

# Machine Learning based Intelligent Framework for Data Preprocessing

Sohail Sarwar<sup>1</sup>, Zia UI Qayyum<sup>2</sup>, and Abdul Kaleem<sup>1</sup>

<sup>1</sup>Department of Computing, Iqra University Islamabad, Pakistan

<sup>2</sup>Department of Computer Science, National University of Computing and Emerging Sciences Islamabad, Pakistan

**Abstract:** Data preprocessing having a pivotal role in data mining ensures reduction in cost by catering inconsistent, incomplete and irrelevant data through data cleansing to assist knowledge workers in making effective decisions through knowledge extraction. Prevalent techniques are not much effective for having more manual effort, increased processing time, less accuracy percentage etc with constrained data volumes. In this research, a comprehensive, semi-automatic preprocessing framework based on hybrid of two machine learning techniques namely Conditional Random Fields (CRF) and Hidden Markov Model (HMM) is devised for data cleansing. Proposed framework is envisaged to be effective and flexible enough to manipulate data set of any size. A bucket of inconsistent dataset (comprising of customer's address directory) of Pakistan Telecommunication Company (PTCL) is used to conduct different experiments for training and validation of proposed approach. Small percentage of semi cleansed data (output of preprocessing) is passed to hybrid of HMM and CRF for learning and rest of the data is used for testing the model. Experiments depict superiority of higher average accuracy of 95.50% for proposed hybrid approach compared to CRF (84.5%) and HMM (88.6%) when applied in separately.

**Keywords:** Machine learning, hidden markov model, conditional random fields, preprocessing.

Received March 14, 2015; accepted June 14, 2016

## 1. Introduction

Mammoth masses of data are piling up every day at a great pace making it a challenging task to extract useful information (i.e., data mining) aiding in intelligent decision making through machine learning techniques. Effectiveness of these techniques in making right decisions greatly relies on quality of data instilled into data repositories (i.e., data warehouses) where data needs to be atomic, correct, consistent, less redundant and structured in a meaningful way [3, 17, 21].

Presence of these features in data can be ensured through incorporation of certain data preprocessing techniques (i.e., Data Cleansing) [4, 13]. Some tools provided by Oracle, Premier International, Structured Query Language-Power (SQL-Power), Informatica [1, 2, 8, 10] are also available for preprocessing of data but they are expensive and not comprehensive enough to cater large data size and variety of cleansing situations (more suited to European standards of data expression) since these tools operate in a rule-based fashion i.e., if-then-else.

Currently there is no effective mechanism adopted by Pakistan Telecommunication Company Limited (PTCL) that can overcome above problems in customer addresses. Each customer's address in whole directory is stored in single column of database table, without following any standard and pattern. So it appears to be quite challenging task to correct millions of prior records congregated over years. Manifold approach,

with mutual comparison of probabilistic/statistical techniques, has been designed to address the above issues in order to cleanse the subject dataset of enormous customer base.

A semi-automatic AI techniques driven model is employed to pre-process the data with minimum time, human effort and improved accuracy level. This solution targets data preprocessing of a Telco's address directory named PTCL providing voice and data services across the Pakistan with approx 3200 telephone exchanges having 7 million customers [15].

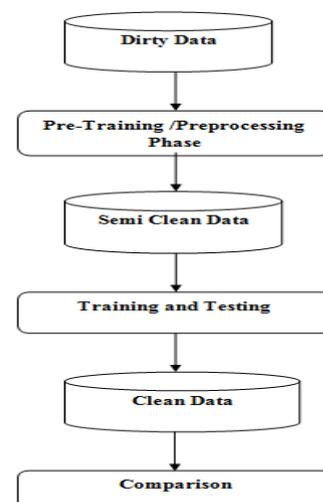


Figure 1. Overview of proposed approach.

Main steps of proposed approach are illustrated in Figure 1 starting with Pre-training/preprocessing deals with correction of data is followed by training through different Machine Learning (ML) techniques such as Hidden Markov Model (HMM) [14, 21], Conditional Random Fields (CRF) [3] and combination of both alternatively.

Comparison of stated ML techniques helps to decide which one of these is better for data cleansing, segmenting and correction especially on addresses from Asian locale as well as European style. The intended data sets to be used for testing and validation are from PTCL.

Rest of the paper is organized as follows: section 2 covers the literature review of different data cleansing techniques, such as labelling sequences, probabilistic and non-probabilistic record matching techniques, section 3 provides the details of proposed solution and, implementation of proposed solution and results are discussed in section 4 followed by conclusion and future work in section 5.

## 2. Literature Review

Before describing the method used in conducting this research, modules being investigated, their aims with corresponding pros and cons are given in the following: Labeling of a sequence [12, 17] can be considered as a set of classification tasks (independent from each other), one per each member of the sequence. In order to attain higher accuracy, the best way is to make the most appropriate (best) label for a given member depending upon the appropriate choices for adjacent member(s). It is done with the help of algorithms which are specialized in choosing the best set of labels  $S(\text{globally})$  for full sequence.

- *Inputs:*  $I = (i_1, \dots, i_n)$
- *Labels:*  $L = (l_1, \dots, l_n)$
- *Typical goal:* Given  $i$ , predict  $L$

Supervised learning [5] is used in non-probabilistic way i.e., an algorithm based on machine learning used to generate rules for matching. When an algorithm has been selected then parameters are pruned, as a result a less complex matching rule generated. After development on a sample dataset of the improved matching rule, then it is used on the original large dataset. The CART algorithm is used in [14] produces linear combination of the parameters for data classification. It is a nearest neighbor algorithm's generalized form. The names and addresses parameters used with descriptors from a small sample of dataset, in order to build the matching rules.

The distance based matching techniques is all about finding the distance between two records which depends upon the weighted sum of the records [7]. The weighted sum is calculated from the distance between

weighted sums of the records' attribute values. In [11], limitation of probability models is addressed that is in case of missing accurate estimates and counts of probability parameters due to absence of manually matched training data. It uses simply distance based technique to overcome this problem.

Using probabilistic techniques, labelling sequences can be dividing into following components [6]: that is, what are the probabilities among states for an observation with conditions. Next is the identification of a procedure which finds efficiently best output labelled sequence from all possible candidates. There are two different probabilistic methods discussed: HMM and CRF. HMM is the directed graph, where as conditional random fields is an undirected graph.

An HMM is fully represented (mathematically) by two variable,  $x$  and  $y$  and two probability distributions  $A$  and  $B$  [2, 6] as shown in Figure 2. Following is the explanation of the two variables:

- $x$  represents the number of states given in the model.
- $y$  represents number of distinct observation symbols
- In this case  $A$  represents the state transition probability distribution and  $A = a_{ij}$  (here 'a<sub>ij</sub>' denotes the probability of transition from state 'i' to state 'j'.)

$B$  represents probability distribution of the observation symbol and  $B = b_j(k)$ . Here, the symbol set 'b<sub>j</sub>(k)' is the probability of emitting the 'k th' dictionary symbol in state 'j'

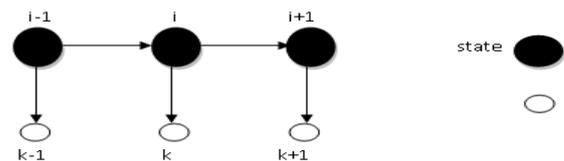


Figure 2. Hidden markov model as a directive graph.

CRFs are categorized as probabilistic models for labeling sequence based on conditions that is attribute set of an observation [4]. Technically these attributes called feature functions as illustrated in Figure 3. CRFs are undirected graphical models [17, 19]. When CRF models the conditional probability, there is a single joint probability distribution over the entire label sequence given in the observation sequence. This is in contrast to defining per-state distributions over the next states given in the current state [18].

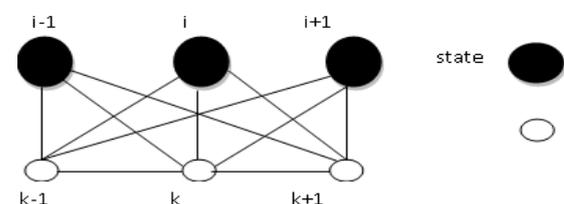


Figure 3. Conditional random fields as an un-directive graph.

As the output sequence may be created using multiple paths with each path containing some value of probability. If there are 'x' states and 'y' represents length of the sequence, then there may be  $O(yx)$  possible paths that would generate the given sequence. The exponential complexity for finding the most probable path may be a bit high and it can be reduced to  $O(kn^2)$  by the Viterbi Algorithm based on dynamic programming. Instead of summing up probabilities from different paths coming to same destination Viterbi picks up the best path and remembers it [6, 8, 14].

Sorted Neighborhood (SN) is one of the most famous approaches. The working goes like initially a blocking key 'K' is defined for each of 'n' entities. In general the blocking key is formed by concatenating the pre-fix of some attributes. In next step, sorting is performed on entities using this blocking key. Then a fixed size frame 'f' is applied over the records (already sorted) and in each step, all entities falling within the frame 'f' are compared. However, the range of distance would be f-1. [12, 16].

### 3. Proposed Architecture

The proposed solution as shown in Figure 4, has two phases i.e., 1<sup>st</sup> phase is pre-training cleansing and Tagging and 2<sup>nd</sup> phase is of Training and Testing

#### 3.1. Pre-Training Cleansing and Tagging (Phase-1)

This phase removes un-necessary words and symbols such as "no.", "#", "-", ",", for not providing meaningful information but making the training cumbersome by adding noise. Moreover, anomalies of multiple abbreviations of same word like "ISB", "ISL", "IBD" etc for "Islamabad" is addressed in addition to fixing incorrect spellings like "Islamabd" for "Islamabad" or "UMS" for "LUMS" followed by "tagging" of words.

So we end up with semi-cleansed database having no unwanted words/delimiters and most of the words with correct spellings.

These addresses are suitable input for training and testing phase carried out through AI techniques. Each process of this phase is briefly discussed in the following:

##### 3.1.1. Distinct Words Dictionary

This process tokenizes all addresses on the basis of characters, digits, delimiters and strings and creates a Dictionary. Dictionary contains all the distinct tokens [i.e. words, numbers, delimiters, characters etc] with their frequencies. Dictionary is used to get the names of Area, streets, companies, Roads and buildings etc., There is a large number of distinct tokens (words) extracted out of massive large data set. Linear search, especially in case of alphas, slowed overall processing

and performance of the system in replacement, training and testing phase.

Therefore, a 3-dimensional data structure has been employed for search optimization where two dimensions of the structure have fix length that is dictionary [1], each index representing one alphabet; the third dimension grows dynamically which contains singly linked list of distinct words. A typical representation of 3-D data dictionary is given in Figure 5.

##### 3.1.2. Data Cleansing

This component of the system is used to assign a standard word to different variations of that word and eliminates the useless words from the records. Basic use of sorted neighborhood method is to find different variation of single word by creating keys. For identification of the different variation of a single value (word), Sorted Neighborhood is one of the best approaches. The brief description of Sorted Neighborhood used into model is:

- *Keys Creation:* A key is generated for every distinct word. There are two steps for key generation. In first step all vowels are eliminated and in second step repetition of consonants is eliminated.
- *Sort data:* On the basis of key, records are sorted and grouped.
- *Merge:* Here, all words having the same key are observed and if there was any incorrect word(s), then a correct alternate value (word) is assigned for replacement.

##### 3.1.3. Replacement

While using sorted neighbourhood technique different tokens are shown to user on the basis of same key, which provides ease to correct them manually by identifying incorrect spellings in pre-training process. Once all incorrect distinct tokens are checked and correct alternate is assigned then words into database are automatically replaced by replacement process.

All the non-standard variations of a word will be replaced with standard assigned word into database by the end of this process as shown in Tables 1 and 2 respectively.

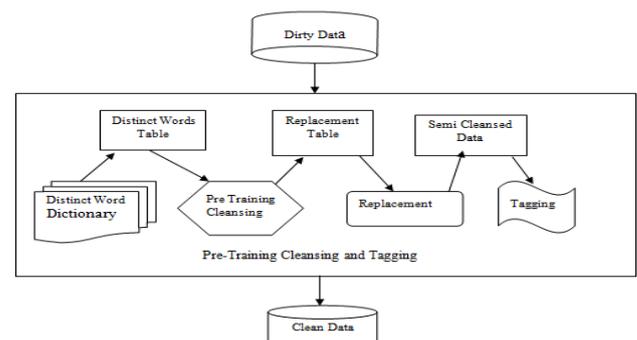


Figure 4. Proposed framework based on hybrid of ML techniques.

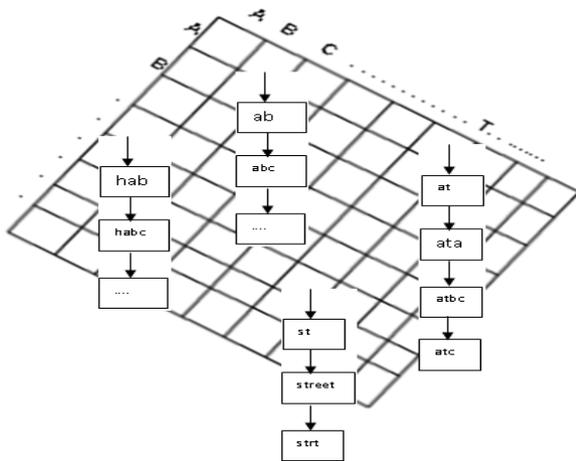


Figure 5. Representation of data dictionary.

Table 1. Key generation by ignoring vowels and repetition of consonants.

Word	Key
Colony	Clny
Rawlpindi	Rwlpnd
Muhmmad	Mhmd
Mohmmad	Mhmd
Clony	Clny
Clny	Clny
Mohmad	Mhmd

Table 2. Replacement assigned to non-standard variations of a word.

Word	Key	Replacement
Rawlpind	Rwlpnd	Rawlpindi
Rawlpindi	Rwlpnd	Rawlpindi
Colony	Clny	-
Clony	Clny	Colony
Clny	Clny	Colony
Muhmmad	Mhmd	Mohammad
Mohmmad	Mhmd	Mohammad
Mohmad	Mhmd	Mohammad

### 3.1.4. Tagging

Tags are assigned to all the distinct replaced words into the dictionary according to their nature. It creates look-up tables used in training to identify tags for the tokens in the records.

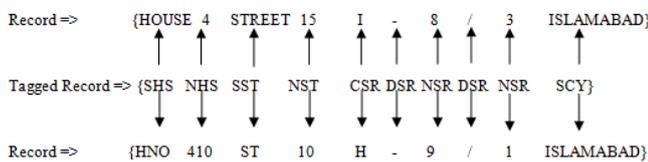


Figure 6. Tags assigned to records.

It is the most critical process. The accuracy of result is directly related to tagging; if a word is tagged wrongly then it would affect training of the machine learning techniques which would produce inaccurate result as illustrated in Figure 6.

### 3.2. Machine Learning (Phase-2)

This phase consists of training and subsequent testing of machine learning models (HMM, CRF, Hybrid-A (employing HMM followed by CRF) and Hybrid-B

(employing CRF followed by HMM)) of both) on data produced in phase-I that is asserted as training set and testing set respectively. Different sets of data buckets are used in training and testing phases. However, data for all three models remains the same in respective training and testing phases. Effectiveness of these techniques is evaluated based on degree of their elementization i.e. address is segmented into its atomic unit such that values are placed in their appropriate field (column) correctly.

## 4. Implementation and Evaluation

For implementation of HMM, Microsoft (MS) Visual studio 2010 is used as Integrated Development Environment (IDE) with MS Access as backend database in MS Windows 7 environment. Training of HMM is followed by testing where most of the steps are performed automatically or with less manual effort.

CRF code is taken from the website www.crfsharp.codeplex.com. Only one modification is made into code that is test/ output data is also stored into MS access database besides the notepad file for the comparison among HMM, CRF and their hybrid as shown in Figure 7.

At first, unclean, un-segmented and incorrect 5000 addresses were taken for training and testing of the model. All pre-processing steps were performed on data set that was transformed and semi-cleansed. It was observed that all the addresses were combination of all or some of the 5 states (i.e., HOUSE, STREET, SECTOR, ROAD and CITY). 10%, 20% and 40% addresses were selected randomly from dataset for training purpose, and rest of addresses were dedicated for testing through HMM, CRF and their hybrid. Accuracy in terms of corrected segmentation from HMM appeared to be 84.42%, 83.55% and 86.75%; accuracy for CRF accuracy was 87.10%, 89.13% and 88.24%; accuracy for Hybrid-A (HMM followed by CRF) accuracy was 87.15%, 86.34%, 88.40%; accuracy for Hybrid-B (CRF followed by HMM) was 94.50%,96.35%,95.45% given the training percentages as stated above respectively. This shows that increase in training percentage does not make any immense difference in accuracy but it will increases the training effort.

After finding facts from the dataset above, same pre-processing activities were performed with 5000 addresses. All the models (HMM, CRF, Hybrid-A and Hybrid-B) were trained with 10% of data while adhering to tenfold cross validation which gave the average accuracy of 84.52% and 88.60%, 89.22%, 95.50% respectively. Performance of model Hybrid-A and that of model Hybrid-B are comparable but when comes the point of learning cost for converging to more accurate solution, Hybrid-B appears to be a better choice as shown in Figure 8.

It is also observed that increase in size of training would not improve the results. HMM, CRF and Hybrids are fast learning models but they are based on probability so after attaining certain level of accuracy further increasing the size of training set would not affect degree of accuracy.

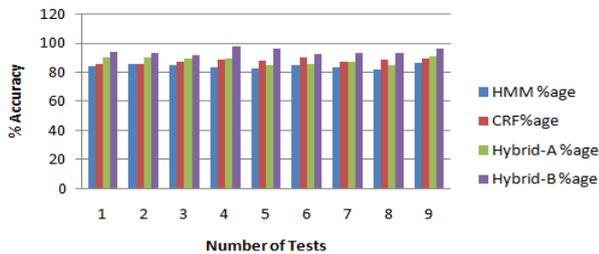


Figure 7. Performance comparison of CRF, HMM, hybrid-a and hybrid-b Tests.

Moreover, concise number of states in a dataset gives more accurate result compared to dataset having more states, though both the datasets have same number of records.

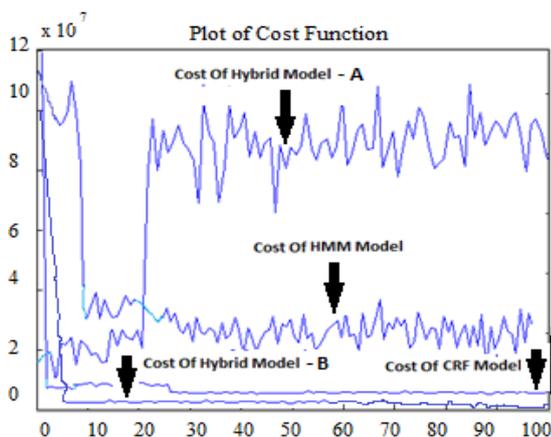


Figure 8. Cost Function of CRF, HMM, Hybrid-A and Hybrid-b Test.

## 5. Conclusions

The proposed model of data cleansing is capable of cleansing any type of large set of data especially addresses. A semi-automatic hybrid mechanism is developed that makes use of two probabilistic machine learning techniques i.e., CRF and HMM. Proposed technique is capable of generating high accuracy with minimum human effort where hybrid of CRF and HMM is found to be more accurate than peer techniques. Another major advantage that can be exploited by using proposed approach is its effectiveness for Asian-style addresses as well as for European-style addresses.

The testing of proposed model can be further improved in various ways such as correction of reference data (addresses) for comparison is made manually and for large dataset, more time will be consumed. There is a need of automatic solution which

trains the machines and lessens the training effort. This may be achieved by training of machines by using distinct observation: (based on non-repetitive tags) and / or by setting the weight of the tags.

## References

- [1] Ahmed M. and Zaman M., "Data Quality Tools for Data Warehousing: Enterprise Case Study," *IOSR Journal of Engineering*, vol. 3, no. 1, pp. 75-76, 2013.
- [2] "Address Cleansing," [www.sqlpower.ca/page/dqguru](http://www.sqlpower.ca/page/dqguru), Last Visited, 2014.
- [3] Banerjee S., Kulia A., Roy A., Naskar S., Rosso P., and Bandyo S., "A Hybrid Approach for Transliterated Word-Level Language Identification: CRF with Post-Processing Heuristics," in *Proceedings of the Forum for Information Retrieval Evaluation*, Bangalore, pp. 54-59, 2015.
- [4] Borkar V., Deshmukh K., and Sarawagi S., "Automatically Extracting Structure from Free Text Addresses," *IEEE Data Engineering Bulletin*, vol. 23, no. 4, pp. 27-32, 2000.
- [5] Breiman L., Jerome F., Richard A., and Charles J., *Classification and Regression Trees*, Chapman and Hall, 1984.
- [6] Canisius S., Bosch A., and Daelemans W., "Discrete Versus Probabilistic Sequence Classifiers for Domain-specific Entity Chunking," in *Proceedings of the 18<sup>th</sup> Belgium-Netherlands Conference on Artificial Intelligence*, Namur, pp. 175-186, 2006.
- [7] Christen P., *Data Matching: Concepts and Techniques for Record Linkage, Entity Resolution, and Duplicate Detection*, Springer Publishing Company, 2012.
- [8] Christen P., "Automatic Record Linkage using Seeded Nearest Neighbor and Support Vector Machine Classification," in *Proceedings of the 14<sup>th</sup> ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, Las Vegas, pp. 151-159, 2008.
- [9] Chu H., He Y., Chakrabarti K., and Ganjam K., "TEGRA: Table Extraction by Global Record Alignment," *ACM SIGMOD*, in *Proceedings of International Conference on Management of Data*, Victoria, pp. 1713-1728, 2015.
- [10] "Data Cleansing Tools," [www.premier-international.com/solutions\\_name\\_and\\_address\\_data\\_cleansing.aspx](http://www.premier-international.com/solutions_name_and_address_data_cleansing.aspx), Last Visited, 2013.
- [11] Galhardas H., Florescu D., Shasha D., Simon E., and Saita C., "Declarative Data Cleaning: Language, Model, and Algorithms," in *Proceedings of 27<sup>th</sup> International Conference on Very Large Databases*, Roma, pp. 371-380, 2001.

- [12] Kolb L., Thor A., and Rahm E., "Parallel Sorted Neighborhood Blocking with MapReduce," *Journal of Computer Science-Research and Development*, vol. 27, no. 3, pp. 45-63, 2012.
- [13] Kulkarni P. and Bakal L., "Article: Hybrid Approaches for Data Cleaning in Data Warehouse," *International Journal of Computer Applications*, vol. 88, no. 18, pp. 7-10, 2014.
- [14] Noor S. and Bashir S., "Evaluating Bias in Retrieval Systems for Recall Oriented Documents Retrieval," *The International Arab International Journal of Information Technology*, vol. 12, no. 1, pp. 53-59, 2015.
- [15] "PTCL," [www.ptcl.com.pk/pd-content.php?pd\\_id=41](http://www.ptcl.com.pk/pd-content.php?pd_id=41), Last Visited, 2014.
- [16] Raham E. and Hai H., "Data Cleaning: Problems and Current Approaches," *Bulletin of the Technical Committee on Data Engineering*, vol. 23, no. 4, pp. 3-13, 2000.
- [17] Shuxin Z., Zhonghong X., and Yuehong C., "Information Extraction from Research Papers based on Conditional Random Field Model," *International Journal Electrical Engineering and Computer Science*, vol. 11, no. 3, pp. 1213-1220, 2013.
- [18] Tan Y., Yao T., Chen Q., and Zhu J., "Applying Conditional Random Fields to Chinese Shallow Parsing," in *Proceedings of International Conference on Intelligent Text Processing and Computational Linguistics*, Mexico City, pp. 167-176, 2005.
- [19] "Tools for Data Warehousing," [www.ctg.albany.edu](http://www.ctg.albany.edu), Data Quality-A Small Sample Survey, Center for Technology in Government, Last Visited, 2014.
- [20] Torra V., "Information Fusion-Methods and Aggregation Operators," in *Proceedings of Data Mining and Knowledge Discovery Handbook*, Boston, pp. 999-1008, 2010.
- [21] Wallach H., "Efficient Training of Conditional Random Fields," in *Proceedings of 6<sup>th</sup> Annual Computational Linguistics U.K. Research Colloquium*, 2002.



**Sohail Sarwar** is a PhD student at Department of Computing, Iqra University Islamabad in Pakistan. His domain of interest is application of machine learning techniques such as data mining, software testing and e-learning.



**Zia Ul Qayyum** is a Professor in Department of Computing, Iqra University. His area of research covers application of AI techniques in image processing and classification, natural language processing, data preprocessing, semantics and recommender systems.



**Abdul Kaleem** is a Master's student in Iqra University. He has research interest in applying knowledge engineering in big data. He has keen interest in building web based software applications pattern mining and business intelligence.