



# Classification approach based customer prediction analysis for policy preferences of life insurance customers.

S.Balaji<sup>1</sup>, Dr.S.K.Srivatsa<sup>2</sup>

<sup>1</sup> Research Scholar, Vels University Email: srisaibalaji@rediffmail.com

<sup>2</sup> Senior Professor, St. Joseph Eng. College Chennai-600100

## Abstract:-

Prediction analysis is a definite need of any business sector for retaining and attracting the most valuable customers. It is considered as a major challenge facing companies in this information age. Data mining enables companies, in the context of defined business objectives, discover new knowledge, to explore, visualize and understand their data, and to identify patterns, relationships and dependencies that impact on business outcomes. The main focus of this paper concerned with Naive Bayesian classification algorithm for customer classification and prediction on Life Insurance dataset.

Keywords: Data Mining, Naïve Bayesian, customer relationship Management.

## 1.Introduction

Insurance industry in India aims to protect the interest of and secure fair treatment to policyholders and to bring about speedy and orderly growth of the insurance industry (including annuity and superannuation payments), for the benefit of the common man, and to provide long term funds for accelerating growth of the economy. In life insurance business, India ranked 9th among the 156 countries, for which data are published by Swiss Re. During 2010-11, the estimated life insurance premium in India grew by 4.2 per cent (inflation adjusted). However, during the same period, the global life insurance premium expanded by 3.2

per cent. The share of Indian life insurance sector in global market was 2.69 per cent during 2010, as against 2.45 per cent in 2009. Data Mining can help life insurance companies to make crucial business decisions and turn the new found knowledge into actionable results in business practices. Insurance firms can increase profitability by identifying the most lucrative customer segments and then prioritize marketing campaigns accordingly. Problems with profitability can occur if life insurance companies do not offer the right policy or the right rate to the right customer segment at the right time. The success of the life insurance profession depends, above all, upon the knowledge and integrity of the people who advise customers – and are their first, and most important point of contact. The term Insurance

agents coined as Insurance advisor .An insurance advisor has the unique opportunity to earn the gratitude of people in addition to highly rewarding one, both in terms of money and in terms of prestige and satisfaction. At the IRDA, the regulator's goal is to see that life insurers are increasingly able to attract, motivate and retain outstanding people, committed to adopting a 'needs-based' approach to financial advice. With DM operations insurance firms can now utilize all of their available information to better develop new products and marketing campaigns. Life insurance industry recorded a premium income of 2,91,605 crore during 2010-11 as against 2,65,447 crore in the previous financial year, registering a growth of 9.85 per cent. Data mining techniques can help insurance companies to guide the potential insurance advisors and customers and to map exact policy for proposal. Data mining enables companies, in the context of defined business objectives, discover new knowledge, to explore, visualise and understand their data, and to identify patterns, relationships and dependencies that impact on business outcomes.

## 2.Literature survey

Data Mining is a crucial step in the Knowledge Discovery in Database (KDD) process that consists of applying data analysis and knowledge discovery algorithms to produce useful patterns (or rules) over the datasets. Using data mining technology can filtrate and classify customer resources of insurance, divide credit customers into several grades, to predict the customer risk, thus investigating customer material of the low forecasted degrees of comparison can avoid deceiving policy effectively, and avoid service risk. Marisa .S.Viveros[1996] addresses the effectiveness of two data mining techniques in analyzing and retrieving unknown behavior patterns from gigabytes of data collected in the health insurance industry. Mittal & Kamakura (2001) find the link between customer satisfaction and retention to be moderated by customer characteristics. kanwal garg(2008) find decision tree method for identifying customer behaviour of investment in life insurance sector. Patrick A Rivers( 2010) examined some of the benefits and challenges of using data mining processes within the health-care arena.

## 3.Methodology:-

Data mining methodology can often improve existing actuarial models by finding additional important variables, by identifying interactions, and by detecting nonlinear relationships. Insurance Market is purely based on customer penetration. Navie Bayes is the basis for many machine learning and data mining methods. In Bayesian classification is a classification method is applicable for huge dataset. Naive Bayesian classifier works with hypothesis H such as that the data tuple X belongs to a specified class C. The determination of  $P(H/X)$  that the hypothesis H holds given the evidence or observed data tuple X.  $P(H/X)$  is the posterior probability of H conditioned on X. Bayes' theorem is useful in that it provides a way of calculating the posterior probability,  $P(H/X)$ , from  $P(X/H)$  and  $P(X)$ ,

Bayes theorem is  

$$P(H/X) = P(X/H)P(H)/P(X).$$

### 3.1 Naive Bayesian Classification Algorithm

The naive Bayesian classifier, or simple Bayesian classifier, works as follows:

1. Let D be a training set of tuples and their associated class labels. As usual, each tuple is represented by an n-dimensional attribute vector,  $X = (x_1, x_2, \dots, x_n)$ , depicting n measurements made on the tuple from n attributes, respectively,  $A_1, A_2, \dots, A_n$ .
2. Suppose that there are m classes,  $C_1, C_2, \dots, C_m$ . Given a tuple, X, the classifier will predict that X belongs to the class having the highest posterior probability, conditioned on X. That is, the naive Bayesian classifier predicts that tuple x belongs to the class  $C_i$  if and only if

$$P(C_i|X) > P(C_j|X) \text{ for } 1 \leq j \leq m, j$$

$\neq i$

Thus we maximize  $P(C_i|X)$ . The class  $C_i$  for which  $P(C_i|X)$  is maximized is called the maximum posteriori

hypothesis. By Bayes' theorem

3. As  $P(X)$  is constant for all classes, only  $P(X|C_i) P(C_i)$  need be maximized. If the class prior probabilities

are not known, then it is commonly assumed that the classes are equally likely, that is,  $P(C_1) = P(C_2)$

$= \dots = P(C_m)$ , and we would therefore maximize  $P(X|C_i)$ . Otherwise, we maximize  $P(X|C_i)P(C_i)$ . Note that

the class prior probabilities may be estimated by  $P(C_i) = |C_i, D| / |D|$ , where  $|C_i, D|$  is the number of training

tuples of class  $C_i$  in  $D$ .

4. Given data sets with many attributes, it would be extremely computationally expensive to compute  $P(X|C_i)$ .

In order to reduce computation in evaluating  $P(X|C_i)$ , the naïve assumption of class conditional

independence is made. This presumes that the values of the attributes are conditionally independent of one

another, given the class label of the tuple (i.e., that there are no dependence relationships among the

attributes). Thus,

$$= P(x_1|C_i) \times P(x_2|C_i) \times \dots \times P(x_m|C_i).$$

We can easily estimate the probabilities  $P(x_1|C_i)$ ,  $P(x_2|C_i)$ , ...,  $P(x_m|C_i)$  from the training tuples.

Recall that

here  $x_k$  refers to the value of attribute  $A_k$  for tuple  $X$ . For each attribute, we look at whether the attribute is

categorical or continuous-valued. For instance, to compute  $P(X|C_i)$ , we consider the following:

(a) If  $A_k$  is categorical, then  $P(X_k|C_i)$  is the number of tuples of class  $C_i$  in  $D$  having the value  $x_k$  for  $A_k$ ,

divided by  $|C_i, D|$ , the number of tuples of class  $C_i$  in  $D$ .

(b) If  $A_k$  is continuous valued, then we need to do a bit more work, but the calculation is pretty straightforward. A continuous-valued attribute is typically assumed to have a Gaussian distribution with a mean

$\mu$  and standard deviation  $\sigma$ , defined

by

$$g(x, \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

So that

$$P(x_k|C_i) = g(x_k, \mu_{ci}, \sigma_{ci})$$

ci)

We need to compute  $\mu_{ci}$  and  $\sigma_{ci}$ , which are the mean and standard deviation, of the values of attribute

$A_k$  for training tuples of class  $C_i$ . We then plug these two quantities into the above equation.

5. In order to predict the class label of  $X$ ,  $P(X|C_i)P(C_i)$  is evaluated for each class  $C_i$ . The classifier

predicts that the class label of tuple  $X$  is the class  $C_i$  if and only if

$P(X|C_i)P(C_i) > P(X|C_j)P(C_j)$  for  $1 \leq j \leq m, j$

$\neq i$

In other words, the predicted class label is the class  $C_i$  for which  $P(X|C_i)P(C_i)$  is the maximum

#### 4.Experiments and its results

4.1.Data Set:-A IRDA Data set of Life Insurance Corporation Of India provided its transaction data for analyses. LIC contributed most of the business procured in this portfolio by garnering `123 crore of individual premium from 29.51 lakh policies and 138 crore of group premium

The entire data set covers the period from January 2011-December 2011.The dataset,containing 10,000 customer proposals.The dataset of customers contained customer product type and the plan preference under that category in addition to proposer's risk coverage details.classification of customers based on data for the study. It is known that too many attributes involved will very possibly result in discovered information that is difficult to interpret, or even meaningless. Therefore, by in-depth discussion with domain managers, we eliminated some of the attributes and finally came to a conclusion of 7 attributes, namely 1)Gender 2)age 3)Marital status 4)No of Kids-the customers of marital status married having possible values one or two or three5)Service category :-customer may be a minor or major .Non-minor customer may be in still service or retired 6)Product type –Unit Linked or Traditional product 7) Plan types –Savings plans,Protection plans,Pension plans and Child Plans. For the data set age values may be binned in to the following categories Unmarried in service,Unmarried not in service,Newly Married-in service-without kids,Married-inservice-withchildren,Married not in service,without children ,Married not in service with children. Their preference over the policy is transformed as A,B,C and D.In the above class Not in sevice refres to minor or jobless or retired from service.

Figure :-4.1 Sample records of insurance data set

sample attributes of the records of the dataset

Married ?	Service	Kids	Age subsection(<40)	Product Type	Policy type
NO	NO	NO	NO	B	A
NO	YES	NO	NO	B	B
NO	YES	NO	YES	A	B
YES	YES	NO	NO	A	C
YES	YES	YES	YES	A	D
YES	YES	YES	NO	B	D
YES	NO	NO	NO	B	A
YES	NO	YES	NO	A	C
YES	NO	YES	YES	B	C
NO	NO	NO	YES	A	A

Product type A refers to Unit-Linked and B refers to Traditional one.

Policy type refers to A For Savings,B Class refers to Protection plans,C refers to Pension Plans and D for Child plans.

The maximum number of levels in the tree in the tree was limited to four and the minimum number of records in a node was set to 4000,inorder to prevent the Decision Tree from becoming very complex.

The Bayesian approach having the beauty that the probability of the dependent attribute can be estimated by computing estimates of the probabilities of the independent attributes.

There are 10000 (s=10000) samples and four classes.

The frequency of classes as A=3000,B=2000,C=3000,D=2000

The prior probabilities are obtained by dividing these frequencies by the total number in the training data.

$$P(A)=0.3,P(B)=0.2,P(C)=0.3,P(D)=0.2$$

Computation of posterior probabilities for the four classes,namely that the customer with attribute values X has Policy preference of Policy type A or type B or type C or type D.

The computation of  $P(X|Ci)P(Ci)$  for each of four classes are as  $P(A)=0.3,P(B)=0.2,P(C)=0.3,P(D)=0.2$  and these values are basis for comparing the four classes.

Predicting a class label using naïve Bayesian classification,we wish to predict the class label of a tupe using naïve Bayesian classification from the sample data .The data tuples are described by the attributes Married,service,age\_sub\_section ,kids and product type.

The policy type attribute has four distinct attributes namely { A,B,C,D}.

Let C1 coresspond to policy type preference A, C2 coresspond to policy type preference B, C3 coresspond to policy type preference C, C4 coresspond to policy type preference D.

I.The tupe to classify is

$$X=\{ \text{Married=yes.service=no,kids=YES,agesubsection}(<40)=\text{NO,producttype=A} \}$$

$$P(X/A)=1000/3000*3000/3000*3000/3000*2000/3000 *1000/3000$$

$$P(X/A)=1.2$$

$$P(X/B)=0*0*0*1000/2000*1000/2000=0$$

$$P(X/C)=3000/3000*2000/3000*2000/3000*2000/3000=48/15=3.2$$

$$P(X/D)=2000/2000*0*2000/2000*1000/2000*1000/2000=0$$

$P(X/Ci)P(Ci)$  is being computed to find the class that maximizes  $Ci$ ,

$$P(X/\text{Policy type =A})=0.3*1.2=0.36$$

$$P(X/\text{Policy type=B})=0.2*0=0$$

$$P(X/\text{Policy type=C})=0.3*3.2=0.96$$

$$P(X/\text{Policy type=D})=0.2*0=0$$

Therefore the naïve Bayesian classifier predicts Policy type=C ( Pension Plans) for tuple X.

Where as X is observed as Married=yes.service=no,kids=YES,agesubsection(<40)=NO,producttype=A }

Bayes theorem assumes that all attributes are independent and that the sample is good enough to estimate probabilities

II. The tupe to classify is

$$X=\{ \text{Married=yes.service=no,kids=YES,agesubsection}(<40)=\text{NO,producttype=B} \}$$

$$P(X/A)=1000/3000*3000/3000*3000/3000*2000/3000 *2000/3000=36/15=2.4$$

$$P(X/B)=0*0*0*1000/2000*1000/2000=0.0$$

$$P(X/C)=3000/3000*2000/3000*2000/3000*2000/3000*1000/3000=24/15=1.6$$

$$P(X/D)=2000/2000*0*2000/2000*1000/2000*1000/2000=0$$

$P(X/C_i)P(C_i)$  is being computed to find the class that maximizes  $C_i$ ,

$$\begin{aligned} P(X/\text{Policy type}=A) &= 0.3 * 2.4 = 0.72 \\ P(X/\text{Policy type}=B) &= 0.2 * 0 = 0 \\ P(X/\text{Policy type}=C) &= 0.3 * 1.6 = 0.48 \\ P(X/\text{Policy type}=D) &= 0.2 * 0 = 0 \end{aligned}$$

Therefore the naïve Bayesian classifier predicts Policy type=A (Savings Plans) for tuple X, which maximized  $C_i$ . The value X is observed as  $X = \{\text{Married}=\text{yes}, \text{service}=\text{no}, \text{kids}=\text{YES}, \text{agesubsection}(<40)=\text{NO}, \text{producttype}=\text{B}\}$ .

III. The subsequent tuple to classify is  $X = \{\text{Married}=\text{NO}, \text{service}=\text{no}, \text{kids}=\text{YES}, \text{agesubsection}(<40)=\text{NO}, \text{producttype}=\text{A}\}$

$$\begin{aligned} P(X/A) &= 2000/3000 * 3000/3000 * 3000/3000 * 2000/3000 * 1000/3000 = 36/15 \\ P(X/A) &= 2.4 \\ P(X/B) &= 2000/3000 * 0 * 0 * 1000/2000 * 1000/2000 = 0 \\ P(X/C) &= 0 * 2000/3000 * 2000/3000 * 2000/3000 * 2000/3000 = 48/15 = 0 \\ P(X/D) &= 0 * 0 * 2000/2000 * 1000/2000 * 1000/2000 = 0 \end{aligned}$$

$P(X/C_i)P(C_i)$  is being computed to find the class that maximizes  $C_i$ ,

$$\begin{aligned} P(X/\text{Policy type}=A) &= 0.3 * 2.4 = 0.72 \\ P(X/\text{Policy type}=B) &= 0.2 * 0 = 0 \\ P(X/\text{Policy type}=C) &= 0.3 * 0 = 0 \\ P(X/\text{Policy type}=D) &= 0.2 * 0 = 0 \end{aligned}$$

Therefore the naïve Bayesian classifier predicts Policy type=A (Savings Plans) for tuple X. Whereas X is observed as  $\text{Married}=\text{no}, \text{service}=\text{no}, \text{kids}=\text{YES}, \text{agesubsection}(<40)=\text{NO}, \text{producttype}=\text{A}$

IV. The tuple to classify is  $X = \{\text{Married}=\text{NO}, \text{service}=\text{no}, \text{kids}=\text{YES}, \text{agesubsection}(<40)=\text{NO}, \text{producttype}=\text{B}\}$

$$\begin{aligned} P(X/A) &= 2000/3000 * 3000/3000 * 3000/3000 * 2000/3000 * 2000/3000 = 72/15 = 4.8 \\ P(X/B) &= 2000/2000 * 0 * 0 * 1000/2000 * 1000/2000 = 0.0 \\ P(X/C) &= 0 * 2000/3000 * 2000/3000 * 2000/3000 * 1000/3000 = 0.0 \end{aligned}$$

$$P(X/D) = 0 * 0 * 2000/2000 * 1000/2000 * 1000/2000 = 0.0$$

$P(X/C_i)P(C_i)$  is being computed to find the class that maximizes  $C_i$ ,

$$\begin{aligned} P(X/\text{Policy type}=A) &= 0.3 * 4.8 = 14.4 \\ P(X/\text{Policy type}=B) &= 0.2 * 0 = 0 \\ P(X/\text{Policy type}=C) &= 0.3 * 0 = 0 \\ P(X/\text{Policy type}=D) &= 0.2 * 0 = 0 \end{aligned}$$

Therefore the Naïve Bayesian classifier predicts Policy type=A (Savings Plans) for tuple X. Which maximized  $C_i$ . The value X is observed as  $X = \{\text{Married}=\text{no}, \text{service}=\text{no}, \text{kids}=\text{YES}, \text{agesubsection}(<40)=\text{NO}, \text{producttype}=\text{B}\}$ .

Bayesian classifier approach to insurance dataset observes customer preference towards the Savings Plans policy type based on attributes characterized mainly based on significant contribution of Marital status of the customer. The naïve bayesian classifier makes the assumption of class conditional independence, that is given the class label of a tuple, the values of the attributes are assumed to be conditionally independent of one another. The above analysis will help insurance advisors

## 5. Conclusion

Classification approach is applied for predicting the customer behaviour in Insurance domain for their preference towards life insurance products. Posteriori classification process is applied by looking at the data. Naïve bayes classification method is used to conduct policy preferences of life insurance customers. Naïve bayes classifier is one of the effective classifier in comparison to decision tree and neural network classifier have found it to be comparable to all other classifiers. The result of the analysis demonstrate that Naïve bayes can potentially be effective in conducting customer preference analysis over life insurance products. However, this paper observed the KDD/DM application in insurance domain, other issues may also be significant, such as considering other customer attributes (age group wise) policy preferences toward the insurance products.

## 7. References:

1. [Marisa S. Viveros, BM Research Division T. J. Watson Research Center ]Applying Data Mining Techniques to a Health Insurance Information System, Proceedings of the 22nd VLDB Conference Mumbai (Bombay), India, 1996
2. Kamakura, Wagner A. & Michel W. (2000). Factor analysis and missing data, *Journal of Marketing Research*, Vol.37, pp. 490–498.
3. Berry A.J. Michael , & Linoff Gordon S. (2000) *Mastering Data Mining: The Art and science of Customer Relationship Management* , Wiley, New Jersey
4. Saundra Glover , Patrick A Rivers , Derek A Asoh , Crystal N Piper and Keva Murph , Data mining for health executive decision support: an imperative with a daunting future., *Health Services Management Research* ,2010 Volume 23, Number 1 > Pp. 42-46
5. Data Mining concepts and Techniques , Jiawei Han and Micheline Kamber, Morgan Kaufmann Publishers, 2006.
6. Billungual report of Insurance Regulatory authority of India -2010-2011
7. Zhiyuan Yao, Annika H. Holmbom, Tomas Eklund and Barbro Back, Combining Unsupervised and Supervised Data Mining Techniques for Conducting Customer Portfolio Analysis, *ICDM 2010. LNAI 6171*, pp.292-307, 2010
8. Khurana Sunayna. (2008). "Customer references in Life Insurance Industry in India", *The ICFAI University Journal of Services* vol. 6(3), 61-68.
9. Chien-Hsing Wu,, Shu-Chen Kao, Yann-Yean Su, Chuan-Chun Wu, Targeting customers via discovery knowledge for the insurance industry, 0957-4174/\$ - see front matter q 2005 Elsevier Ltd.
10. Hokey Min, Developing the Profiles of Supermarket Customers through Data Mining[J], *The Service Industries Journal*, Vol.26, No.7, October 2006, pp.747–763