# Measuring Interesting Amino Acid Patterns for Alzheimer's Disease Related Studies Targets on the Binding Site Using Association Rule Mining

D. Ponmary Pushpa Latha[1*] and D. Joseph Pushpa Raj[2]

[1]Department of Computer Applications, Karunya University, Coimbatore-641 114, India.
[2]Department of Information Technology, Francis Xavier Engineering College, Tirunelveli- 627 011, India.

## ARTICLE INFO

## ABSTRACT

Data mining techniques are used in various areas like stock exchange, education, bioinformatics, health care etc. The main purpose of data mining techniques is used to extract the useful and interesting information. Association Rule Mining (ARM) associates the different attributes and gives the most suitable rules from large database. Protein ligand binding is an important step in enzymatic mechanisms and in drug discovery. This research work gives the association rules for amino acid residues which are present in the binding site of Alzheimer's Disease Related Studies targets. The data are collected from Protein Data Bank. Association rule mining is applied in the Alzheimer's Disease Related Studies protein and the interesting rules for the amino acid residues present in the Alzheimer's Disease Related Studies are formed. Ile and Ser are the amino acid residues which are having major role in the precedence of association rules of Alzheimer's Disease Related Studies. This research work may support in identify new binding protein-ligand pairs and predict protein ligand binding in particular diseases.

## INTRODUCTION

Alzheimer's disease (AD) is a brain disorder that brings disturbances in reasoning, planning, language and perception. Increased age, High blood pressure, coronary artery disease and diabetes are the risk factor for Alzheimer's disease.

In the current scenario, 18 million people are affected by Alzheimer's disease in India. In the year 2025, the induced rate of Alzheimer's disease in India will be 34 million (Rao and Shaji, 2007).

Data mining tools and techniques provides significant knowledge for the various applications. The knowledge which is produced by the data mining techniques is massively useful in the disease related studies. Most of the data mining algorithms results in patterns. A pattern is a set of attributes which are supported by a number of transactions.

These patterns are useful to find the association between the attribute. Association rule mining is one of the important techniques in data mining which helps to produce the hidden knowledge in the massive data set (Al-Shalabi, 2011).

_____

* Corresponding Author
*Department of Computer Applications, Karunya University, Coimbatore-641 114, India, Email: josesph_raj@yahoo.com*

## Association Rule Mining (ARM)

ARM is used to find the correlation relationship among the data items in the large data set. In the sales application or retail industry, consider that let D be the set of n transactions such that D={T1, T2, T3,…. Tn} Where Ij= I and I be the set of items. I =(i1, i2…..im) Let X, Y and Z be the three item sets in the I. If an association rule is formed like $X \cap Y \Rightarrow Z$ then X and Y are antecedent and Z is called as consequent. Some of the validation measures like support, confidence and lift are used in the ARM. Support is the number of transaction of X=>Y in the XUY number of transaction. Confidence is the number of transaction which satisfies the individual association rule in the total number of transaction. The relationship between X and Y is quantified in lift (Ramaraj *et al.,* 2009).

## ARM in Bioinformatics

(Sallab *et al.,* 2004) developed a tool which is named by QuantMiner. This tool is used to mine the quantitative rules for atheoscderosis data set. This tool is designed in such a way that to support both categorical and numerical attribute (Ordonex *et al.,* 2006 & 2000), is utilized association rules to find absence or existence of heart diseases by giving the set of rules with high-qualitymetrics.

(Gasmi *et al.,* 2005) extracted the association rules from SAGE data set. (Kwasnicka and Switalski, 2006) merged the association rules with genetic algorithms for medical database. (Basemann *et al.,* 2004) applied functions region association in Protein-Protein Interaction.

(Lopez *et al.,* 2007) generated fuzzy based association rules for medical domains. (Li *et al.,* 2005) mined the risk patterns for medical domains. (Li *et al.,* 2005) applied association rule mining in biological duplicate detection. (Nehemiah *et al.,* 2007) applied association rule mining in medical datasets in order to support physician in decision making. (Ohsaki *et al.,* 2003) formed interesting rules from chronic hepatitis dataset. (Gupta *et al.,* 2006; Gupta and Agrawal 2009) applied association rule mining in the amino acid residues.

## ARM in Amino Acid Residues

The size of structural information deposited in the Protein Data Bank (Berman *et al.,* 2000) increases day by day. It leads to have automated *in silico* studies involved thousands of protein ligand complexes and binding site. The structural classifications of binding sites on protein-surfaces are applicable for the predication and modeling of protein ligand interactions. Since many known biologically active compounds are ligands bound to proteins, this research work is important in the mathematical foundations of drug discovery and drug design. In the present work, data-mining techniques are applied for the amino acids set, shaped from the residues at each active sites present in the disease specific approach.

Analysis, classification and characterization of binding sites are important in predicting and designing enzymatic mechanisms, since protein ligand complex is a key step in enzymatic mechanisms. (Chen *et al.,* 2004) presented a novel unsupervised learning approach to discover frequent patterns in the protein families, based on biochemical, geometric and dynamic features. Without any prior knowledge of functional motifs, the method finds the frequent patterns for each type of amino acid and identifies the conserved residues in three protease subfamilies; chymotrypsin and subtilisin subfamilies of serine proteases and papain subfamily of cysteine proteases. The catalytic triad residues are notable by their strong spatial coupling (high interconnectivity) to other conserved residues.

Although the spatial arrangements of the catalytic residues in the two subfamilies of serine proteases are similar, their frequent patterns are found to be quite different. The present approach appears to be a promising tool for detecting functional patterns in rapidly growing structure databases and providing insights into the relationship among protein structure, dynamics and function.

(Gabor Ival *et al.,* 2007) analyzed the residue composition of the binding sites in the entire PDB for frequency and for unseen association rules. The following are the results of the paper: (i) the cleaning and repairing algorithm (ii) redundancy elimination from the data (iii) application of association rule mining to the cleaned non-redundant data set. Gobal Ivan created

numerous significant relations of the residue-composition of the ligand binding sites on protein surfaces (Kuo *et al.,* 2011) propose a method to find out the association relationship among amino acid residues on binding sites.

Such knowledge of binding sites is very obliging in predicting protein-protein interactions. Protein complexes which have protein-protein recognition are focused to find out the association relationship among amino acid residues. The association rule mining technique is used to discover geographically adjacent amino acids on a binding site of a protein complex.

(Pant *et al.,* 2012) present association based rules formulation for the most frequently occurring amino acids in HIV viruses to analyze the functioning of this virus. Most people in the world are affected by HIV disease and efforts are taken throughout the world to build new vaccines and drugs. Apriori algorithm is applied to find the frequent item set in the amino acid residues, since these residues can be a source for good drug targets. Severities in liver disease are varying in a huge number of patients due to the plasma amino-acid concentrations.

In chronic active hepatitis, the following amino acid patterns found in Liver diseases: Plasma concentrations of aspartate, threonine, serine, methionine, and the aromatic amino-acid tyrosine were significantly raised, while concentrations of proline and of the three branched chain amino-acids valine, isoleucine, and leucine were significantly reduced (Marsha, *et al.,* 1982). The present research work deals about the Alzheimer's disease.

## METHODOLOGY

The aim of this research work is to find Alzheimer's disease patterns of amino-acid residues in the protein binding site using association rule mining (Fig. 1, 2). The following are the steps to find the association rules.

Step 1: The disease causing targets are segregated from the Protein Data Bank and categorize the proteins into domain.
Step 2: A tool is developed to automate data extraction from the Protein Data Bank.
Step 3: Calculate the protein ligand interactions within 6 Angstrom.
Step 4: Remove the redundant data from the protein ligand interaction by finding the distinct amino acid residues.
Step 5: Association Rule Mining models are used to mine the protein ligand interaction data to find the amino acid association.

The following tools and techniques are used to obtain the results. Java (Schildt, 2003) is used to calculate the protein ligand interaction. MySQL (http://dev.mysql.com) is used as the back end to store the data.

Rapid miner tool (http://rapid-i.com) is used to find the association rules. Frequency of amino acid residues are analyzed by the polyAnalyst Tool (www.megaputer.com).

```
┌─────────────────────────────────────────┐
│        Collect disease causing target     │
└─────────────────────────────────────────┘
                    │
                    ▼
┌─────────────────────────────────────────┐
│   Data Extraction and Translation into   │
│           Relational Database             │
└─────────────────────────────────────────┘
                    │
                    ▼
┌─────────────────────────────────────────┐
│           Find the Binding Site           │
└─────────────────────────────────────────┘
                    │
                    ▼
┌─────────────────────────────────────────┐
│        Eliminate the redundant data       │
└─────────────────────────────────────────┘
                    │
                    ▼
┌─────────────────────────────────────────┐
│   Apply association rule mining to        │
│   discover amino acid patterns            │
└─────────────────────────────────────────┘
```
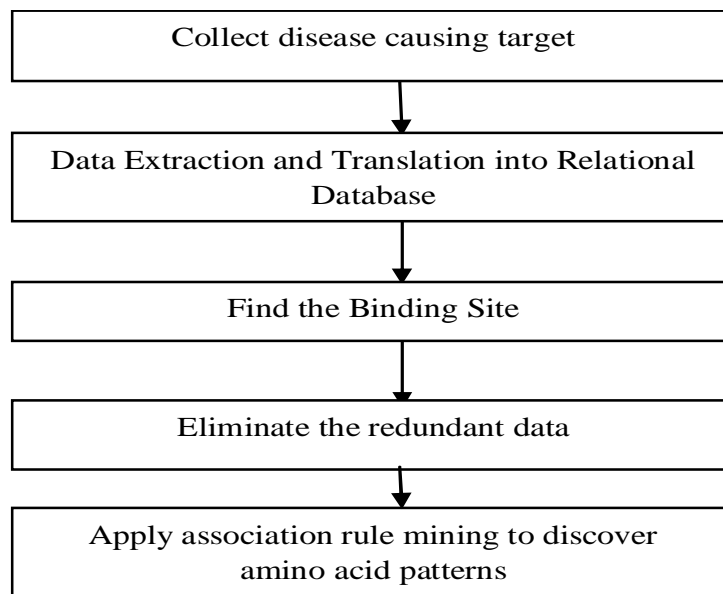
**Fig. 1:** Stages involved in finding disease specific patterns of amino-acid residues in the protein binding site using association rule mining.
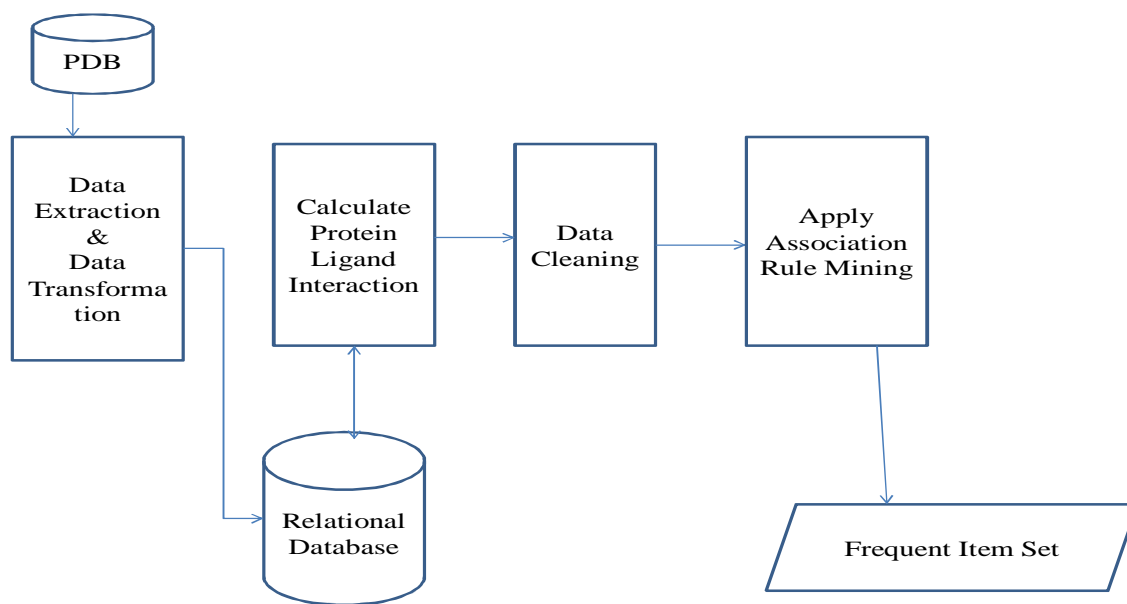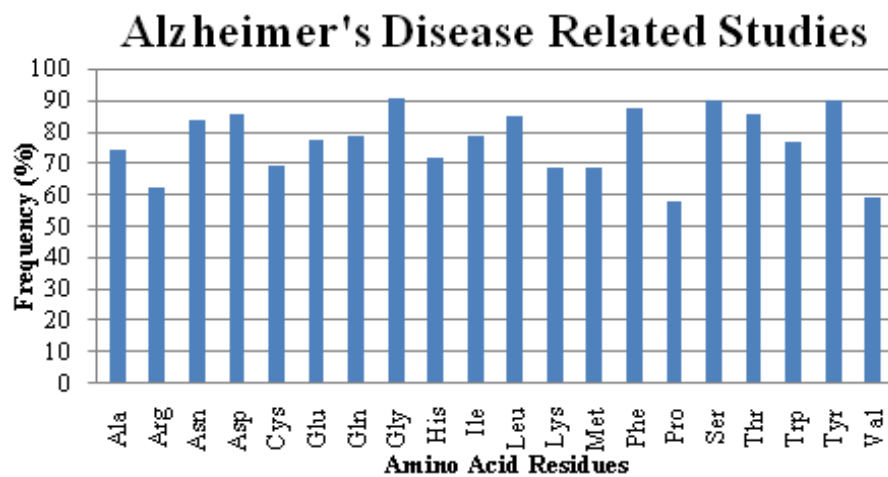


**Fig. 2:** Design Model.



**Fig. 3:** Composition of Amino Acid Resdues.

**Table . 1:** Data set of 98 Protein.

| | | | | | |
|---|---|---|---|---|---|
| 1BPT | 1CBW | 1T7C | 1T8L | 1T8M | 1T8N |
| 1T8O | 1TAW | 2V00 | 1OQN | 2IPT | 2IPU |
| 2IQ9 | 2IQA | 2R0W | 2R0Z | 3BKC | 3BKM |
| 1DX6 | 1E3Q | 1E66 | 1EVE | 1GPK | 1GPN |
| 1GQR | 1GQS | 1H23 | 1OCE | 1ODC | 1QTI |
| 1UT6 | 1VOT | 1W4L | 1W6R | 1W75 | 1W76 |
| 2BAG | 2CEK | 2CKM | 2CMF | 2J3D | 2J3Q |
| 2V96 | 2V97 | 2V98 | 2VA9 | 2VJA | 2VJB |
| 2VJC | 2VJD | 2VQ6 | 2VT6 | 2VT7 | 1AQC |
| 1FKN | 1H4W | 1M4H | 1MX1 | 1SO8 | 1U7T |
| 1W50 | 1W51 | 1X11 | 1XN2 | 1XN3 | 1XS7 |
| 1YM2 | 1YM4 | 2DYQ | 2F3E | 2F3F | 2FK1 |
| 2FK2 | 2FK3 | | | | |

**Table. 2:** The frequency of residues on the binding sites.

| Residue | Frequency (%) | Residue | Frequency (%) | Residue | Frequency (%) |
|---|---|---|---|---|---|
| Ala | 74.4 | Gly | 90.8 | Pro | 58.1 |
| Arg | 62.2 | His | 71.4 | Ser | 89.7 |
| Asn | 83.6 | Ile | 78.5 | Thr | 85.7 |
| Asp | 85.7 | Leu | 84.6 | Trp | 76.5 |
| Cys | 69.3 | Lys | 68.3 | Tyr | 89.7 |
| Glu | 77.5 | Met | 68.3 | Val | 59.1 |
| Gln | 78.5 | Phe | 87.7 | | |

**Table 3:** Association Rules for Alzheimer's Disease Related Studies Target

| SNo. | Premises | Conclusion | Support | Confidence | Laplace | Gain | p-s | Lift |
|---|---|---|---|---|---|---|---|---|
| 1 | Trp | Ser | 0.765306 | 1 | 1 | -0.76531 | 0.078092 | 1.113636 |
| 2 | Phe, Ile | Ser | 0.744898 | 1 | 1 | -0.7449 | 0.07601 | 1.113636 |
| 3 | Phe, Gln | Ser | 0.72449 | 1 | 1 | -0.72449 | 0.073928 | 1.113636 |
| 4 | Phe, Trp | Ser | 0.744898 | 1 | 1 | -0.7449 | 0.07601 | 1.113636 |
| 5 | Thr, Asp | Ser | 0.77551 | 1 | 1 | -0.77551 | 0.079134 | 1.113636 |
| 6 | Thr, Ile | Ser | 0.744898 | 1 | 1 | -0.7449 | 0.07601 | 1.113636 |
| 7 | Thr, Trp | Ser | 0.72449 | 1 | 1 | -0.72449 | 0.073928 | 1.113636 |
| 8 | Asp, Gln | Ser | 0.714286 | 1 | 1 | -0.71429 | 0.072886 | 1.113636 |
| 9 | Asp, Trp | Ser | 0.72449 | 1 | 1 | -0.72449 | 0.073928 | 1.113636 |
| 10 | Leu, Ile | Ser | 0.734694 | 1 | 1 | -0.73469 | 0.074969 | 1.113636 |
| 11 | Leu, Gln | Ser | 0.755102 | 1 | 1 | -0.7551 | 0.077051 | 1.113636 |
| 12 | Leu, Trp | Ser | 0.683673 | 1 | 1 | -0.68367 | 0.069763 | 1.113636 |
| 13 | Ile, Gln | Ser | 0.714286 | 1 | 1 | -0.71429 | 0.072886 | 1.113636 |
| 14 | Ile, Trp | Ser | 0.714286 | 1 | 1 | -0.71429 | 0.072886 | 1.113636 |
| 15 | Gln, Trp | Ser | 0.683673 | 1 | 1 | -0.68367 | 0.069763 | 1.113636 |
| 16 | Phe, Thr, Asp | Ser | 0.744898 | 1 | 1 | -0.7449 | 0.07601 | 1.113636 |
| 17 | Phe, Thr, Ile | Ser | 0.714286 | 1 | 1 | -0.71429 | 0.072886 | 1.113636 |
| 18 | Phe, Thr, Gln | Ser | 0.693878 | 1 | 1 | -0.69388 | 0.070804 | 1.113636 |
| 19 | Phe, Thr, Trp | Ser | 0.704082 | 1 | 1 | -0.70408 | 0.071845 | 1.113636 |
| 20 | Phe, Asp, Ile | Ser | 0.704082 | 1 | 1 | -0.70408 | 0.071845 | 1.113636 |
| 21 | Phe, Asp, Gln | Ser | 0.683673 | 1 | 1 | -0.68367 | 0.069763 | 1.113636 |
| 22 | Phe, Asp, Trp | Ser | 0.704082 | 1 | 1 | -0.70408 | 0.071845 | 1.113636 |
| 23 | Phe, Leu, Ile | Ser | 0.704082 | 1 | 1 | -0.70408 | 0.071845 | 1.113636 |
| 24 | Phe, Leu, Gln | Ser | 0.704082 | 1 | 1 | -0.70408 | 0.071845 | 1.113636 |
| 25 | Phe, Leu, Trp | Ser | 0.663265 | 1 | 1 | -0.66327 | 0.06768 | 1.113636 |
| 26 | Phe, Ile, Gln | Ser | 0.693878 | 1 | 1 | -0.69388 | 0.070804 | 1.113636 |
| 27 | Phe, Ile, Trp | Ser | 0.693878 | 1 | 1 | -0.69388 | 0.070804 | 1.113636 |
| 28 | Phe, Gln, Trp | Ser | 0.663265 | 1 | 1 | -0.66327 | 0.06768 | 1.113636 |
| 29 | Thr, Asp, Leu | Ser | 0.714286 | 1 | 1 | -0.71429 | 0.072886 | 1.113636 |
| 30 | Thr, Asp, Ile | Ser | 0.704082 | 1 | 1 | -0.70408 | 0.071845 | 1.113636 |
| 31 | Thr, Asp, Gln | Ser | 0.673469 | 1 | 1 | -0.67347 | 0.068721 | 1.113636 |
| 32 | Thr, Asp, Trp | Ser | 0.683673 | 1 | 1 | -0.68367 | 0.069763 | 1.113636 |
| 33 | Thr, Leu, Ile | Ser | 0.714286 | 1 | 1 | -0.71429 | 0.072886 | 1.113636 |
| 34 | Thr, Leu, Gln | Ser | 0.72449 | 1 | 1 | -0.72449 | 0.073928 | 1.113636 |
| 35 | Thr, Leu, Trp | Ser | 0.663265 | 1 | 1 | -0.66327 | 0.06768 | 1.113636 |
| 36 | Thr, Ile, Gln | Ser | 0.693878 | 1 | 1 | -0.69388 | 0.070804 | 1.113636 |
| 37 | Thr, Ile, Trp | Ser | 0.683673 | 1 | 1 | -0.68367 | 0.069763 | 1.113636 |
| 38 | Thr, Gln, Trp | Ser | 0.653061 | 1 | 1 | -0.65306 | 0.066639 | 1.113636 |
| 39 | Asp, Leu, Ile | Ser | 0.693878 | 1 | 1 | -0.69388 | 0.070804 | 1.113636 |
| 40 | Asp, Leu, Gln | Ser | 0.693878 | 1 | 1 | -0.69388 | 0.070804 | 1.113636 |
| 41 | Asp, Ile, Gln | Ser | 0.673469 | 1 | 1 | -0.67347 | 0.068721 | 1.113636 |
| 42 | Asp, Ile, Trp | Ser | 0.673469 | 1 | 1 | -0.67347 | 0.068721 | 1.113636 |
| 43 | Leu, Ile, Gln | Ser | 0.704082 | 1 | 1 | -0.70408 | 0.071845 | 1.113636 |
| 44 | Leu, Ile, Trp | Ser | 0.673469 | 1 | 1 | -0.67347 | 0.068721 | 1.113636 |
| 45 | Leu, Gln, Trp | Ser | 0.663265 | 1 | 1 | -0.66327 | 0.06768 | 1.113636 |
| 46 | Ile, Gln, Trp | Ser | 0.673469 | 1 | 1 | -0.67347 | 0.068721 | 1.113636 |

| 47 | Thr, Gln, Trp | Ile | 0.653061 | 1 | 1 | -0.65306 | 0.139942 | 1.272727 |
|----|----|----|----|----|----|----|----|----|
| 48 | Leu, Gln, Trp | Ile | 0.663265 | 1 | 1 | -0.66327 | 0.142128 | 1.272727 |
| 49 | Phe, Thr, Asp, Leu | Ser | 0.683673 | 1 | 1 | -0.68367 | 0.069763 | 1.113636 |
| 50 | Phe, Thr, Asp, Ile | Ser | 0.673469 | 1 | 1 | -0.67347 | 0.068721 | 1.113636 |
| 51 | Phe, Thr, Asp, Gln | Ser | 0.653061 | 1 | 1 | -0.65306 | 0.066639 | 1.113636 |
| 52 | Phe, Thr, Asp, Trp | Ser | 0.663265 | 1 | 1 | -0.66327 | 0.06768 | 1.113636 |
| 53 | Phe, Thr, Leu, Ile | Ser | 0.683673 | 1 | 1 | -0.68367 | 0.069763 | 1.113636 |
| 54 | Phe, Thr, Leu, Gln | Ser | 0.683673 | 1 | 1 | -0.68367 | 0.069763 | 1.113636 |
| 55 | Phe, Thr, Ile, Gln | Ser | 0.673469 | 1 | 1 | -0.67347 | 0.068721 | 1.113636 |
| 56 | Phe, Thr, Ile, Trp | Ser | 0.663265 | 1 | 1 | -0.66327 | 0.06768 | 1.113636 |
| 57 | Phe, Asp, Leu, Ile | Ser | 0.663265 | 1 | 1 | -0.66327 | 0.06768 | 1.113636 |
| 58 | Phe, Asp, Leu, Gln | Ser | 0.663265 | 1 | 1 | -0.66327 | 0.06768 | 1.113636 |
| 59 | Phe, Asp, Ile, Gln | Ser | 0.653061 | 1 | 1 | -0.65306 | 0.066639 | 1.113636 |
| 60 | Phe, Asp, Ile, Trp | Ser | 0.653061 | 1 | 1 | -0.65306 | 0.066639 | 1.113636 |
| 61 | Phe, Leu, Ile, Gln | Ser | 0.683673 | 1 | 1 | -0.68367 | 0.069763 | 1.113636 |
| 62 | Phe, Leu, Ile, Trp | Ser | 0.653061 | 1 | 1 | -0.65306 | 0.066639 | 1.113636 |
| 63 | Phe, Ile, Gln, Trp | Ser | 0.653061 | 1 | 1 | -0.65306 | 0.066639 | 1.113636 |
| 64 | Thr, Asp, Leu, Ile | Ser | 0.673469 | 1 | 1 | -0.67347 | 0.068721 | 1.113636 |
| 65 | Thr, Asp, Leu, Gln | Ser | 0.663265 | 1 | 1 | -0.66327 | 0.06768 | 1.113636 |
| 66 | Thr, Asp, Ile, Gln | Ser | 0.653061 | 1 | 1 | -0.65306 | 0.066639 | 1.113636 |
| 67 | Thr, Leu, Ile, Gln | Ser | 0.683673 | 1 | 1 | -0.68367 | 0.069763 | 1.113636 |
| 68 | Thr, Leu, Ile, Trp | Ser | 0.653061 | 1 | 1 | -0.65306 | 0.066639 | 1.113636 |
| 69 | Thr, Gln, Trp | Ser, Ile | 0.653061 | 1 | 1 | -0.65306 | 0.146606 | 1.289474 |
| 70 | Ser, Thr, Gln, Trp | Ile | 0.653061 | 1 | 1 | -0.65306 | 0.139942 | 1.272727 |
| 71 | Thr, Ile, Gln, Trp | Ser | 0.653061 | 1 | 1 | -0.65306 | 0.066639 | 1.113636 |
| 72 | Asp, Leu, Ile, Gln | Ser | 0.663265 | 1 | 1 | -0.66327 | 0.06768 | 1.113636 |
| 73 | Leu, Gln, Trp | Ser, Ile | 0.663265 | 1 | 1 | -0.66327 | 0.148896 | 1.289474 |
| 74 | Ser, Leu, Gln, Trp | Ile | 0.663265 | 1 | 1 | -0.66327 | 0.142128 | 1.272727 |
| 75 | Leu, Ile, Gln, Trp | Ser | 0.663265 | 1 | 1 | -0.66327 | 0.06768 | 1.113636 |
| 76 | Phe, Thr, Leu, Ile, Gln | Ser | 0.663265 | 1 | 1 | -0.66327 | 0.06768 | 1.113636 |

## RESULT AND DISCUSSION

There are totally 98 PDB files (Table 1) are segregated from protein data bank, the segregated PDB files are related with Alzheimer's disease related studies. The following residues frequency (Fig 3) are high in the binding site of Alzheimer's disease related studies (Table 2): Asn (83.6%), Asp (85.7%), Gly (90.8%), Leu (84.6%), Phe (87.7%), Ser (89.7%), Thr (85.7%) and Tyr (89. 7%).

The following residues frequency are low in the binding site of Alzheimer's disease related studies: Pro (58.1%), Val (59.1%) and Arg (62.2%). Association rules are formed with the following cut off value (Table 3): confidence value is equal to one, the support value is above six and the laplace value is equal to one. In the formed association rules, Ser and Ile are conclusions for many rules. The lift value is above or equal to 1.1 for all the rules which were formed for the confidence value is equal to one.

## CONCLUSION

In this study, Association rule mining is applied for the Alzheimer's Disease Related Studies in the binding site. Interesting amino acid patterns are found using this study. {Pro} and {Val} are the fewest amino acid residues in the binding site. {Arg}, {Lys} and {Met} are also rated low in the probability of appearance. {Phe}, {Ser} and {Tyr} are high in the binding site. Totally there are 397 association rules are formed for the cutoff value which is greater than or equal to 0.95. Out of 397 association rules, 76 rules are having confidence value is one.   The association rules which  are having confidence value is one are listed in the table 2.

The present study gives the relationship among the amino acid residues for the Alzheimer's Disease Related Studies in the binding site, This will be helpful in the computation design of new drugs.

## REFERENCES

Al-Shalabi L. Knowledge discovery process: Guide lines for new researchers. J Artifi Intell,  2011;4:21-28.

Basemann C, Denton, A, Yekkirala A. Differential association rule mining for the study of protein-protein interaction networks, Proceedings of the 4th Workshop on Data Mining in Bioinformatics, Seattle, 2004; Aug. 22, 72-80.

Berman HM, Westbrook J, Feng Z, Gilliland G, Bhat TN, Weissig H, Shindyalov IN, Bourne PE. The Protein Data Bank. *Nucleic Acids Research,* 2000; 28 (1): 235-242.

Gasmi GT,  Hamrouni S,  Abdelhak S. Ben Yahia, Nguifo EM. Extracting generic basis of association rules from SAGE data. Proceedings of the 8th International ECML/PKDD Workshop Discovery Challenge, Porto, Portugal, 2005; Oct. 7,1-6.

Gupta NN,  Mangal K, Tiwari and  Mitra P. Mining Quantitative Association Rules in Protein  Sequences. Lecture Notes on Artificial Intelligence, 2006;  3755: 273-281.

Gupta RK, Agrawal DP. Improving the performance of association rule mining algorithms by filtering  insignificant transactions dynamically. Asian J. Inform. Manage, 2009; 3:7-17.

Kwasnicka H,  Switalski  K. Discovery of association rules from medical data-clasical and evolutionary approaches. Proceedings of the 21th Autumn Meeting of Polish Information Processing Society Conference, 2006; 163-177.

Li J,  Wai-Chee Fu A,  He H,  Chen J,  Jin H. Mining risk patterns in medical data.  Proceedings of the 11th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2005; Aug. 21-24, 770-775.

Lopez, FJ, Blanco A,  Garcia F,  Marin A. Extracting biological knowledge by fuzzy association rule mining. Proceedings of the IEEE International Conference on Fuzzy Systems,  2007; July 23-26, 1-6.

Marsha Y Morgan, Marshall AW, Judith P  Milsom, Sheila Sherlock.  Plasma amino-acid patterns in liver disease. Gut, 1982; 23, 362-370.

Nehemiah HK, Kannan A,  Vijaya K,  Jane YN,  Merin JB. Employing clinical data sets for intelligent temporal rule mining and decision making, a comparative study.  ICGST-BIME Int J Bioinform Med Eng, 2007; 7, 37-45.

Ohsaki MY, Sato  H, Yokoi,  Yamaguchi T. A rule discovery support system for sequential medical data-in the case study of a chronic hepatitis dataset.  Proceedings of the  workshop on Discovery Challenge during the ECML/PKDD, 2003; Sept. 22-26, 1-12.

Ordonez C,  Santana C, Braal L.  Discovering interesting association rules in medical data. Proceedings of the ACM SIGMOD Workshop on Research Issues in Data Mining and Knowledge Discovery,  2000; May 14, 78-85.

Ordonez C,  Ezquerra  N, Santana C. Constraining and summarizing association rules in medical data. Knowledge Inf Syst, 2006;  9, 1-2.

Rao TSS, Shaji KS. Demographic aging: Implications for mental health. *Indian J Psychiatry*, 2007; 49, 78–80.

Ramaraj  ET,  Kumar KR. NVApriori: A novel approach to avoid irrelevant rules in association rule mining using n-cross validation technique. Int J Adv Soft Comput Appl, 2009; 1,132-150.

Salleb A, Turmeaux T,  Vrain C,  Nortet C.  Mining quantitative association rules in a atherosclerosis dataset. Proceeding of the 6th European Conference on Principles and Practice of  Knowledge Discovery in Databases, 2004;  Sept. 20-24, 98-103.

Schildt, H. (2003) "The Complete Reference Java 2", In Tata McGraw-Hill Edition, "An Overview of Java", 17-39.

Gábor Iván, Zoltán Szabadka, Vince Grolmusz. Being a binding site: Characterizing residue composition of binding sites on proteins. Bioinformation, 2007; 2(5), 216 - 221.

Huang-Cheng Kuo, Jung-Chang Lin, Ping-Lin Ong, Jen-Peng Huang. Discovering amino acid patterns on binding sites in protein complexes, Bioinformation,  2011;  6(1), 10 - 14.

Kumud Pant, Bhasker Pant, Shweta Negi. Association Rule Mining to Deduce the Most Frequently Occurring Amino Acid Patterns in HIV. International Journal of Computer Applications, 2012;  51(9),  77 - 85.

Shann-Ching Chen, Ivet Bahar. Mining frequent patterns in protein structures: a study of protease families. Bioinformatics, 2004; 20(1), 77 - 85.