

Sales Prediction : Analysis of Time Series Data Using K-Means

Based Smooth Subspace Clustering

M. Jayasri*, S.Santhoshkumar

Department of Computer Science, Alagappa University, Karaikudi, Tamil Nadu, India

ABSTRACT

The large database to mine information that is a Data mining process and convert it into a reasonable structure for further use. Launched and order to support the organization is decision making, business planning and Data mining techniques. Data analysis, increase profitability, innovation, efficiency in resource utilization is based on important management tool in data mining. Today companies gains competitive advantage from collecting past data and using for future forecasting. Past data and information based on future estimates. In this paper, the research subject is selected as the data of a consumer electronics store company. Two year Time series sales amount data of consumer electronics was used and grouped as four quarters in a year. Next year's regression equations and naive bayes classifier methods and comprised by real sales amounts using the first quarter sales are forecasted. The real amounts and seasonal factors are really important to some product ranges that are near the sales forecasts results. In this context, various campaigns and marketing approaches have been proposed for the sales of company products by evaluating the forecast results.

Keywords: Data mining, K-means, Clustering, Time series data.

I. INTRODUCTION

The large datasets is an effort of discover patterns that is a connection of the process in data mining field (the KDD process is one of the Data mining analysis step). It use method at the relationship of statistics, database system, machine learning and artificial intelligence. The on the whole the Data mining is used to extract hidden and useful information from various databases and convert it into a reasonable structure for further use [15]. The data mining tasks to verify of the preprocessing step. That is considered as the time series segmentation [12].

Regression analysis is used to the dependent variable are related to among the independent variables and the relationship between the independent and depend variables. There is different Data Mining techniques are built-up and used in Data Mining analysis work. In recent times mutually with prediction, clustering, classification, association and sequential patterns. Data mining process is to identify the data points is called prediction. It is explanation of another related data value. Generally prediction is used for the Regression analysis.



Classification based on the machine learning and that is a one of the classic Data mining technique. It is used to classify each time in a set of data into one of predefined set of classes or groups [13]. Classification methods apply of mathematical technique such as linear programming, decision tree, neural network and statistics [13].The clustering is Different from classification [14].it is creating a meaningful and useful cluster of objects and clustering is a Data mining technique. So, in this paper proposed work a stock marketing for analyze the time series data.

II. RELATED LITERATURE

2.1 Paul Goodwin, Karima dyussekeneva and Sheik meeran "The use of analogies in forecasting the annual sales of new electronics products" 2012. One of the most popular is the bass model. The Mathematical models are often used to describe the sales and adoption patterns of products following in this year. New products are problematical because there are main time series data to evaluation the model parameters, so this model to forecast sales time series. Earlier time period to launch the sales time series of related products that is one possible solution is to fit the model and to assume that the parameter values identified for the analogy are applicable to the new product. This paper focus on the study of analogies to find the efficiency on electronic products marketing in USA for forecasting methods. It is found that all the methods have tendency to lead to forecast with high absolute percentage faults, which is consistent with other studies of new product sales forecasting. The use of published parameter values for analogies led to higher errors than the parameters we estimated from our own data [16]. When using our own data, averaging the parameter values of multiple analogies rather than relying on a single most-similar, product led to improved accuracy. However, there was little to be gained by using more than five or six analogies.

2.2. Michael Schaidnagel_, Christian Abele_, Fritz Lauxy, Ilia Petrov "Sales Prediction with

Parameterized Time Series Analysis "2013.When forecasting sales figures, not only the sales history but also the future price of a product will influence the sales quantity. At first sight, multivariate time series seem to be the appropriate model for this task. Nonetheless, in real life history is not always repeatable, i.e. in the case of sales history there is only one price for a product at a given time. This complicates the design of a multivariate time series. However, for some seasonal or perishable products the price is rather a function of the expiration date than of the sales history. This additional information can help to design a more accurate and causal time series model. The proposed solution uses a univariate time series model but takes the price of a product as a parameter that influences systematically the prediction [20]. The price influence is computed based on historical sales data using correlation analysis and adjustable price ranges to identify products with comparable history. Compared to other techniques this novel approach is easy to compute and allows to preset the price parameter for predictions and simulations. Tests with data from the Data Mining Cup 2012 demonstrate better results than established complicated time series methods.

2.3 Nirav Shah, Mayank Solanki, Aditya Tambe, Dnyaneshwar Dhangar "Sales Prediction Using Effective Mining Techniques" 2015. Mining many item sets from the large transactional database is a very serious and main task. Applications requiring large amount of data processing, consists of two huge problems, one is high storage and its management and the other one is the processing time, due to increase in data. Distributed databases do the work of solving the first problem to a tremendous extent but second problem boosts. As current era is of communication and association and people are interested in storing large data on networks, and hence, researchers are introducing various algorithms to boost the throughput of output data over distributed databases. In our research, we are introducing a algorithm to practice large amount of data at the various servers of same company that lie at different locations and collecting the practiced data on main server machine as much as admin is requiring[19]. The local copy of found data is provided to the users if he/she needs it again, this allows causing a proxy server where constantly searched items can be saved with the density of their access. This not only grants affording fast access to the data but will also afford to maintain list of recurrently accessed data. There are several approaches for accessing the data from the various servers, such as direct networked access, mobile agents, client-server techniques and LAN etc. We have used multi-threaded environment to calculate various distributed servers to gather data. For processing of data at the server end, the use of Apriori Algorithm has been done to get the outputs, which are then addressed to the client. At client data from various servers is assembled and then disciplined into data format. As an association rule mining is defined as the relation between various item sets. Association rule mining takes part in pattern discovery techniques in knowledge discovery and data mining (KDD). The frequent item sets mining, is depends of association rule mining as performance. To mine frequent item set efficiently and necessary. Thus is necessary to mine frequent item set efficiently. Association rules arrange information of this type in the form of "ifthen" statements. These rules are count from the data and, inconsistent the if-then order of logic, association rules are probabilistic simultaneously, the antecedent (the "if" part) and the consequent (the "then" part), an association rule has two numbers that explicit the degree of ambiguity about the rule.

2.4 Xin Xu Lei Tang Venkat Rangan " Hitting your number or not? A Robust & Intelligent Sales Forecast System" 2017.Enterprise sales forecasting is the use of historical data to determine the future outcome (typically, sales revenue) of target month, quarter or year. Accurate forecast helps companies efficiently allocate workforce or funds, pinpoint the revenue growth bottleneck, and strategically pivot the company. Furthermore, reporting accurate projections to the board of directors or Wall Street is a key factor in stock price stability and the overall health of the business. In this work, we present one system to provide accurate, robust forecast for enterprise sales, enabling real-time insights for business decisions. Before we present the system, we will first introduce background of enterprise sales and terms commonly used in sales domain, then discuss about requirements and challenges associated with the forecast problem [18]. Enterprise selling, also known as business-tobusiness (B2B) selling, typically is much more complex than B2C (business-to-consumer) sales. The B2B sales normally go through a funnel, from prospect or lead, to opportunity, and then to customer. Prospects are identified through many marketing campaigns. Events like website visit, eBook download, dropping business cards at some event imply a propensity in seller's product or service. Leads are normally generated by identifying similar companies or organizations to existing clients or prospects. Similarity can be defined in a variety of aspects. Common ones are based on firm graphic information like company size, target vertical, technology stack, and so on.

2.5 Vicky Chrystian Sugiarto1, Riyanarto Sarno,

Dwi Sunaryono " Sales Forecasting Using Holt-Winters in Enterprise Resource Planning At Sales and Distribution Module " 2016.Enterprise Resource Planning is a system used for managing all of the resources owned by the company, business activities, and information used to make a good business process. One of the modules in the ERP is sales and distribution. Sales and distribution is a module that handles the sale and delivery goods to customers to achieve their business objective The use sales and distribution module is intended to simplify the process of selling to customers in accordance with the interests of customers for goods and services, makes it easy to check the sale of goods and delivery of goods, and facilitate the collection of customers who make a purchase of goods or services provided by the company, determines the appropriate services to customers, and forecast the amount of demand for goods desired by the customer[17]. Sales forecasting is calculating the expected sales of a specific product and predicting future sales of the product [17]. It helps in making informed business decisions. In this article some several of sales forecasting methods used, Holt-Winters Additive Method, and Holt-Winters Multiplicative Method applied in ERP. The implementation of sales forecasting using Holt-Winters method in the ERP can make the system more accurate and efficient in the determination of the amount of goods demanded by customers.

III. EXISTING WORK

Classification similarity based the stock prediction problem can be mapped. A set of vectors mapped into the historical stock data and the test data. Every vector denote N dimension for each features. The parallel metric such as Euclidean distance is computed to take a choice. In this part, an explanation of KNN is provided, the highest similarity to the test into closest k records of the training data set. (I.e. query record) then a maximum vote is performed among the selected k records to verify the class label and then assigned it to the query record.

Limitations

- ✓ Need to determine value of parameter K (number of Nearest neighbors)
- ✓ Distance based learning is not understandable which type of distance to apply and which attribute to use to produce the top result.
- ✓ Computational Cost is high because we want to compute distance of every query instance to all training samples.

 ✓ Can be used for both binary and multiclass classification problems. Regression analysis,

IV. PROPOSED WORK

This paper proposed the K-means clustering algorithm with regression for handling time series resources or data efficiency. The initial step for this research work is to omit the noisy data and inconsistency of data from the Cylinder-Bell-Funnel (CBF) dataset. The next step, the input data is grouped in the form is called clustering based on same time along with the usage of K-means algorithm, the result of clustered data is a sales marketing data. The experimental result of clusters constructs. In order to invent the closest level of sales time analyzing points, used for the regression using statistical tool for regression analysis that is used to investigate the association through the variables. The analyzer seeks to taken the impact of one variable upon then another is focusing on the association ship between a dependency variable and one or more independency variables. Usage of previous sales records the K-means algorithm makes clusters and attributes such as volume of sales, price or cost of sales, profit and volume of seasonal sales volume data are invented from the sales numbers of same quantization.

$$J(\mathcal{V}) = \sum_{i=1}^{c} \sum_{j=1}^{c_i} (\|\mathbf{x}_i - \mathbf{v}_j\|)^2$$

Advantages

- Very simple, easy to implement and fast.
- Need less training data.
- Highly scalable. It scales linearly with the number of predictors and data points.

$$J(V) = \sum_{i=1}^{c} \sum_{j=1}^{c_i} \left(\left\| x_i - v_j \right\| \right)^2$$

Advantages

- 1. Very simple, easy to implement and fast.
- 2. Need less training data.
- 3. Highly scalable. It scales linearly with the number of predictors and data points.
- 4. Can be used for both binary and multiclass classification problems.
- 5. Can make probabilistic predictions

V.RESULTS AND DISCUSSION

Today data's are most factors in the real world. The k means algorithm is based on data analysis the accurate data. It is used in sales marketing and stock marketing. The correct and efficient data calculate based on these figures.

							- 17 E			
	SELECT DATASET FILE		E Veyesi p	ojectiTime serie	s k -means/Tim	ac 🔲	BROWSE			
Ť	1	2	3	4	5	6	1	8	9	10
1	-2.2296	-1.2768	-1.9605	-1.7440	-2.1088	-2.1964	-1.2023	-1.1945	-1.7623	-1.39
2	-1.5653	-1.1942	-1.0642	-0.5903	-1.4952	0.7258	-0.9999	-1.0233	0.6366	-0.10
3	-0.8592	-1.0128	-1.2683	-1.2222	-0.3406	-1.2058	-1.4650	-1.9021	-1.1335	-0.43
4	-1.0089	-0.8100	-0.6691	-1.5807	-1.9342	-0.9122	-1.0436	-0.3631	-1.0051	457
5	-1.9016	-1.5264	-1.4133	-1.2672	-1.6306	-1.1906	-2.0094	-1.6504	+1.1518	-1.40
6	-1.2019	-1.3107	-1.1319	-0.6014	-0.6089	-0.4910	-1.2715	-1.0151	-1.9914	-0.45
7	-4.8637	-0.5702	-0.4542	-0.1720	-1.3312	-0.6455	-0.6061	-1.2033	-4 5493	-1.92
8	-4 9782	-0.2074	-0.2886	0.2711	-0.3120	-0.1989	-1.3647	-0.5604	-0.8683	-0.79
9	4.6782	-1.5484	-0.9602	-1.5109	-1,2772	-1.4002	-1.0413	-4.9486	-0.7061	-0.60
10	-1.4148	-1.1366	-0.5780	-0.9379	-0.4310	-1.1414	-0.9421	-1.4739	-1.6612	-1.23
11	-1.8136	-0.9060	-0.7822	-1.4620	-1.2398	-0.9643	-0.9622	+1.2578	-1.1065	-0.85
12	-4.9370	-1.8022	-0.8443	-4.9057	-1.6075	-1.0549	-0.8575	-1.0006	-2.0369	-0.79
13	-0.9361	-0.8501	-1.0282	-1.3606	-1.1033	-0.7588	-1.9615	-0.8002	-1.0371	-1.53
14	-1.6950	-0.1547	-0.4708	-4.2696	-4.6872	-1.4831	0.1847	-0.0564	-2.6865	-4.95
15	4 6395	-0.5275	-0.6045	-4.6204	-1.1412	-0.1902	-1.0023	-4.5506	-1 2707	-0.98
16	-1.9053	-1.4883	-0.8540	-2.1176	-4.8665	-0.7132	-0.6780	-0.9585	-1.0480	-0.89
17	49428	-0.6154	-0.5827	-0.5391	-0.9041	-0.9547	-0.4771	-0.4760	-0.7446	-0.49
18	-1.0756	-0.2637	-0.4091	-4.3237	-1.0352	-0.7282	-1.1438	-1.1462	-0.2914	-0.24
19	-1.0224	-0.9892	-1.0411	-0.2778	-0.1075	-0.1492	-0.6558	-1.4620	-1.0398	-1.29 -
19	-1 0/56	-0.2637	-0.4091	-0.2778	-1.0252 -0.1075	-0.1492	-0.6558	-1.1462	-1 0398	

Figure 1 : Select data and Preprocessing The Fig (1) is put to select the CBF data and it is utilize the initiate the process of the data

	1	2	3	4	5	6	7	1		10
1	2.4547	.3.6579	-2.1163	-2.2097	-2.6509	.2.4291	-2.0754	-1.8921	.17562	-2.80 +
2	-23754	-2.5797	-2.4960	-1.7434	-1.7753	-1.5716	-2.5166	-2.1932	-2.0404	-1.65
1	-2.5412	-1.6395	-2.3610	-3.0361	-2.7965	-1.9036	-1.4341	-2.1681	-2.6247	-1.70
1	-2203	-1.5357	-1.8325	-1.6402	-1.6701	-2.4223	-1.5139	-1.1288	-1.3842	-2.02
5	-2.2296	-12768	-1.9605	-1.7440	-2.1088	-2.1964	-1 2023	-1.1946	-1.7623	-1.39
6	-2.2057	-2.6437	-1.5156	-1.2548	-1.7899	-1.9516	-1.5405	-1.5401	-1.3510	-1.32
1	-2.1825	-1.3343	-0.6940	-1.1085	-4.7997	4.8470	-2.1144	-1.7096	-6-6267	-0.41
1	-2.1683	-1.8207	-1.2923	-2.7240	-2.1501	-2.6282	-2.1038	-1.1271	-1.4168	-1.50
2	-2.1453	-11062	-2.0217	-1.5396	-26113	-2.8040	-2.9099	-1.3711	-1.8062	-2.17
0	-2.1302	-3.2287	-2.3709	-1.6274	-2.1329	-2.4257	-1.8467	-2.3155	-2.0901	-2.61
1	2.1218	-1.6788	-1.5592	1.0459	-5.1105	-1.6540	-1.0404	-1.6110	-1.0041	4.22
2	-2 0900	-1.8596	-1.9443	-2.9401	-1.5962	-1.7853	-2.0252	-1.9019	-1.6893	-1.54
3	-2.0218	-5.5902	-1.3200	-1.5860	-0.7408	-2.2791	-1.5658	-1.1219	-2.2914	-1.50
4	-20151	-1.9021	-2.1871	-1.9155	-2.2699	-2.2638	4.7701	-2.0588	-1.3886	-1.95
5	-2.0105	+1.2309	+1.5756	-1.1572	-1.6459	+1.2568	-2.0165	-1.5007	-2.1157	-1.62
6	-1.9079	-2.5253	-1.9097	-1.9829	-20804	2.1441	-2.5321	-1.0001	-1.3667	-1.25
1	-1.9031	-1.5386	-1.5892	-1.8024	-2.1141	-2.2196	-2.0224	-5 7331	-2.5122	-1.09
1	-1.9013	-2.4229	-1,7295	-1.3167	-3.5473	-1.9633	-2.4301	-2.2455	-2.0628	-2.54
9	-1.9757	-1.8671	0.2541	-1.2433	-1 2010	-19957	-1.0530	-1.1003	-1.2277	-126.*

Figure 2 : Clustering the data

The fig (2) is discussed of Selected data's are applied for preprocessing.



Figure 3:Timeseries of cluster1



Figure 4 :Timeseries of cluster 2



Figure 5 :Timeseries of cluster3

The cluster data's are all time series CBF file data clustering for three times above the figures 3, 4, 5. And finally figure 6 is explained and nalyzes the CBF file data. Finally to produce the cluster and analyze the CBF file data.



Figure 6: Timeseries of CBF Data

II. CONCLUSION

Although a conclusion may review the main points of the paper, do not replicate the abstract as the conclusion. A conclusion might elaborate on the importance of the work or suggest applications and extensions. Authors are strongly encouraged not to call out multiple figures or tables in the conclusion—these should be referenced in the body of the paper.

III. REFERENCES

 Karnani A. (2007), The Mirage of Marketing to the Bottom of the Pyramid, Global Competition, The William Davidson Institute at the Michigan State University, p. 90-113

- [2]. Kress, G. (1985). Linguistic processes in sociocultural practice. Victoria: Deakin University Press, p. 76-77
- [3]. Mentzer J.T. and Kahn K.B, (1995), Journal of Forecasting, Volume 14, Issue 5, p. 465–476
- [4]. Greene K.C. and Armstrong J. S. (2007), The Ombudsman: Value of Expertise for Forecasting Decisions in Conflicts. 1–12.
- [5]. Delurgio S.A. and Bhame C.D. (1991), Forecasting systems for operations management, p. 648 -649,
- [6]. Freedman D.A. (2005), Statistical Models: Theory and Practice, Cambridge University Press, p. 36
- [7]. Armstrong, J. S. (2012), Illusions in Regression Analysis, International Journal of Forecasting, p.689-696
- [8]. Seal H. L. (1967), The historical development of the Gauss linear model, Biometrika, volume 54, p.1–24
- [9]. Stuart R.,Norvig P. (2003), Artificial Intelligence: A Modern Approach, Prentice Hall, p. 126-132
- [10]. Murty N., Susheela D. (2011), Pattern Recognition: An Algorithmic Approach, Springer Publishing, p. 86 -102
- [11]. Bednarz T. F., (2011): Sales Forecasting: Pinpoint Sales Management Skill Development Training Series, Majorium Busine
- [12]. Tai.Yu, Chiu, ting-Chie Hsu and Jia-Shung Wang (2015) ,Interpolation based consensus clustering for gene expression time series.
- [13]. Paul cotofrei, kilian Stoffel (2000), Classification time=Temporal Rules.
- [14]. N.Karthikeyeni Visalakshi, J.Suguna 28th North American Fuzzy Information Processing Society Annual Conference, Cincinnati, ohio,USA (June 14-17,2009).
- [15]. Data mining clustering Lecturers jery stefanowski, Institute of Computer Science Poznan University of Technology Poznan, Poland (2008/2009).
- [16]. Paul goodwin, Karima dyussekeneva and Sheik meeran "The use of analogies in forecasting the annual sales of new electronics products" 2012
- [17]. Vicky Chrystian Sugiarto1, Riyanarto Sarno, Dwi Sunaryono " Sales Forecasting Using Holt-

Winters in Enterprise Resource Planning At Sales and Distribution Module " 2016

- [18]. Xin Xu Lei Tang Venkat Rangan " Hitting your number or not? A Robust & Intelligent Sales Forecast System" 2017
- [19]. Nirav Shah, Mayank Solanki, Aditya Tambe, Dnyaneshwar Dhangar "Sales Prediction Using Effective Mining Techniques" 2015.
- [20]. Michael Schaidnagel_, Christian Abele_, Fritz Lauxy, Ilia Petrov "Sales Prediction with Parameterized Time Series Analysis "2013.R.M.Sahu, Akshay Godase, Pramod CONTROL ENGINEERING, Vol. 4.
- [21]. Kanchan Mahajan, Proff.J.S.Chitode, "Waste Bin Monitoring