

## Jackknife Estimation of Incomplete Data for Data Marts for Customer Relationship Management

Venkata Naresh Mandhala<sup>1</sup>, N Lakshmipathi Anantha<sup>2</sup>,  
Vijay Krishna Dhulipalla<sup>3</sup> and Hye-Jin Kim<sup>4\*</sup>

<sup>1</sup>*Department of Computer Science and Engineering,  
KL University, Vaddeswaram, AP, 522502, India*

<sup>2</sup>*Department of Information Technology, VFSTR University,  
Vadlamudi-522213, Guntur, India*

<sup>3</sup>*Department of Management Studies, VFSTR University,  
Vadlamudi-522213, Guntur, India*

<sup>4</sup>*Sungshin W. University, 2,Bomun-ro 34da-gil, 0Seongbuk-gu, Seoul*

<sup>1</sup>*mvnaresh.mca@gmail.com, <sup>2</sup>anlakshmipathi@gmail.com,*

<sup>3</sup>*taketovijay@gmail.com, <sup>4\*</sup>hyejinaa@daum.net*

### Abstract

*The commitment of measurements to information mining can be followed back to the work by Bayes in 1763. The business organizations gather information and offer it to the Data Marts. The individuals who run little and medium association needs to set up information warehousing to touch base, best case scenario arrangement. Such datasets contain part of missing qualities, at some point the missing qualities range from 10% to 33%. A portion of the information might be fundamental; to recall such information is a troublesome undertaking and this kind of datasets won't yield better arrangement, to take care of this issue the Expectation Maximization (EM) calculation gauges missing qualities. Utilizing EM Algorithm the outcomes are supplanted in the missing positions of the specific information which serves to exact conclusion. In this paper, point estimators were connected, among which EM calculation gives best gauge. It is watched that the more straightforward models by and large yield the best results.*

**Keywords:** Data Marts, Expectation Maximization, Machine Learning, Data Mining

### 1. Introduction

There are a wide range of strategies used to perform information mining undertakings. These systems require particular sorts of information structures, as well as suggest certain sorts of algorithmic methodologies. An examination of the effect measurements has had on the advancement of information mining strategies. Numerous information mining strategies discover their introduction to the world in the machine learning field. In this paper we quickly present a portion of the normal information mining methods [1].

Information Marts purchase information from businessperson and offers information to an amateur individual or organization, who needs to build up his organization. At the point when the client enters the shop, the shop proprietor gathers some individual data about the client for anticipating the Market interest for specific thing in a specific spot at specific time [2]. It is entirely difficult to gather all the accurate data about every one of the clients, for example, client name, area, age, sex, conjugal status, umber of kids, has auto, has house, salary, has telephone, occupation and so forth., To gather each data is unrealistic in light of the accompanying reasons:

---

\* Corresponding Author

1. Customer may hesitant to give all data, because of mystery.
2. Shop proprietor may miss to gather each information amid crest hours of business.
3. Shop proprietor may enter a few fields, for example, age, wage and so on around by seeing the client.

The Data Marts needs to supply data to the penniless individuals at any rate with evaluated or approximated information. Subsequently we apply guide estimators toward take care of this kind of issue [3]. Astounding investigations of the relationship between machine learning and information mining were assessed.

### **1.1. Parametric Models**

It depicts the relationship amongst info and yield using mathematical conditions where a few parameters are not indicated. These unspecified parameters are controlled by giving information illustrations. Despite the fact that parametric modeling is a decent hypothetical theme and can infrequently be utilized, regularly it is either excessively oversimplified or requires more information about the information required than is accessible. In this way, for certifiable issues, these parametric models may not be valuable [6-7].

### **1.2. Non-Parametric Procedures**

These are more suitable for information mining applications. A non-parametric model is one that is information driven. No unequivocal conditions are utilized to determine the model. This implies the displaying procedure adjusts to the current information. Not at all like parametric demonstrating, where a particular model is accepted early, the non-parametric systems make a model in view of the info. While the parametric strategies require more learning about the information before the displaying procedure, the non-parametric method requires a lot of information as contribution to the demonstrating procedure itself. The demonstrating procedure then makes the model by moving through the information. Late non-parametric strategies have utilized machine learning methods to have the capacity to learn powerfully as information are added to the info. Hence, the more information, the better the model made. Likewise, this dynamic learning process permits the model to be made persistently as the constant information is food as info. These components make non-parametric systems especially appropriate to database applications with a lot of powerfully changing normal for information. Non-parametric procedures incorporate neural systems, choice trees, and hereditary calculations [8].

### **1.3. Organization of this Paper**

Area 2 talks about factual viewpoint on information mining. Area 2.1 arrangements, with Point-estimation. In Section 2.2, we talk about, means squared mistake strategy utilizing certainty interim. Area 2.3, we consider the Root Means Squared technique. In Section 2.4, we manage interim evaluation with one or all the more missing qualities. In Section 2.5, we talk about the Maximum Likelihood Estimator Method for assessing missing quality. In Section 2.5.1, we propose EM calculation for evaluating deficient information. Area 3 gives future extent of the paper. Area 4, determines conclusion.

## **2. Statistical Data Mining**

There have been numerous measurable ideas that are the premise for information mining procedures. We quickly survey some of these ideas.

## 2.1. Point Estimation

Guide estimation alludes toward the way toward evaluating a populace parameter,  $\theta$ , by an assessment of the parameter,  $\hat{\theta}$ . This should be possible to gauge mean, difference, standard deviation, or some other factual parameter. Frequently the estimation of the parameter for an all-inclusive community might be ascertained from the parameter esteem for a populace test. An estimator strategy may likewise be utilized to anticipate the benefit of missing information. The predisposition of an estimator is the distinction between the normal estimation of the estimator and the genuine quality:

$$Bias = E(\hat{\theta} - \theta)^2 \quad (1)$$

A fair-minded estimator is one, whose predisposition is 0. While point estimators for little information sets may really be unprejudiced, for bigger database applications we would expect that most estimators are one-sided.

## 2.2. Mean Squared Error (MSE) Method

One measure of the adequacy of an evaluation is the mean squared mistake (MSE), which is characterized as the normal estimation of the squared distinction between the appraisal and the real esteem:

$$MSE(\theta) = E(\hat{\theta} - \theta)^2 \quad (2)$$

The squared blunder is regularly inspected for a particular expectation to gauge precision as opposed to take a gander at the normal distinction. For instance, if the genuine worth for a trait was 10 and the forecast was 5, the squared mistake would be  $(5 - 10)^2 = 25$ . The squaring is performed to guarantee that the measure is constantly positive and to give a higher weighting to the appraisals that are terribly mistaken. As we will see, the MSE is generally utilized as a part of assessing the adequacy of information mining forecast strategies. It is likewise essential in machine learning. Now and again, rather than foreseeing a straightforward point gauge for a parameter, one may decide a scope of qualities inside which the genuine parameter worth ought to fall. This reach is known as a certainty interim.

## 2.3. Root Mean Square (RMS)

RMS may likewise be utilized to gauge blunder or as another measurement to depict an appropriation. Ascertaining the mean does not demonstrate the greatness of the qualities [9]. The RMS can be utilized for this reason. Given an arrangement of  $n$  qualities  $X = \{x_1, x_2, \dots, x_n\}$ , the RMS is defined by

$$RMS = \sqrt{\left\{ \frac{\sum_{j=1}^n x_j^2}{n} \right\}} \quad (3)$$

The next utilization is to gauge the extent of the blunder. The root mean square blunder (RMSE) is found by taking the square foundation of the MSE [9].

## 2.4. Interval Estimate

We talk about an Interval Estimator,  $i^{\text{th}}$  this approach, the evaluation of a parameter  $\mu_i$  is gotten by precluding one worth or more values from the arrangement of watched

qualities is called as Jackknife assessment. Assume that there is an arrangement of n qualities  $X = \{x_1, x_2, \dots, x_n\}$ . An estimate for the mean would be,

$$\bar{\mu}_i = \frac{\sum_{j=1}^{i-1} x_j - \sum_{j=i+1}^n x_j}{n-1} \quad (4)$$

Here, the subscript (i) shows that this assessment is acquired by precluding the i<sup>th</sup> esteem. We can exclude one quality or more values[10]. We can likewise have numerous interims, at the end of the day, we can likewise overlook numerous qualities. Given an arrangement of interim evaluations,  $\theta_{(i)}$ , these can in turn be used to obtain an overall estimate.

$$\bar{\theta}_{(i)} = \sum_{j=1}^n \hat{\theta}_j \quad (5)$$

### Example 1

Accept that a coin is heaved observable all around for five times with the results, {1, 1, 0, 1, 1}, where, 1 demonstrates a head and 0 exhibits a tail. If we expect that the coin heave takes after the Bernoulli movement, we understand that

$$f(x_i/p) = p^{x_i} (1-p)^{(1-x_i)} \quad (6)$$

Assuming a perfect coin when the probability of 1 and 0 are both 1/2, the likelihood is then

$$L(p | 1,1,1,1,0) = \prod_{i=1}^n 0.05 = 0.03 \quad (7)$$

However, if the coin is not perfect but has a bias toward heads such that the probability of getting a head is 0.8, the likelihood is

$$L(p | 1,1,1,1,0) = 0.8 \times 0.8 \times 0.8 \times 0.8 \times 0.2 = 0.08 \quad (8)$$

Here it is more likely that the coin is biased toward getting a head than that it is not biased. The general formula for likelihood is

$$L(p | x_1, x_2, \dots, x_5) = \prod_{i=1}^5 p^{x_i} (1-p)^{(1-x_i)} = p^{\sum_{i=1}^5 x_i} (1-p)^{(5-\sum_{i=1}^5 x_i)} \quad (9)$$

By taking log both the sides we get,

$$\log(L(p)) = \sum_{i=1}^5 x_i \log(p) + \left(5 - \sum_{i=1}^5 x_i\right) \log(1-p) \quad (10)$$

and then we take the derivative with respect to p

$$\frac{\partial l(p)}{\partial p} = \left\{ \frac{\sum_{i=1}^5 x_i}{p} \right\} - \left\{ \frac{\left(5 - \sum_{i=1}^5 x_i\right)}{(1-p)} \right\} \quad (11)$$

Setting equal to zero we finally obtain

$$p = \left\{ \frac{\sum_{i=1}^5 x_i}{5} \right\} \quad (12)$$

For this example, the estimate for p is then  $p = 4/5 = 0.8$ . Thus, 0.8 is the value for p that maximizes the likelihood that the given sequence of heads and tails would occur.

## 2.5. Maximum Likelihood Estimate (MLE)

Another procedure for point estimation is known as the MLE. Probability can be characterized as a quality corresponding to the genuine likelihood that with a particular conveyance the given example exists. Thus, the example gives us an evaluation for a parameter from the dissemination. The higher the probability esteem, the more probable the hidden conveyance will deliver the outcomes watched. Given a sample set of values  $X = \{x_1, x_2, \dots, x_n\}$  from a known  $f(x_i/\theta)$  distribution function, the MLE can estimate parameters for the population from which the sample is drawn. Estimating attributes can be done by the analysis and extensions of RELIEF method [5]. The methodology acquires parameter appraisals that augment the likelihood that the example information happen for the particular model. It takes a gander at the joint likelihood for watching the example information by duplicating the individual probabilities [11]. The probability capacity, L, is subsequently characterized as

$$L(\theta | x_1, \dots, x_n) = \prod_{i=1}^n f(x_i/\theta) \quad (13)$$

The value of  $\theta$  that maximizes L is the estimate chosen. This can be found by taking the derivative, after finding the log of each side to simplify the formula with respect to  $\theta$ . Example 1, illustrates the use of MLE.

**2.5.1. The Expectation-Maximization (EM) algorithm:** The EM calculation is a methodology that takes care of the estimation issue with inadequate information. The EM calculation finds a MLE for a parameter, for example, a mean, utilizing a two-stage process: estimation and augmentation. The fundamental Hierarchical mixtures of experts and the EM algorithm is appeared in Algorithm 1 [4]. An underlying arrangement of appraisals for the parameters is acquired. Given these appraisals and the preparation information as information, the calculation then computes a quality for the missing information. For instance, it may utilize the evaluated intend to anticipate a missing worth. These information (with the new esteem included) are then used to decide an evaluation for parameters which are drawn from the probability[12]. These strides are connected iteratively until successive parameter gauges focalize. Any methodology can be utilized to locate the underlying parameter gauges. In Algorithm 1 it is accepted that the info database has genuine watched values  $X_{obs} = \{x_1, x_2, \dots, x_k\}$  as well as values that are missing  $X_{miss} = \{x_{k+1}, x_n\}$ . We assume that the entire database is actually  $X = \{X_{obs} \cup X_{miss}\}$ . The parameters to be estimated are  $\theta = \{\theta_1, \theta_2, \dots, \theta_p\}$ . The likelihood function is defined by

$$L(\theta | X) = \prod_{i=1}^n f(x_i | \theta) \quad (14)$$

We are looking for which maximizes L. The MLE of are the estimates that satisfy

$$\partial \ln L(\theta | X) / \partial \theta_i = 0 \quad (15)$$

The desire part of the calculation evaluates the missing qualities utilizing the present assessments of  $\theta$ . This should at first be possible by finding a weighted normal of the watched information. The amplification step then finds the new gauges for  $\theta$  parameters that amplify the probability by utilizing those assessments of the missing information. An illustrative case of the EM calculation is appeared in Example 2.

**Algorithm 1**

*Input:*  $\theta = \{\theta_1, \theta_2, \dots, \theta_n\}$   
 //Parameters to be estimated  
 $X_{obs} = \{x_1, x_2, \dots, x_k\}$   
 //Input database values observed  
 $X_{miss} = \{x_{k+1}, x_n\}$   
 //Input database values missing  
  
*Output:*  $\theta$  // Estimate for  $\theta$

**EM algorithm:**

*Initialize*  $i = 0$ ;  
 Obtain initial parameter MLE estimate,  $\theta_i$  ;  
 repeat  
     Estimate missing data,  $x_{miss}^2$  ;  
     Increment  $i++$   
     obtain next parameter estimate,  $\theta^2$  to maximize likelihood;  
 until estimate converges;

**Example 2**

We wish to find the mean,  $\mu$ , for data that follow the normal distribution where the known data are {1,5,10, 4} with two data items missing. Here  $n = 6$  and  $k = 4$ . Suppose that we initially guess  $\mu_0 = 3$ . We then use this value for the two missing values. Using this, we obtain the MLE estimate for the mean as

$$\mu^1 = \left\{ \frac{\sum_{i=1}^k x_i}{n} \right\} + \left\{ \frac{\sum_{i=k+1}^n x_i}{n} \right\} \tag{16}$$

We now repeat using this as the new value for the missing items, and then estimate the mean as

$$\mu^2 = \left\{ \frac{\sum_{i=2}^k x_i}{n} \right\} + \left\{ \frac{\sum_{i=k+1}^n x_i}{n} \right\} \tag{17}$$

Repeating we obtain

$$\mu^3 = \left\{ \frac{\sum_{i=3}^k x_i}{n} \right\} + \left\{ \frac{\sum_{i=k+1}^n x_i}{n} \right\} \tag{18}$$

and then

$$\mu^4 = \left\{ \frac{\sum_{i=4}^k x_i}{n} \right\} + \left\{ \frac{\sum_{i=k+1}^n x_i}{n} \right\} \tag{19}$$

The expected values are tabulated in Table (1).

### 3. Future Scope of the Paper

These Point estimators can be used to estimate on large sample of data. It can be further applied on Distributed Data Marts. While collecting data, we can apply on Real-Time.

**Table 1. Expected Values**

μ values	Expected value calculation	Expected values
1	3.33 + (3+3)/6	4.33
2	3.33+(4.33+4.33)/6	4.77
3	3.33+(4.77+4.77)/6	4.92
4	3.33+(4.92+4.92)/6	4.97

We decide to stop here because the last two estimates are only 0.05 apart. Thus, our estimate is  $\mu = 4.97$ .

### 4. Conclusion

When we purchase information from information bazaars it contains inadequate or missing information. These approximated information's are evaluated effectively in our paper utilizing different point estimators, for example, MSE, RMSE, Interval appraisal, MLE and EM Algorithm. Among which EM calculation gives best gauge for inadequate information and the outcomes are organized. One of the fundamental rules in evaluating is that less complex models by and large yield the best results.

### References

- [1] A. P. Dempster, N. M. Laird and D. B. Rubin, "Maximum likelihood from incomplete data via the EM algorithm", Journal of the royal statistical society. Series B (methodological), (1977), pp. 1-38.
- [2] J. Elder and D. Pregibon, "A statistical perspective on KDD", Advances in knowledge discovery and data mining, (1996), pp. 83-116.
- [3] I. H. Witten and E. Frank, "Data Mining: Practical machine learning tools and techniques", Morgan Kaufmann, (2005).
- [4] M. I. Jordan and R. A. Jacobs., "Hierarchical mixtures of experts and the EM algorithm", Neural computation, vol. 6, no. 2, (1994), pp. 181-214.
- [5] I. Kononenko, "Estimating attributes: analysis and extensions of RELIEF", European conference on machine learning, Springer Berlin Heidelberg, (1994), pp. 171-182.
- [6] L. R. Dysert, "Developing a parametric model for estimating process control costs", AACE International Transactions, T11, (1999).
- [7] D. Eck, B. Brundick, T. Fettig, J. Dechoretz and J. Ugljesa, "Parametric estimating handbook", The International Society of Parametric Analysis (ISPA), (2009).
- [8] Z. Zhao, "Parametric and nonparametric models and methods in financial econometrics", Statistics Surveys, vol. 2, (2008), pp. 1-42.

- [9] K. Kelley and K. Lai, "Accuracy in parameter estimation for the root mean square error of approximation: Sample size planning for narrow confidence intervals", *Multivariate Behavioral Research*, vol. 46, no. 1, (2011), pp. 1-32.
- [10] D. Gilliland and V. Melfi, "A note on confidence interval estimation and margin of error", *Journal of Statistics Education*, vol. 18, no. 1, (2010), pp. 1-8.
- [11] F. W. Scholz, "Maximum likelihood estimation", *Encyclopedia of statistical sciences*, (1985).
- [12] M. R. Gupta and Y. Chen, "Theory and use of the EM algorithm", *Now Publishers Inc*, (2011).