Vol. 6, No. 2, 2021

RECOMMENDATION SYSTEM FOR PURCHASING GOODS BASED ON THE DECISION TREE ALGORITHM

Yurii Kohut, Iryna Yurchak

Lviv Polytechnic National University, 12, S. Bandery Str., Lviv, 79013, Ukraine. Author`s e-mail: yurii.kohut@ukr.net

https://doi.org/10.23939/acps2021.02.121

Submitted on 03.10.2021

© Kohut Y., Yurchak I., 2021

Abstract: Over the past few years, interest in applications related to recommendation systems has increased significantly. Many modern services create recommendation systems that, based on user profile information and his behavior. This services determine which objects or products may be interesting to users. Recommendation systems are a modern tool for understanding customer needs. The main methods of constructing recommendation systems are the content-based filtering method and the collaborative filtering method.

This article presents the implementation of these methods based on decision trees. The content-based filtering method is based on the description of the object and the customer's preference profile. An object description is a finite set of its descriptors, such as keywords, binary descriptors, etc., and a preference profile is a weighted vector of object descriptors in which scales reflect the importance of each descriptor to the client and its contribution to the final decision. This model selects items that are similar to the customer's favorite items before. The second model, which implements the method of collaborative filtering, is based on information about the history of behavior of all customers on the resource: data on their purchases, assessments of product quality, reviews, marked product. The model finds clients that are similar in behavior and the recommendation is based on their assessments of this element. Voting was used to combine the results issued by individual models - the best result is chosen from the results of two models of the ensemble. This approach minimizes the impact of randomness and averages the errors of each model. The aim: The purpose of work is to create real competitive recommendation system for short period of time and minimum costs.

Index Terms: Recommendation system, Decisions tree, Quick decision making, Big Data, Analysis of Big Data.

I. INTRODUCTION

With the development of globalization, the spread of IT technology and general computerization [1], which leads to the creation of large data sets [2], there is a need to analyze them and make quick, but rather "instant" solutions. This process has contributed to the emergence of tools and algorithms for fast information processing. One such tool is the "decision tree" algorithm, which is an effective tool for data mining and predictive analytics [3]. A recommendation system is a tool that actively finds information that may be of interest to a user from a large amount of information [4]. Building a system that supports online user decisions, recommending a personalized, highly-matched product or project is a core issue in the recommended system area [5]. It can be traced back to cognitive science, approximation theory, information retrieval, prediction theory, management science, and customer selection models in the market [6].

In view of the theoretical and practical application value of the recommendation system, this paper reviews the research progress of the recommendation system, and attempts to lay a foundation for further research on the recommendation system theory and the expansion of its application field.

The purpose of this work is to implement models of the recommendation system based on decision trees. This approach is appropriate because the decision tree algorithm is easy to implement, capable of processing large amounts of non-normalized data.

II. FORMULATION OF THE PROBLEMS

Due to the rapid transition of business to the Internet, marketers who have not yet considered it appropriate to conduct a marketing strategy on the Internet, the question "how to find out what a site visitor (potential customer) wants without asking him?". In search of an answer to this question, companies have accumulated terabytes of information "who, what, when, where and how" buys, in the future with the help of algorithms to determine what their customer wants. This is roughly how the concept of Big Data came about.

Big data allows us to see and understand the connections between pieces of information that until recently we were just trying to capture [7]. Due to the rapid spread of smart and interconnected devices and systems, the amount of data collected is growing at an alarming rate. In some industries, about 90% of data is stored unstructured, and their volume increases by 50% annually. With regard to big data analysis and other analytical tasks, current solutions do not provide the system response speed required to work with analysis tasks, which reduces user productivity and delays the decision-making process [8]. Business methods are changing. Consumer behavior is changing. Consumers

themselves are changing. To stay competitive, businesses seek to know in real time when customers are buying something, where they are buying, and even what they think before going to the store or visiting a Web site. Big data, Big Data analysis, and an integrated platform for business intelligence (BI) and Big Data analysis can help [9].

III. THE PROCESS OF BUILDING A DECISION TREE

The decision tree is a method of representing the decision rules in a hierarchical structure consisting of two elements: nodes (Fig. 1) and leaves (Fig. 2). In nodes there are decisive rules, and check of conformity of examples to this rule on any attribute of set is carried out.

In the simplest case, according to the results of the test, most of the examples that fall into the node are divided into two categories: those that fall into the examples, and those that do not satisfy the condition.





Fig. 2. Leaves

Then the rule is applied to each category again and the procedure is repeated recursively until the condition of stopping the algorithm is reached. As a result, in the last node check and division is not carried out and it is declared by the letter. The letter determines the solution for each example in it. For a classification tree, it is a class associated with a node, and for a regression tree, it is a letter-modal interval of the target variable.

The main area of application of decision trees is to support management decision-making processes used in statistics, data analysis and machine learning. The tasks that are solved with this device are:

• Classification – assigning objects to one of the previously known classes. The target variable must have discrete values.

• Regression (numerical prediction) – prediction of the numerical value of the independent variable for a given input vector.

• Object description – a set of rules in the decision tree allows you to compactly describe objects.

Therefore, instead of complex structures that describe objects, you can store decision trees.

A. THE DECISION TREE DESIGNING.

Algorithms for designing decision trees consist of stages "construction" or "creation" of a tree (tree building) and "reduction" of a tree (tree pruning) [5–9]. Tasks of choosing the criterion for splitting the algorithm and stopping learning are resolved during the creation of the tree, if it is provided by algorithm. During the stage of tree reduction, the issue of cutting off some of its branches is solved.

The process of creating a tree is from top to bottom, i.e. is descending. During the process, the algorithm must find such a splitting criterion, sometimes also called breakdown criterion, to split the set into subsets that would be associated with a given validation node. Each test node must be marked with a specific attribute. There is a rule for selecting an attribute: it must split the original data set so that the objects of the subsets resulting from this breakdown are members of the same class or are as close as possible to that breakdown. The last phrase means that the number of objects from other classes, the so-called "impurities", in each class tended to a minimum. There are different criteria for cleavage. The most famous are the measure of entropy and the Gini index. Some methods use the so-called measure of attribute subspaces, which is based on the entropy approach and is known as the "information gain measure" or entropy measure, to select the split attribute. Another splitting criterion proposed by Breiman et al. Is implemented in the CART algorithm and is called the Gini index. With this index, the attribute is selected based on the distances between class distributions. If a given set T, including examples of n classes, the index Gini, ie gini (T), is determined by the formula:

$$gini(T) = 1 - \sum_{j=1}^{n} p_j^2$$

T – the current node, pj – the probability of class j in the node T, n – the number of classes.

B. MAIN STAGES OF CONSTRUCTION.

During the construction of the decision tree you need to solve several main problems, each of which involves a corresponding step in the learning process [10]:

• Select the attribute that will be partitioned in this node (partition attribute).

• Choice of the criterion of stopping learning.

• Choice of method of cutting off branches (simplification).

• Estimation of accuracy of the constructed tree.

C. CHOICE OF ATTRIBUTE.

When creating a rule for partitioning in the next node of the tree, you must select the attribute by which it will be done. The general rule for this can be formulated as follows: the selected attribute must break the set of observations in the node so that the resulting subsets contain examples with the same class labels, or were as close as possible, ie the number of objects from other classes ("impurities") in each of these sets was as small as possible. Various criteria were chosen for this, the most popular of which were theoretical-informational and statistical.

D. ALGORITHM STOP CRITERION.

Theoretically, the algorithm for learning the decision tree will work until the result is absolutely "pure" subsets, each of which will be examples of one class. However, it is possible that a tree will be built, in which a separate sheet will be created for each example. Obviously, such a tree will be useless, because it will be retraining – each example will correspond to its unique path in the tree, and hence the set of rules relevant only for this example.

Retraining in the case of the decision tree leads to the same consequences as for the neural network – accurate recognition of examples involved in learning and complete failure on new data. In addition, tree retraining has a very complex structure and is therefore difficult to interpret.

The obvious solution is to force the tree to stop building until it has been retrained. The following approaches have been developed for this purpose.

1. Premature stop – the algorithm will be stopped as soon as the specified value of a criterion is reached, such as the percentage of correctly recognized examples. The only advantage of the approach is the reduction of training time. The main disadvantage is that early stopping is always done to the detriment of the accuracy of the tree, so many authors recommend that you cut off the branches.

2. Limiting the depth of the tree – the task of the maximum number of breaks in the branches, after which the training stops. This method also reduces the accuracy of the tree.

3. The task of the minimum number of examples per node is to prohibit the algorithm to create nodes with the number of examples less than the specified (for example, 5). This will avoid the creation of trivial partitions and, consequently, insignificant rules.

All these approaches are heuristic, and do not guarantee a better result or work only in some cases. Therefore, their use should be approached with caution. At present, there are no reasonable recommendations as to which method works best. Therefore, analysts have to use the method of trial and error.

E. CUTTING OF BRANCHES.

As mentioned above, if the "growtha" of the tree is limited, the result will be a complex tree with a large number of nodes and leaves. As a result, it will be difficult to interpret. At the same time, the decisive rules in such trees, which create nodes, which fall into two or three examples, are insignificant from a practical point of view.

It is much better to have a tree consisting of a small number of nodes, which would correspond to a large number of examples from the training sample. Therefore, an alternative approach to early stopping is to build all possible trees and choose the one that at a reasonable depth provides an acceptable level of recognition error, ie to find the most favorable balance between the complexity and accuracy of the tree.

An alternative approach is the so-called pruning. It contains the following steps:

1. Construct a complete tree (so that all the leaves contain examples of one class).

2. Determine two indicators: the relative accuracy of the model – the ratio of the number of correctly recognized examples to the total number of examples, and the absolute error – the number of incorrectly classified examples.

3. Remove leaves and knots from the tree, the cutting of which will not significantly reduce the accuracy of the model or increase the error.

The cutting of branches, obviously, is carried out in the direction opposite to the direction of growth of a tree, ie from the bottom to the top, by consecutive transformation of knots into leaves. The advantage of cutting off branches compared to early stopping is the ability to find the optimal ratio between accuracy and clarity of the tree. The disadvantage is more training time due to the need to first build a complete tree.

F. EXCERPT FROM THE RULES.

Sometimes even a simplified decision tree is still too complex to visually perceive and interpret. In this case, it may be useful to extract the decision rules from the tree and organize them into sets describing the classes.

To remove the rules, you need to trace all the paths from the root node to the leaves of the tree. Each such path will give a rule consisting of a set of conditions that represent a check in each node of the path.

Visualization of complex decision trees in the form of decision rules instead of a hierarchical structure of nodes and leaves may be more convenient for visual perception.

G. ADVANTAGES AND DISADVANTAGES.

Advantages:

• Intuitiveness of decision trees. The classification model, presented as a decision tree, is intuitive and simplifies the understanding of the problem to be solved.

• The result of the algorithms for designing decision trees, in contrast to, for example, neural networks, which are "black boxes", is easily interpreted by the user [11]. This property of decision trees is not only important when assigning a new object to a particular class, but is also useful in interpreting the classification model as a whole. The decision tree allows you to understand and explain why a particular object belongs to a particular class.

• Decision trees allow you to retrieve rules from the database in plain language. Example rule: If Age > 35 and Income > 200, then issue a loan.

• Decision trees allow you to create classification models in areas where it is difficult for the analyst to formalize knowledge.

• The algorithm for constructing the decision tree does not require the user to select the input attributes of independent variables). At the input of the algorithm you can submit all existing attributes, the algorithm will choose the most significant among them, and only they will be used to build a tree. Compared to, for example, neural networks, this greatly facilitates the user's work, because in neural networks, the choice of the number of input attributes significantly affects the learning time.

• Accuracy of models created with the help of decision trees in comparison with other methods of construction of classification models (statistical methods, neural networks).

• Developed a number of scalable algorithms that can be used to build decision trees on ultra-large databases. Scalability here means that as the number of database examples or records increases, the time spent learning, ie building decision trees, increases linearly. Examples of such algorithms: SLIQ, SPRINT.

• Fast learning process. It takes much less time to build classification models using decision tree construction algorithms than, for example, to train neural networks.

• Most algorithms for constructing decision trees have the ability to specifically process missing values.

• Many classical statistical methods used to solve classification problems can only work with numerical data, while decision trees work with both numerical and categorical data types.

• Many statistical methods are parametric, and the user must have certain information in advance, for example, know the type of model, have a hypothesis about the type of relationship between variables, assume what kind of distribution the data have. Solution trees, in contrast to such methods, build nonparametric models. Thus, decision trees are able to solve such Data Mining problems in which there is no a priori information about the type of relationship between the studied data.

Disadvantages:

• Decision trees are sensitive to noise in the input data. Small changes in the training sample may lead to

global adjustments to the model, which will affect the change in the rules of classification and interpretation of the model.

• The dividing boundary has certain limitations, due to which the decision tree on the quality of classification is inferior to other methods.

• It is possible to relearn the decision tree, which is why you have to resort to the method of "cutting off branches", setting the minimum number of elements in the tree leaf or the maximum depth of the tree.

• Complex search for the optimal solution tree: this leads to the need to use heuristics such as no feature search with the maximum increase in information, which ultimately does not give a 100 percent guarantee of finding the optimal tree.

• The decision tree makes a constant prediction for objects that are in the feature space outside the parallelepiped, which does not cover all objects in the training sample.

Despite the shortcomings, and due to the main advantages, decision trees are an important tool in the work of every specialist in data analysis.

IV. POTENTIAL DEVELOPMENT DIRECTION

A. VARIETY

The data available in the recommendation system is complex and complex. For example, information in the social network, location information, and other contextaware information are taken into account. Not only does the amount of data increase, but the computational complexity also increases exponentially. In addition, the recommendation system research involves privacy protection.

How to ensure personalized recommendation and protect user's privacy is a confrontational problem, which brings great challenges to researchers and developers.

B. INTERPRETABLE

As an important product in the field of artificial intelligence, the recommendation system is widely accepted and applied. The core of the recommendation is the rationality of the high recommendation result, which also requires the recommendation result to be well interpretable, although this has long been known. It is realized, but the special research on interpretability is still lacking. In the current research, the interpretability discussion of the recommendation algorithm is generally the selection process after the algorithm evaluation. With the high demands of users, the research of "recommendation reasons" has received more and more attention in industry and academia.

V. SYSTEM DEVELOPMENT

Let's imagine that we have opened an online store and promote it on Facebook. Let's start with advertising, for men there will be one link, for women another (Fig. 3) As everyone knows, the consumer psychology of women and men is very different (Fig. 4, Fig. 5), so the existence of identical or similar catalogs of goods and the principle of their submission is impossible.

After successful order of goods, appears a window which invites buyer to leave a feedback (Fig. 6).



Fig. 3. Links for fe(male)



Fig. 4. Algorithm of buying for man



Fig. 5. Algorithm of approximately buying by women



Fig. 6. Feedback button

The algorithm has to take into account everything:

- Recommendations (reviews);
- number of visits to the page;
- opening the page of a particular product;
- sequence of opening goods;
- trends;
- the number of preferences (Fig. 7);
- newest goods (Fig. 8);
- the greatest demand;

- the balance of the goods;
- going by category;
- gender of the buyer;
- the number of points for an individual product;
- the number of products in the basket;



Fig. 7. Button 'Leave a like'



Fig. 8. Newest goods

Summarizing all the above, the principle of the recommendation system will be as follows (Fig. 9):



Fig. 9. Operating principle of recommender system

VI. COMPARATIVE ANALYSIS OF EXISTING RECOMMENDATION SYSTEMS

It is not news that the largest sales platforms have developed and use their own algorithms to calculate user needs. The principle of operation of these algorithms is kept secret by everyone. It is also known that there is no one universal algorithm that would be suitable for all areas of business. For this reason, it is not possible to consider analogues. However, there are some common principles.

A. CONTENT-BASED RECOMMENDATION

The first step in a typical Content-Based approach is to build a User Profile. A simpler construction method is to consider all the items that the User has ever scored, and make a weighted average of the Item Profiles of these items as the UserProfile of the User. Obviously, the strategy for building a User Profile can be complex. For example, we can consider the time factor and calculate the Profile of User in different time periods to understand the changes in User's preference on historical data. With the User Profile, we can start recommending. The simplest recommendation strategy is to calculate the similarity between all the items that the user has not tried and the user's User Profile, in order of similarity. A list of recommendations is generated and output as a result. Besides, the recommendation strategy can also be very complicated, such as considering the real-time interaction data collected during the user interaction process on the data source to determine the ordering, using the decision tree and artificial neural network on the model, but the core of these methods The links are all calculated using the similarity between the User Profile and the Item Profile.

B. COLLABORATIVE FILTERING-BASED

Collaborative Filtering-Based Recommendation refers to collecting the past behavior of the user to obtain explicit or implicit information about the product, that is, according to the user's preference for the item or information, discovering the relevance of the item or the content itself, or the relevance of the user, and then Based on these associations, recommendations are made. According to the foregoing, recommendations based on collaborative filtering can be based on User-based Recommendations, based on Item-based Recommendations, and based on subclasses such as Model based Recommendation. The user's preference or scoring matrix for the item is often a large sparse matrix. In order to reduce the amount of calculation, clustering items for Collaborative Filtering can be used. Recommendation system framework classification (Fig. 10) consist of application field and data mining technology.

C. COMPARISON OF EFFICIENCY OF ANALOGUE ALGORITHMS

The effectiveness of algorithms that work on a similar principle can be said only on the basis of comprehensive research on the effectiveness of business systems, analyzing the period of development of the business structure before the introduction of the recommendation system and the period after introduction. Given the fact that it is impossible to judge the effectiveness of the algorithm based on sales figures after its introduction, because the effectiveness of the business structure is influenced not only by this factor, it seems impossible to truly assess the impact of the algorithm. However, as amateurs it can be assumed that the efficiency is calculated as follows: $\varepsilon = \alpha / \beta \times 100 \%$, \mathcal{E} – efficiency, α – number of sales after putting the algorithm into action, β – number of sales before putting the algorithm into action.

Of course, it can be told about these numbers only by analyzing the working business structures, such as the hypermarket chain Epicenter, Metro, or Target, which is an American hypermarket chain. Unfortunately,

information about Epicenter and Metro remains closed. According to research by Charles Duhigg [12], Target sales efficiency increased by 1470 %, thus gaining more than 1000 % of new customers, profits increased from 44 billion to 65 billion after the introduction of the recommendation system compared to the same period until the system worked, which can be considered an incredible breakthrough in the field of sales psychology. Not to mention the speed with which In addition, we should mention such a criterion for evaluating efficiency as the cost of the system. This cost includes the cost of research, the cost of equipment, and the cost of creating the system, in addition, an important factor is the duration of these studies. In comparison, Target (from 2002 to 2009) spent 7 years and \$ 1.500.000 to build its system. We should not forget the fact that any information about the research and development of referral systems is a secret of each company, and none of them will disclose their cards about how they obtain personal data and methods of data processing, and the cost of such work, therefore, it is very difficult to objectively judge the effectiveness and benefits of some recommendation systems over others.



Fig. 10. Classification framework of recommender system

VII. CONCLUSIONS

Modern recommendation systems are ensembles of models. Two levels are used: data processing in each model of the ensemble and the application of rating selection of the best results. Combining models helps increase accuracy and the ability to flexibly customize different customer groups.

The article presents two models that are implemented as decision trees. The first model is used to pre-determine customer preferences based on his profile data. The customer profile contains information about those items that have been marked or purchased. A time factor is involved to understand changes in the client's preferences for historical data. The similarity between all elements which the client did not try, and preferences of the user in an order of similarity is calculated. The system generates a list of recommendations.

The second model filters according to the behavior of other customers to obtain explicit or implicit information about the element, according to the preferences of customers to the element. The system provides recommendations based on these associations as well. The received recommendations from two models are compared and are chosen more probable for the client. The preferences of other customers of this resource are taken into account, which may affect the correction of recommendations for a particular customer.

The obtained results showed that the proposed approach allows to increase the accuracy of the recommendation for a particular customer, as it takes into account his personal preferences on this resource: deferred or purchased item, product quality assessment, positive feedback about the product. The data are historical, as the time factor is important for making recommendations.

The solution tree algorithm was used to implement the models. The algorithm is simple to implement and allows you to process large amounts of data, such information as sex, age, origin, payment method, transition sequence, the number of goods in the check, and etc., without prior procedures. The algorithm provides an opportunity to assess the accuracy of the model using statistical tests.

Describes and analyzes the criteria for assessing the effectiveness of recommendation systems, information on the creation of which companies keep secret. Based on the results of work for 3 months, a real working recommendation system was created, the cost of which is \$ 400, the effectiveness of this development can be judged only by testing in a real working business structure for at least a year.

REFERENCES

 Mayer-Schoenberger W. Big data. A revolution that will change the way we live, work and think / Victor Mayer-Schoenberger, Kenneth Kukier; lane. with English Inna Gaidyuk. – M.: Mann, Ivanov and Ferber, 2014. – p. 240. Available at: https://doi.org/10.1016/j.jvcir.2019.102705. (Accessed: 18 November 2021).



Yurii Kohut received a bachelor's degree in computer engineering in 2020. From 2020 he has been receiving a master's degree. Now he is second year computer engineering student at Lviv Polytechnic National University. His research interests include image processing and segmentation, object identification using neural networks.

- [2] A. Pal, P. Parhi and M. Aggarwal. An improved content based collaborative filtering algorithm for movie recommendations. 2017 Tenth International Conference on Contemporary Computing (IC3), 2017, pp. 1–3, doi: 10.1109/IC3.2017.8284357.
- [3] Ajah, I.A. Nweke, H.F. Big Data and Business Analytics: Trends, Platforms, Success Factors and Applications. Big Data Cogn. Comput. 2019, 3, 32. https://doi.org/10.3390/bdcc3020032.
- [4] Y. Roh, G. Heo and S. E. Whang. A Survey on Data Collection for Machine Learning: A Big Data – AI Integration Perspective. in IEEE Transactions on Knowledge and Data Engineering, vol. 33, no. 4, pp. 1328–1347, 1 April 2021, doi: 10.1109/TKDE.2019.2946162.
- [5] Lin, J., Zhong, C., Hu, D., Rudin, C. & Seltzer, M. (2020). Generalized and Scalable Optimal Sparse Decision Trees. Proceedings of the 37th International Conference on Machine Learning, in Proceedings of Machine Learning Research 119:6150–6160. Available at: https://proceedings.mlr.press/ v119/lin20g.html. (Accessed: 18 November 2021).
- [6] Avellaneda, F. (2020). Efficient Inference of Optimal Decision Trees. Proceedings of the AAAI Conference on Artificial Intelligence, 34(04), 3195–3202. https://doi.org/10.1609/ aaai.v34i04.5717.
- Kingsford, C., Salzberg, S. What are decision trees?. Nat Biotechnol 26, 1011–1013 (2008). https://doi.org/10.1038/ nbt0908-1011.
- [8] Tanha, J., van Someren, M. & Afsarmanesh, H. Semi-supervised self-training for decision tree classifiers. Int. J. Mach. Learn. & Cyber. 8, 355–370 (2017). https://doi.org/10.1007/s13042-015-0328-7.
- [9] Breiman, L., Friedman, J. H., Olshen, R. A., & Stone, C. J. (1984). Classification And Regression Trees (1st ed.). Routledge. https://doi.org/10.1201/9781315139470.
- [10] R. Rivera-Lopez and J. Canul-Reich. Construction of Near-Optimal Axis-Parallel Decision Trees Using a Differential-Evolution-Based Approach. in IEEE Access, vol. 6, pp. 5548– 5563, 2018, doi: 10.1109/ACCESS.2017.2788700.
- [11] Oksana Svystun, Iryna Yurchak. Recommendation Dialog System for Selecting the Computer Hardware Configuration. Advances in Cyber-Physical Systems, Volume 6, Number 1, 2021, pp. 70–76. ISSN: 2524-0382 (print), 2707-0069 (online) DOI: https://doi.org/10.23939/acps2021.01.070.
- [12] M.Mohri, A.Rostamizadeh, A.Talwalkar. Foundations of Machine Learning, second edition. MIT Press, Second Edition, 2018. Available at: https://doi.org/10.1016/j.jvcir.2019.102705. (Accessed: 18 November 2021).



Iryna Yurchak received B.S. and M.S. degrees at Lviv Polytechnic State University, Lviv, in 1987. She received the Ph.D. degree in automated control systems and advanced information technologies from State Research Institute of Information Infrastructure, Lviv, in 1999.

Her research interests: artificial intelligence systems, intelligent computing systems, neural networks, genetic algorithms, fuzzy logic, recognition systems, prediction problems, computer graphics, computer modeling and animation, web design.