

Data Mining for Evolving Fuzzy Association Rules for Predicting Monsoon Rainfall of India

C.T. Dhanya and D. Nagesh Kumar

Department of Civil Engineering, Indian Institute of Science, Bangalore-560012, India

ABSTRACT

We used a data mining algorithm to evolve fuzzy association rules between the atmospheric indices and the Summer Monsoon Rainfall of All-India and two homogenous regions (Peninsular and West central). El Nino and Southern Oscillation (ENSO) and Equatorial Indian Ocean Oscillation zonal wind index (EQWIN) indices are used as the causative variables. Rules extracted are showing a negative relation with ENSO index and a positive relation with the EQWIN index. A fuzzy rule based prediction technique is also implemented on the same indices to predict the summer monsoon rainfall of All-India, Peninsular, and West central regions. Rules are defined using a training dataset for the period 1958-1999 and validated for the period 2000-2006. The fuzzy outputs of the defined rules are converted into crisp outputs using the weighted counting algorithm. The variability of the summer monsoon rainfall over the years is well captured by this technique, thus proving to be efficient even when the linear statistical relation between the indices is weak.

KEYWORDS

frequent pattern algorithm, weighted counting algorithm, ENSO, EQWIN, Indian summer monsoon rainfall

Reprint requests to: D. Nagesh Kumar, Professor, Department of Civil Engineering,
Indian Institute of Science, Bangalore-560012, India

1. INTRODUCTION

The influence of various climatic and atmospheric indices on Indian summer monsoon rainfall (ISMR) from June to September has been studied intensively by various studies for the last few decades. Such studies are of great significance because a major portion of the annual Indian rainfall is contributed by the summer monsoon rainfall. The predictor variables used for the study of ISMR include Darwin sea-level pressure, Latitudinal position of 500 mb ridge along 75° E, Arabian sea SST, Indian Ocean SST, Quasi-Biennial Oscillation, Sea-surface temperature anomalies over different Nino regions, Western Pacific region SST, Eastern Indian Ocean region SST, Eurasian surface temperature and Indian Surface temperature, Equatorial East Indian Ocean sea surface temperature, Nino 3.4, Equatorial Indian Ocean Oscillation zonal wind index (EQWIN), Eurasian snow cover and NW Europe temperature. The complex relation between these large scale circulation patterns and ISMR leads to a poor performance of the models used so far to forecast ISMR (Gadgil et al., 2004; Rajeevan et al., 2004).

Among these indices, El Nino and Southern Oscillation (ENSO) and Equatorial Indian Ocean Oscillation (EQUINOO) together can explain much of the ISMR variability (Gadgil, 2003). Gadgil et al. (2004) found a composite index that is a linear combination of the EQWIN and the ENSO index (Nino 3.4 SST), correlates better with an excess or deficit in ISMR than either index alone. The authors separated the excess and deficit ISMR events by a line determined by a linear combination of the EQWIN and ENSO index.

However, because the correlation between ISMR and ENSO index and also EQUINOO index is very small (0.33 and 0.19 respectively), traditional linear statistical modeling approaches may not prove useful, even at seasonal timescale (Gadgil et al., 2005). In the present study, we use a data mining algorithm viz., frequent pattern growth algorithm to evolve fuzzy association rules between ENSO, EQUINOO indices, and ISMR. Rules are also extracted for the summer monsoon rainfall of two homogenous regions of India—Peninsular and West central. These regions cover about 50% of the total area of India as can be seen from figure 1. Also, a fuzzy rule-based technique is used to quantify the relation of ENSO and EQUINOO index with summer monsoon rainfall of three regions and thereby reproduce the variability of rainfall in June to September (JJAS) months.

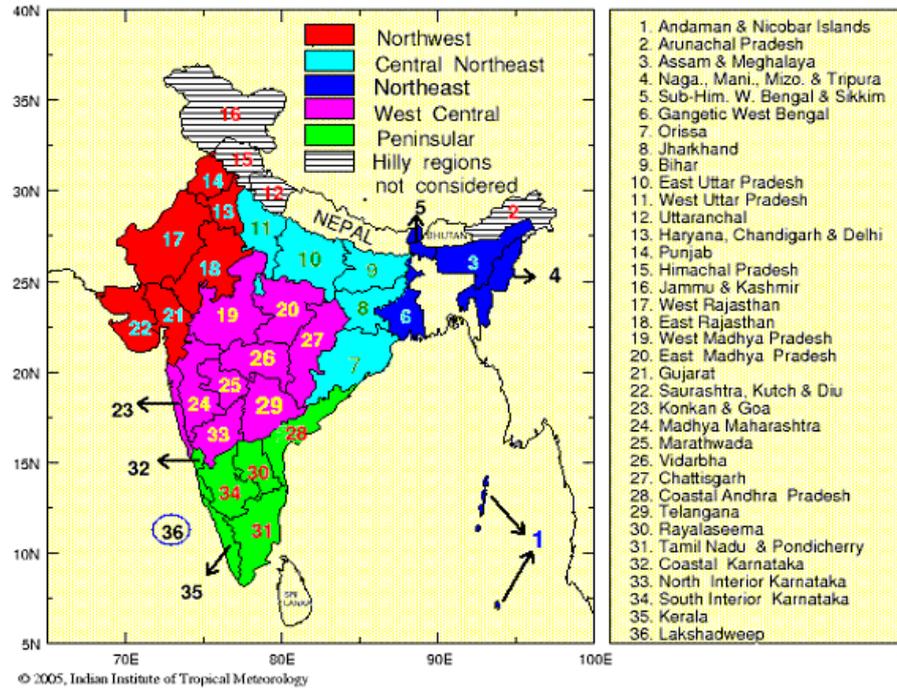


Fig. 1: Homogenous Monsoon Regions of India as defined by Indian Institute of Tropical Meteorology

2. DATA USED

1. ENSO index: Sea surface temperature anomaly from Nino 3.4 region (5°S – 5°N, 170°E – 120°W). Monthly sea surface temperature data from Nino 3.4 region for the period, January 1958 to December 2006, are obtained from the Website of National Weather Service, Climate Prediction Centre of NOAA (<http://www.cpc.noaa.gov/data/indices/>).
2. EQWIN index: Negative of zonal wind anomaly over equatorial Indian Ocean region (60°E – 90°E, 2.5°S – 2.5°N). Monthly surface wind data for the period, January 1958 to December 2006, are obtained from National Center for Environmental Prediction (<http://www.cdc.noaa.gov/Datasets>).
3. Rainfall data: Monthly rainfall data over entire India and over the two homogenous regions are obtained for the period, January 1958 to December

2006, from the Website of Indian Institute of Tropical Meteorology (IITM), Pune, India (<http://www.tropmet.res.in/data.html>).

According to Maity et al. (2007), ENSO and EQWIN index at a lag of one month can better separate the extreme positive and negative anomalies of monsoon rainfall. Most of the extreme positive anomalies occurred when the ENSO index < EQWIN and negative anomalies occurred when the ENSO index > EQWIN. Hence for the present study, standardized values of the ENSO and EQWIN index for the months May to August are considered as the antecedents. Standardized values of the rainfall anomalies for the monsoon months (June to September) are selected as the consequents.

3. FUZZY ASSOCIATION RULES

Association rules indicate whether or how much the values of an attribute depend on the values of the other attributes in the data set. A rule consists of a left-hand side proposition (antecedent) and a right hand side proposition (consequent). The rule states that when the antecedent occurs (is true), then the consequent also occurs (is true). The conditional probability of the occurrence of the consequent given the antecedent is referred to as the confidence of the rule. For example, if a pattern “B follows A” occurs n_1 times and the pattern “C follows B follows A” occurs n_2 times, then the association rule “whenever B follows A, C will also follow” has a confidence of (n_2/n_1) . The interestingness of a rule is usually measured in terms of its confidence.

In association rule mining, subdivision of the quantitative values into crisp sets would lead to over or underestimating values near the borders. Some of the rules may have been omitted because the combination of the antecedents may not be falling in the same discrete range, but they may be tending toward the nearby discrete states. Fuzzy sets can overcome this problem by allowing partial memberships to each of the different sets. Caulfield (1997) and Chen and Chen (2007) have shown the efficacy of fuzzy partition over crisp partition. Hence, a better extraction of the rules may be possible if the classification of the indices is done in a fuzzy manner and not in a crisp manner. The approach used to discover fuzzy association rules in the present study is described below:

3.1 Fuzzy Set Construction

All the three data sets are divided into five classes: (1) Lowest, (2) Low, (3) Medium, (4) High, and (5) Highest. The values and overlaps used for the classification are shown in figure 2.

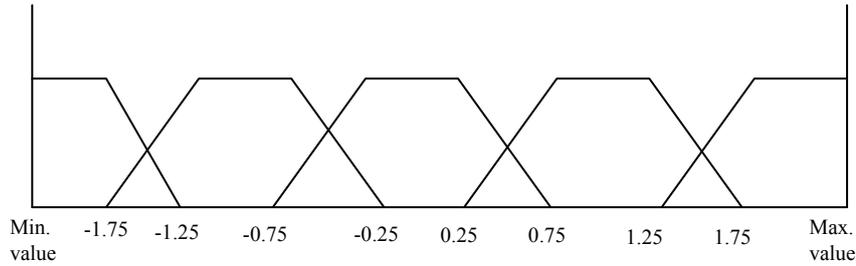


Fig. 2: Boundary values and overlap used for fuzzy set classification

The membership functions are computed based on whether it is near the upper border of a set or near the lower border.

Near the upper border,

$$\mu(x) = \frac{ub(f^n) - x}{ub(f^n) - lb(f^{n+1})} \tag{1}$$

Near the lower border,

$$\mu(x) = \frac{x - lb(f^n)}{ub(f^n) - lb(f^{n+1})} \tag{2}$$

where $lb(f^n)$ is the lower border of a set, $ub(f^n)$ is the upper border and x is the original value of any attribute in the database.

The abbreviations used for all various ranges of indices falling in different fuzzy classes are given in table 1.

TABLE 1

Abbreviations used for the classification

Fuzzy set	ENSO	EQWIN	Rainfall
1. Min value to -1.25	Nino1	Eqwin1	Extreme drought
2. -1.75 to -0.25	Nino2	Eqwin2	Moderate drought
3. -0.75 to 0.75	Nino3	Eqwin3	Normal
4. 0.25 to 1.75	Nino4	Eqwin4	Moderate flood
5. 1.25 to Max. value	Nino5	Eqwin5	Extreme flood

3.2 Preparation of Dataset for Mining

A new data set has to be constructed from the original database based on the definition of fuzzy sets described above. For every fuzzy set defined, there is one row in the new database containing the grade of membership of the single items to the specific set. In the present study, since five fuzzy sets are defined for each of the three attributes, there will be a total of 15 columns, the first 5 columns giving the membership values of first attribute for all five fuzzy sets, the next five of the second attribute and the last five of the third attribute. Thus, the new table will only contain the membership values of these fuzzy sets.

3.3 Support and Confidence of Fuzzy Rules

Triangular norms are used for the calculation of quality measures, support and confidence. A triangular norm is a commutative, associative, non-decreasing function in $T : [0,1]^2 \rightarrow [0,1]$ such that $T(x,1) = x$ for all $x \in [0,1]$. The basic continuous t-norms are the minimum, the product and the Lukasiewicz t-norms.

For an association rule $A \rightarrow B$, the support and confidence are given as:

$$\text{sup}(A \rightarrow B) = \sum_{x,y \in D} T(A(x), B(y)) \quad (3)$$

$$\text{conf}(A \rightarrow B) = \frac{\sum_{x,y \in D} T(A(x), B(y))}{\sum_{x \in D} A(x)} \quad (4)$$

For the present study Lukasiewicz t-norm is used for the calculation of support and confidence, since it is the most popular method for calculating fuzzy operations. Hence, support and confidence can be expressed as,

$$\text{sup}(A \rightarrow B) = \sum_{x,y \in D} \min(A(x), B(y)) \tag{5}$$

$$\text{conf}(A \rightarrow B) = \frac{\sum_{x,y \in D} \min(A(x), B(y))}{\sum_{x \in D} A(x)} \tag{6}$$

3.4 Frequent Pattern (FP) Growth Algorithm

The FP-Growth algorithm allows generating frequent itemsets and unlike the Apriori algorithm, it does not create huge number of candidates. The data are organized in a tree form, called the Frequent Pattern Tree. The algorithm first constructs the tree out of the original data set and then grows the frequent patterns.

The data are preprocessed before applying the algorithm. The dataset is scanned first to compute the support of each individual item. The items that have a support value less than a user specified minimum support are discarded. The remaining items are rearranged in a decreasing order with respect to their support. An example of the dataset preprocessing (Han et al., 1999) is given in table 2.

TABLE 2

Sample database showing re-arranging of the order

Original database	Support	Preprocessed database
abd	sup(b) = 6	bda
bcde	sup(d) = 5	bde
bd	sup(e) = 5	bd
ade	sup(a) = 4	dea
ab	sup(c) = 2	ba
abe	Min.sup = 3	bea
cde		de
be		be

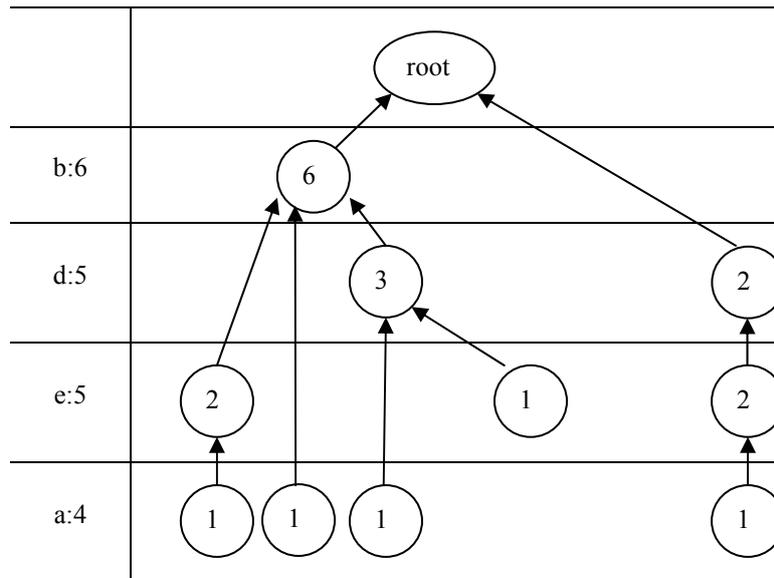


Fig. 3: Frequent Pattern Tree for the sample database

Now, for constructing the FP-tree, a scan over the database is made for adding each item to the tree. The first itemset will be the first branch of the tree. The second transaction shares a common prefix with the already existing set in the tree. In this case, the values along the path of the common prefix will be increased by one, and the remaining items will make new nodes for the tree. In the example given above, the first itemset is *bda*. Hence the first branch of the tree would be the items *b*, *d* and *a*. While adding the second itemset *bde*, a new node is introduced for *e* and also the values along the path are increased by one. The tree for the example database (Han et al., 1999) is shown in Figure 3.

Each node of the FP-Tree consists of three fields (Han et al., 1999):

- item-name*: The name of the item that the node represents is stored.
- count*: The accumulated support of the node within the current path.
- node-link*: Links have to be built between the nodes. It stores the ancestor of the current node, and null if there is none.

For discovering all frequent itemsets, the FP-Growth algorithm takes a look at each level of depth of the tree starting from the bottom and generating all possible itemsets that include nodes in that specific level. A detailed procedure of FP growth algorithm can be found in Han et al. (2004) and Han et al. (1999). The procedure of the algorithm is given below.

Procedure FP-Growth ($Tree, \alpha$)

```

{
1) if  $Tree$  contains a single path  $P$ 
2) then for each combination ( $\beta$ ) of the nodes in the path  $P$  do
3) generate pattern  $\beta \cup \alpha$  with support = minimum support of nodes in  $\beta$ 
4) else for each  $a_i$  in the header of  $Tree$  do {
5) generate pattern  $\beta = a_i \cup \alpha$  with support =  $a_i$  . support;
6) construct  $\beta$ 's conditional pattern base and then  $\beta$ 's conditional FP_Tree,
    $Tree_\beta$ 
7) if  $Tree_\beta \neq \Phi$ 
8) then call FP-growth ( $Tree_\beta, \beta$ )
}
}

```

FP tree algorithm can be used effectively for fuzzy database also. But while dealing with fuzzy database, both support of occurrences and value of membership are to be considered. The membership value is simply added to the overall count of the node.

FP algorithm can be effectively used to find association rules of rare events also, by specifying the rare events as specific target events (Berberidis et al., 2004). FP algorithm is particularly useful for large data bases since it scans the data only twice. However, it is equally efficient for small data bases also since it avoids the candidate generation step of Apriori algorithm. A performance study by Han et al. (2004) has shown that FP growth algorithm is about one order of magnitude faster than the Apriori algorithm and also faster than some recently reported frequent-pattern mining methods, for both short and long frequent patterns.

After constructing the FP tree using the membership values, association rules are generated. The antecedent can consist of any number of items, but in the consequent, only one item is allowed.

In the present study, the ENSO and EQWIN indices are considered as the antecedents and standardized summer monsoon rainfall of the three regions are considered as the consequent. Interestingness of the rules is decided based on the rule confidence. Rules for which the confidence exceeds the minimum confidence threshold are stored.

4. FUZZY RULE BASED PREDICTION

A fuzzy rule based modeling approach is used for the prediction of rainfall by utilizing the linkage between the summer monsoon rainfall and the ENSO-EQWIN indices. The entire dataset from 1958-2006 is divided into two sets—training set (1958-1999) and validation set (2000-2006). Fuzzy rules are constructed using the training set by applying the weighted counting algorithm (Bardossy & Duckstein, 1999). The algorithm consists of the following steps:

1. Classify the attributes (both antecedents and consequent) into five classes. Let the antecedents be X_1 and X_2 and the consequent be Y . The five fuzzy sets for each of the attributes can be expressed as $X_1^1, X_1^2, \dots, X_1^5$; $X_2^1, X_2^2, \dots, X_2^5$ and Y^1, Y^2, \dots, Y^5 .
2. Calculate the membership functions of all the attributes for the entire training set, thus replacing the original training database with a database of membership functions. For each attribute, find the maximum membership value at each data point. Thus, each $X_{i,j}$ ($i=1,2$; $j=1,2,\dots,$ length of training set, n_t) data point possesses a value $M_{i,j}$ which is the maximum membership value. Similarly each Y_j ($j=1,2,\dots, n_t$) possess a value $M_{o,j}$.
3. To combine the effect of all the antecedents, the degree of fulfillment (DOF) of each dataset of the training set is calculated as the product of all M_i 's for each j where $j=1,2,\dots, n_t$. i.e.,

$$DOF_j = \prod_{i=1}^2 M_{i,j} \quad (7)$$

Degree of fulfillment (DOF) indicates the degree of applicability of the rule within the system.

4. For each rule, a weight is assigned to each rule as the product of DOF and the membership value of the consequent. If the rule is repeating in the database, the weights are added upon. Hence, weight for k^{th} rule can be expressed as

$$wt_k = \sum_{j=1}^{n_t} .DOF_j .M_{o,j} \quad (8)$$

After scanning through the entire training set, all the derived rules will have weights assigned to them.

5. Validation of the rules: Calculate the DOF of all attributes for all data points in the validation set. Identify those rules that are showing similar antecedent conditions and also having a minimum DOF value, from the derived ones.
6. Defuzzification: Each rule leads to a fuzzy response. A crisp output can be obtained from all these rules through defuzzification process. In weighted counting algorithm, the center of gravity is commonly used to obtain the estimated value of the consequent variable. Hence, the estimated value can be expressed as

$$\hat{Y}_j = \frac{\sum_{k=1}^n .DOF_k .wt_k .B_k^{(2)}}{\sum_{k=1}^n .DOF_k .wt_k} \quad (9)$$

where k = number of similar rules extracted for a given antecedent combination;
 $B_k^{(2)}$ = most likely value of the consequent of the k^{th} rule, i.e, the value for which the membership function is one.

5. RESULTS

5.1 Fuzzy Association Rules

FP growth algorithm is applied on the entire dataset from 1958 to 2006 to derive fuzzy rules. Rules are extracted for drought, flood and normal conditions based on the confidence measure. The rules extracted and the corresponding confidences for the three regions are shown in Tables 3 to 5.

The rules for the extremes (droughts and floods) are in concordance with the previous works done by Gadgil et al. (2004) and Maity and Nagesh Kumar (2008) combining both ENSO and EQWIN indices for the prediction of ISMR. Considering All-India, no rules are extracted for extreme flood condition, because upon reducing the minimum confidence level, it is found that a combination of Nino1 + Eqwin3 leads to an extreme flood condition, but with a very less confidence of 0.38. Similarly, for other two regions also, rules with high confidence do not exist for some conditions. However, rules for all the three regions are showing a negative relation with Nino index and a positive relation with EQWIN

index.

TABLE 3

Fuzzy association rules for All-India

Antecedent	Consequent	Confidence
Nino5, Eqwin1	Severe drought	0.95
Nino5, Eqwin1	Moderate drought	0.83
Nino5, Eqwin3	Moderate drought	0.66
Nino1, Eqwin4	Moderate flood	1.0
Nino2, Eqwin5	Moderate flood	0.83
Nino1, Eqwin1	Normal	1.0
Nino2, Eqwin1	Normal	0.79
Nino5, Eqwin5	Normal	1.0

TABLE 4

Fuzzy association rules for Peninsular region

Antecedent	Consequent	Confidence
Nino5, Eqwin1	Severe drought	1.0
Nino4, Eqwin2	Moderate drought	0.66
Nino1, Eqwin3	Moderate flood	1.0
Nino1, Eqwin4	Severe flood	0.99
Nino5, Eqwin4	Normal	0.75
Nino3, Eqwin1	Normal	0.68
Nino1, Eqwin1	Normal	1.0

TABLE 5

Fuzzy association rules for West central region

Antecedent	Consequent	Confidence
Nino5, Eqwin1	Moderate drought	1.0
Nino5, Eqwin3	Moderate drought	0.67
Nino1, Eqwin3	Moderate flood	1.0
Nino2, Eqwin5	Moderate flood	0.98
Nino5, Eqwin4	Normal	0.75
Nino2, Eqwin1	Normal	0.81
Nino1, Eqwin1	Normal	1.0

5.2 Fuzzy Rule Based Prediction

Rules are defined based on the training set for the period 1958 to 1999. Validation is done for the period 2000 to 2006. Prediction is done using the extracted rules of the antecedent combinations in the validation set. For example the antecedent combination for the first data point of All-India, in the validation set is Nino2 and Eqwin4. The maximum membership values are 0.78 and 1, respectively. Now from the training database, six rules have been extracted that have similar or nearby antecedent combinations. The rules are given in table 6.

TABLE 6

Selected Rules for Predicting the First Dataset in the Validation database

Antecedent	Consequent	DOF	Weight	B ⁽²⁾
Nino2,Eqwin4	Moderate flood	0.78	1.95	1
Nino2,Eqwin4	Normal	0.81	0.98	0
Nino3, Eqwin4	Moderate flood	0.88	4.97	1
Nino3, Eqwin4	Normal	0.67	5.38	0
Nino3, Eqwin4	Moderate drought	0.57	0.76	-1
Nino3, Eqwin4	Severe flood	0.97	1.67	2

Using the Eqn. (9), the crisp value of the consequent can be calculated as follows:

$$\hat{Y}_j = \frac{(0.78 \times 1.95 \times 1) + (0.81 \times 0.98 \times 0) + (0.88 \times 4.97 \times 1) + (0.67 \times 5.38 \times 0) + (0.57 \times 0.76 \times -1) + (0.97 \times 1.67 \times 2)}{(0.78 \times 1.95) + (0.81 \times 0.98) + (0.88 \times 4.97) + (0.67 \times 5.38) + (0.57 \times 0.76) + (0.97 \times 1.67)} \quad (10)$$

$$= 0.704.$$

The predicted value is much nearer to the original value of 0.86. The consequent values are predicted for the entire validation set in a similar fashion for all the three regions. The predicted values for the period 2000 to 2006 for the three regions are shown in Figures 4 to 6. The correlation coefficients between the observed monthly summer monsoon rainfall and the predicted monthly summer monsoon rainfall and the corresponding root mean square errors are shown in table 7.

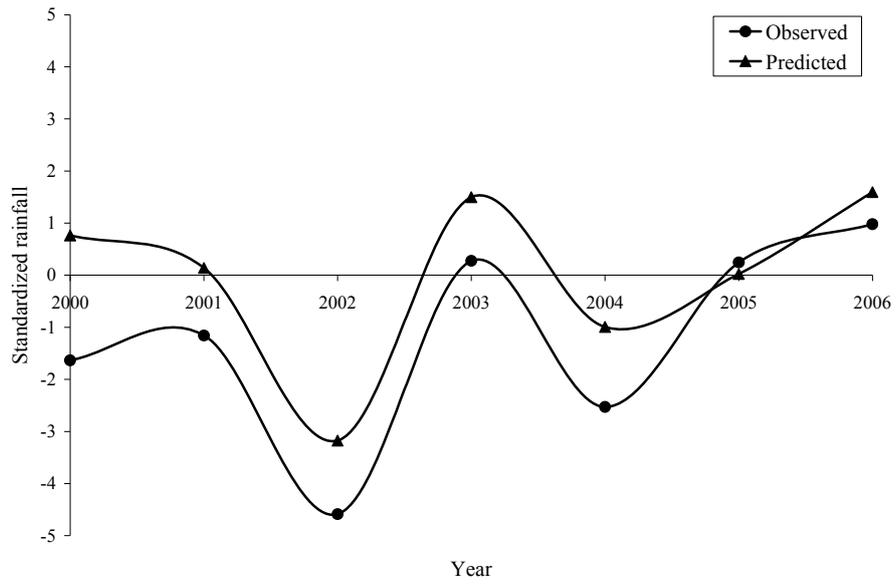


Fig. 4: Observed and predicted summer monsoon rainfall values of All-India region for the validation period

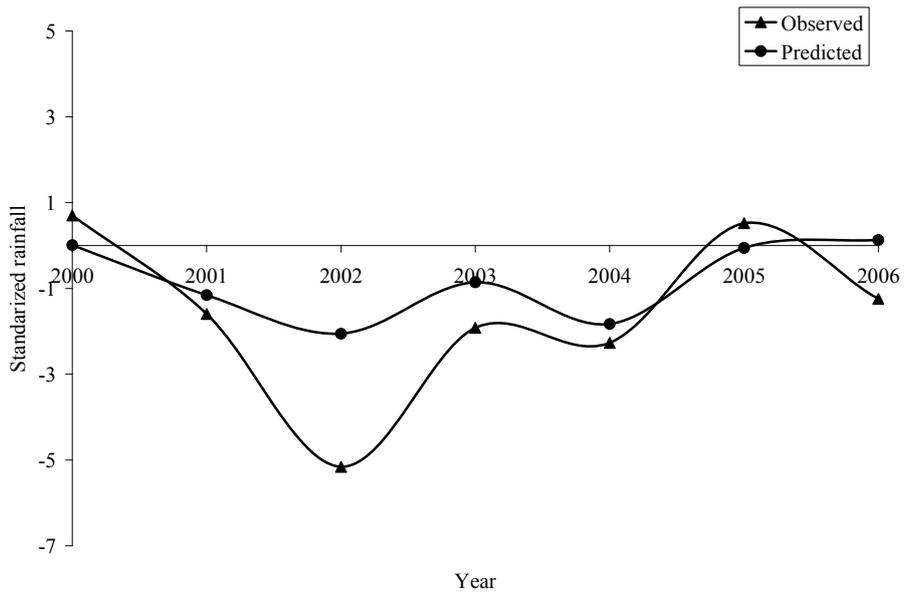


Fig. 5: Observed and predicted summer monsoon rainfall values of Peninsular region for the validation period

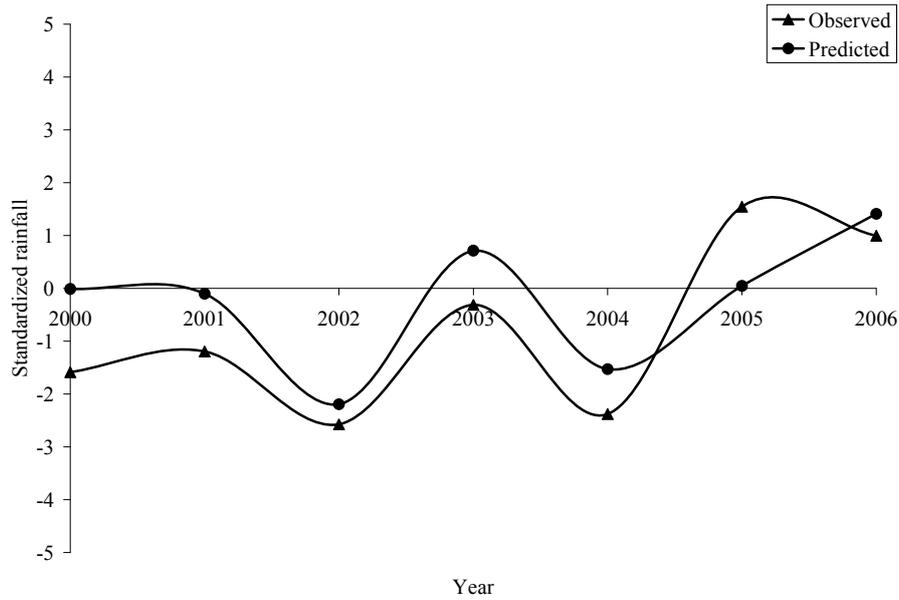


Fig. 6: Observed and predicted summer monsoon rainfall values of West central region for the validation period

TABLE 7

Correlation values and root mean square errors between the observed and predicted summer monsoon rainfall

Region	Correlation		Root mean square error
	Monthly rainfall	Summer monsoon season	
All- India	0.63	0.90	1.02
Peninsular	0.77	0.85	0.71
West central	0.51	0.78	0.99

It is evident from the results that the model performs quite well in the validation period. Since rainfall is a complex phenomena influenced by many climatic and

atmospheric indices, the inclusion of more indices in the antecedent group may improve the performance of the model.

6. CONCLUSIONS

Fuzzy association rule helps in extracting the relation between the variables by overcoming the sharp boundary problem when mining association rules from quantitative data. The rules extracted using the algorithm give a one to one relation of the antecedents with the consequent, without going into the complex mathematical relations between them. The support of the rules is not taken into account because for the extreme rainfall events, support may be very small, when compared with the normal rainfall events and thus may lead to the omission of these infrequent events.

Fuzzy rule based prediction using the weighted counting algorithm provides an excellent tool for predicting the monthly rainfall in monsoon months using the two indices, ENSO and EQWIN. The variability of summer monsoon rainfall for the three regions—All India, Peninsular, and West central—is reasonably captured using this method. This approach is a better substitute for the usual statistical methods, which are limited mostly by the linear statistical relation between the indices.

REFERENCES

- Bardossy, A., Duckstein, L. 1999. *Fuzzy rule-based modeling with applications to geophysical, biological and engineering systems*, CRC, New York.
- Berberidis, C., Angelis, L., Vlahavas, I. 2004. PREVENT An algorithm for mining intertransactional patterns for the prediction of rare events, *Proc. Second Starting AI Researchers' Symposium, Frontiers in Artificial Intelligence and Applications*, **9**.
- Caulfield, H.J. 1996. Fuzzy optical meteorology, *IEEE Transactions on Fuzzy Systems*, **4(2)**.

- Chen, Z., Chen, G. 2007. An approach to classification based on fuzzy association rules, *Advances in Intelligent Systems Research*, International Conference on Intelligent Systems and Knowledge Engineering .
- Gadgil S. 2003. The Indian Monsoon and its variability. *Annual Review of Earth Planetary Science*, **31**, 429-67.
- Gadgil, S., Rajeevan, M., Nanjundiah, R. 2005. Monsoon prediction: why yet another failure, *Current Science*, **88(9)**, 1389-1400.
- Gadgil, S., Vinayachandran, P. N., Francis, P. A., Gadgil, S. 2004. Extremes of the Indian summer monsoon rainfall, ENSO and equatorial Indian Ocean oscillation, *Geophys. Res. Lett.*, **31**, doi:10.1029/2004GL019733.
- Han, J., Pei, J., Yin, Y., Mao, R. 2004. Mining frequent patterns without candidate generation: A frequent-pattern tree approach, *Data Mining and Knowledge Discovery*, **8**, 53-87.
- Han, J., Pei, J., Yin, Y. 1999. Mining frequent patterns without candidate generation, *2000 ACM SIGMOD Intl. Conference on Management of Data*, ACM Press.
- Maity, R., Nagesh Kumar, D., Nanjundiah, R.S. 2007. Review of hydroclimatic teleconnection between hydrologic variables and large-scale atmospheric circulation indices with Indian perspective, *ISH Journal of Hydraulic Engineering*, **13(1)**, 77-92.
- Maity, R., Nagesh Kumar, D. 2006. Hydroclimatic association of monthly summer monsoon rainfall over India with large-scale atmospheric circulation from tropical Pacific Ocean and Indian Ocean region, *Atmospheric Science Letters*, **7(4)**, 101-7.
- Rajeevan, M., Pai, D.S., Diskhit, S.K., Kelkar, R.R. 2004. IMD's new operational models for long range forecast of southwest monsoon rainfall over India and their verification for 2003, *Current Science*, **86(3)**, 422-30.