support: The lead holds with bolster sup in T (the exchange informational index) if sup%  of exchange contains

X£Y.

Sup=Pr(X£Y).

(ii)   Confidence: The administer holds in T with certain conf if conf% of exchange that contains X likewise contain Y.

Conf= Pr( Y| X)

$$\text{Support (A)} = \frac{\text{Number of transaction in which A appears}}{\text{Total number of transactions}}$$

$$\text{Confidence (A} \rightarrow \text{B)} = \frac{\text{Support(AUB)}}{\text{Support(A)}}$$

### B.  APRIORI ALGORITHM

Apriori utilizes a "bottom up" approach, where visit subsets are argumented one thing a once (a stage known as competitor era), and gatherings of applicants are tried against the information. The calculation ends when no further fruitful expansions are found. Apriori utilizes broadness first hunt and a Hash tree structure to number competitor things sets effectively. It produces hopeful thing set of length K from thing set of length (k-1), and afterward it prunes the applicant set contains all incessant K-length thing sets. Fro that point onward, it check the exchang databas to decide visit thing sets among the competitors.

The pseudo code for the calculation is given beneath for an exchange database T, and a bloster limit of epsilon. Common set theoretic documentation is utilized; however take note of that T is a multi set. CK is the applicant set for level K. Produce () Calculation is accepted to create the hopeful sets from the huge thing sets of the first level, paying attention to the descending conclusion lemma. Count[c]. Access a field of the information structure that speaks to hopeful set c, which is at first thought to be zero. Many subtle elements are discarded underneath, generally the most critical piece of the usage is the information discarded underneath, generally the most critical piece of the usage is the information structure utilized for putting away the hopeful sets, and checking their frequencies.

Apriori standard: If a thing set is visit, then the majority of its subsets should likewise be vist. Illustration: Assume { C,D,E} is visit thing set. At that point any exchange that contains {c,D,E} should likewise contain {C,D}, {C,E}, {D,E}, {c}, {D} and {E}.

The general purpose of the calculation (and information mining, as a rule) is to separate valuable data from a lot of information. For instance, the data that a client who buys a console additionally tends to purchase a mouse in the meantime is gained from the pseudocode of Apriori :

*Lk: frequent k-itemset, satisfy minimum support Ck: candidate k-itemset, possible frequent k-itemsets*

*Input : Database*

*Output : Large Itemsets*

```
L₁={frequent 1-itemsets};
for (k=2; Lₖ₋₁ ≠ 0; k++) do begin
    Cₖ=apriori-gen(Lₖ₋₁);
    for each transactions t ∈D do begin //scan DB
    Cₜ=subset(Cₖ, t) //get the subsets of t that are candidates
    for each candidate c ∈Cₜ do
        c.count++;
    end
    Lₖ={c ∈Cₖ | c.count ≥ minsup}
end
Answer=∪ₖLₖ;
```

affiliation lead underneath: Support: The rate of errand significant information exchanges for which the example is valid.

### C. RANWAR ALGORITHM

*Input: Data matrix D (rows = genes, columns = samples), original gene-list A1 according to D, rank-wise gene-list A2 (according to p-values of genes by Limma), flag of sorting the evolved rules sort Flag (w.r.t. either WCs or Wcc), minimum support threshold min wsupp, minimum confidence threshold min wconf.*

*Output: Set of rules Rules, support RuleSupp, confidence RuleConf.*

$$support(I) = \frac{Number\ of\ transactions\ containing\ I}{Total\ number\ of\ transactions}$$

Confidence: The measure of certainty or trustworthiness with each discovered pattern.

$$confidence(X \rightarrow Y) = \frac{Number\ of\ transactions\ containing\ X\ and\ Y}{Number\ of\ transactions\ containing\ X}$$

**procedure RANWAR:**

1. *Normalize the data-matrix D using zero-mean normalization.*

2. *Calculate rank of genes (i.e., rankk (:)) according to original gene list A1.*

3. *Assign weights wt(:) to all genes according to their ranks rank(:).*

4. *Transpose the normalized data-matrix.*

5. *Choose initial seed values for using k-means clustering.*

6. *Discretize the transposed matrix applying standard k-means clustering sample-wise.*

7. *Apply post-discretization technique.*

8. *Initialize k = 1.*

9. *Find frequent 1-itemsets, FIk = {i|i ∈A1 ∧ WCs (i) ≥min wsupp}.*

10. *repeat*

11. *k=k+1.*

12. *Generate candidate item sets, CIk from FIk−1 item sets.*

13. *for each candidate itemset, c ∈ CIk do*

14. *Calculate wcs(c) for each candidate itemset, c.*

15. *if wcs(c) >= min wsupp then*

16. *FIk ← [FIk; c].*

17. *Generate rules, rule(:) from the frequent itemset, c.*

18. *Determine wcc(:) for each rule(:).*

19. *for each evolved rule, r ∈ rule(:) do*

20. *if wcc(r) >= min wconf then*

21. *Store the r in the resulting rule-list Rules with its wcs and wcc;*

    *Rules ← r, RuleSupp ← wcs(r) and RuleConf ← wcc(r).*

22. *end if*

23. *end for*

24. *end if*

25. *end for*

26. *until (FIk = Ø)*

27. *end procedure*

### VI.          RESULTS

In this project, we propose a weighted rule-mining technique (say, RANWAR or Rank-based Weighted Association Rule-mining) to rank the rules using two novel rule-interestingness measures, Rank-based Weighted Condensed Support (WCS) and Weighted Condensed Confidence (WCC) measures to bypass the pro blem. These measures are basically depended on the rank of items (genes). Using the rank, we assign weight to each item. RANWAR generates much smaller number of frequent item sets than the state-of-the-art association rule mining algorithms. Thus, it saves time of execution of the algorithm.

Table 1.    Gene Rank

| Gene | Rank_in_Gene |
|------|--------------|
| CRH | 1 |
| CXCR4 | 3 |
| BRCA1 | 3 |
| CCL14 | 5 |
| ALB | 6 |
| EPB42 | 7 |
| SH3TC1 | 7 |
| SH3BP5 | 9 |
| MLN | 10 |
| BRCA2 | 10 |

Table 2.   Rank Based   Weighted Association    Rule Mining Gene Expression and Methylation Data

| Gene | Rule | WCC (%) | WCS (%) | Conf (%) | Support (%) |
|------|------|---------|---------|----------|-------------|
| ap123 | Apxc123 | 123 | 345 | 1.0 | 0.78 |
| apx123 | apxc | 123 | 123 | 5 | 6 |
| ssay | Ssaaa | 5 | 2 | 2 | 2 |

## VII.        DISCUSSION

In my research paper I have discussed about whenever we run RANWAR on gene expression and methylation datasets. The genes of the top rules are biologically validated by Gene Ontologies (GOs) and KEGG pathway analyses. Many top ranked rules extracted from RANWAR that hold poor ranks in traditional Apriority, are highly biologically significant to the related diseases. Finally, the top rules evolved from RANWAR that are not in Apriority are reported.
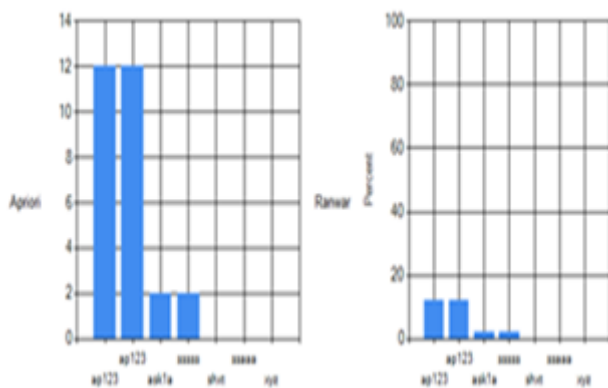


Fig. 2.        Comparison results

## VIII.CONCLUSION

The enormous number of principles advanced by ARM calculations settles upon the perplexity of the leader to choose the superior qualities. Along these lines, in this article, we have proposed two novel rank-based weighted consolidated run intriguing quality measures. A weighted lead mining calculation has been created utilizing the measures particularly for microarray/dot chip information. RANWAR uses an authentic test, Limma to process p-estimation of each quality (thing), and some weight is given to every quality in view of their p-esteem positioning. Is essentially weighted refreshed type of Apriori We utilize two quality expression datasets and two methylation datasets to contrast the execution of RANWAR and the best-in-class ARM calculations? Creates a smaller number of incessant things sets than the others. Hence, it spares time of execution of the RANWAR calculation. Another favorable position of RANWAR is that some most natural critical principles stand beat here which hold low rank in Apriori. The guidelines are approved by GO-terms and KEGG pathways of qualities of the standards. Some top standards separated from that are absent in Apriori, but rather have high natural criticalness, are additionally detailed.

## References

[1]  R. Agrawal, T. Imielinski and A. Swami, Mining Association Rules between Sets of Items in large Databases, Proc. ACM SIGMOD, New York, NY, USA:ACM, pp. 207-216, 2005.

[2]  G. Smyth, Linear Models and Empirical Bayes Methods for Assessing Differential Expression in Microarray Experiments, Statistical Applications in Genetics and Molecular Biology, vol. 3, no. 1, pp. 3, 2004.

[3]  S. Mallik and et al., Integrated Analysis of Gene Expression and Genome wide DNA Methylation for Tumor Prediction: An Association Rule Mining based Approach, CIBCB, IEEE SSCI, Singapore, pp. 120-127, 2013.

[4]  S. Bandyopadhyay, U. Maulik and J.T.L. Wang, Analysis of Biological Data: A Soft Computing Approach, World Scientific, Singapore, 2007.

[5]  M. Anandhavalli, M.K.Ghose and K.Gauthaman, Association Rule Mining in Genomics, International Journal of Computer Theory and Engineering, vol.2, no. 2, pp. 1793-8201,2010.

[6]  D. Arthur and S. Vassilvitskii, k-means++: the advantages of careful seeding, In Proc. of ACM-SIAM SODA 2007, Society for Industrial and Applied Mathematics Philadelphia, PA, USA, pp. 1027-1035, 2007.

[7]  S. Bandyopadhyay and et al., A Survey and Comparative Study of Statistical Tests for Identifying Differential Expression from Microarray Data, IEEE/ACM TCBB, 2014, DOI: 10.1109/TCBB.2013.147.

[8]  [8] A. Thomas and et al., Expression profiling of cervical cancers in Indian women at different stages to identify gene signatures during progression of the disease, Cancer Medicine, 2013, DOI: 10.1002/cam4.152.

[9]  J. Liu and et al., Identifying differentially expressed genes and pathways in two types of non-small cell lung cancer: adenocarcinoma and squamous cell carcinoma, Genet Mol Res, 2014, DOI: http://dx.doi.org/10.4238/2014.January.8.8.

[10]  W. Wei and et al., The potassium-chloride co transporter 2 promotes cervical cancer cell migration and invasion by an ion transport-independent mechanism, J Physiol., 2011, DOI: 10.1113/jphysiol.2011.214635.

J Physiol., 2011, DOI: 10.1113/jphysiol.2011.214635.