

# DESIGN & DEVELOPMENT OF A BIG DATA ALGORITHM OPTIMIZATION TECHNIQUE FOR A SALES SIMULATION SYSTEM OF A BUSINESS ORGANISATION

**Prof. Dr. G.Manoj Someswar,**  
Dean (R &D), Global Research Academy [Autonomous],  
Hyderabad, Telangana State, India  
[manojgelli@gmail.com](mailto:manojgelli@gmail.com)

**Ganji Vivekanand**  
Assistant Professor, Dept. of MCA, Vaageswari College of Engineering,  
Thimmapur, Karimnagar, Telangana State, India.



## Publication History

Research Article | Open Access

Peer-review: Double-blind Peer-reviewed

Article ID: IJIRAE/RS/Vol.08/Issue05/MYAE10085

Received: 07, May 2021

Accepted: 22, May 2021

Published Online: 08, June 2021

Volume 2021 | Article ID MYAE10085 | <https://doi.org/10.26562/ijirae.2021.v0805.004>

Manoj,Ganji (2021). Design & Development of a Big Data Algorithm Optimization Techniques for a Sales Simulation system of a Business Organisation IJIRAE:: International Journal of Innovative Research in Advanced Engineering, Vol: VIII,107-119

doi: <https://doi.org/10.26562/ijirae.2021.v0805.004>

**Editor-Chief:** Dr.A.Arul Lawrence Selvakumar, Chief Editor, IJIRAE, AM Publications, India

Copyright: ©2021 This is an open access article distributed under the terms of the Creative Commons Attribution License; Which Permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited

**Abstract:** When sales representatives and customers negotiate, it should be confirmed that the ultimate deals can render a high enough professional t for the mercantilism company. Massive corporations have completely different strategies of doing this, one amongst that is to run sales simulations. Such simulation systems typically have to be compelled to perform complicated calculations over massive amounts of information that successively needs economical models and algorithms. This research paper intends to judge whether or not it's potential to optimize Associate in Nursing extend an existing sales system known as per centum, that is presently laid low with intolerably high running times in its simulation method. this can be done through analysis of this implementation, followed by improvement of its models and development of economical algorithms. The performance of those optimized and extended models area unit compared to the present one so as to judge their improvement. The conclusion of this research work is that the simulation method in per centum will so be optimized and extended. The optimized models function as a symptom of thought that shows that results just like the first systems are often calculated inside < 1 Chronicles of the first time period for the most important range of shoppers.

**Keywords:** Scaling Simulation, Article Tree, PCT, Discount Inheritance, Discount Thresholds, Simulation Algorithm

## INTRODUCTION

### MOTIVATION

GMS Software Services Private Limited, Hyderabad, Telangana State, India, is a well established IT business consulting firm specializing in business development of top quality system development and data security. One of our IT company's purchasers would like to boost and extend their current order system. This technique is termed percentage that stands for worth revenue management consumer Tool. Percentage is employed by sales representatives to seek out appropriate discount rates for various things once negotiating with customers. As a central a part of this method, a sales representative can run simulations over totally different discount rates so as to judge their expected professional t. When the system performs a simulation, it starts by creating associate estimation of the customer's future order quantities supported their order history. This estimation is a basis once shrewd the expected marginal professional t that the given discount rates can yield. However, running such simulations takes too while in percentage associated per se an optimization of the simulation method is required.

## **BIG DATA**

Big knowledge could be a slightly abstract phrase that describes the relation between knowledge size and processing speed during a system. An obvious definition of the construct is data whose size forces us to seem on the far side the tried-and-true ways that area unit rife at that point." This implies that a state of affairs wherever innovative improvement of each models and algorithms is needed to handle massive amounts {of knowledge of information} would possibly preferably be classified as an enormous data downside. In PCT, the massive knowledge challenge arises from the large amounts of knowledge required so as to run simulations for giant customers. In some cases quite fifty thousand historical order rows could need to be handled, with multiple doable conditions and discount rates applied to each single one amongst them. Whereas the information set itself isn't very massive by today's standards, the advanced operations and calculations that need to be performed on each of them adds new dimensions to the simulation procedure.[1] Discounts area unit for instance transmissible through an outsized tree structure containing tens of thousands of nodes and also the results should be conferred to the user at intervals an inexpensive quantity of your time. The affordable point in time has been American state defined as 10 seconds for the simulation procedure in proportionality. This worth relies on analysis showing that a system user WHO must wait even more for results of advanced calculations can lose focus - one thing that may prove devastating throughout a negotiation with a client. An ideal simulation procedure would forever come back the results within a number of seconds, since {this would this is able to this may this might this may} mean that simulations could occur throughout traditional speech communication while not requiring any waiting in the least.

## **Goals and limitations**

The first goal of this project is to optimize the present discount simulation algorithmic program so as to scale back its period of time. The discount simulation's purpose is to use given discounts to articles and article classes, so as to judge whether or not they can generate an appropriate professional t for the chosen client. The second goal is to make a model with associated algorithms for a scaling extension of the system's simulation practicality. The aim of this extension is to form it doable to use completely different discount rates looking on the amount of individual orders. This may encourage customers to put some massive orders per annum rather than many little ones, so decreasing shipping and warehouse charges for the corporate while not reducing the sales volumes. In order to realize these goals, this report focuses on 2 doable areas of improvement - optimizations of the models and algorithms themselves and enhancements of the underlying SQL information. Different doable enhancements like hardware upgrades on machines running the algorithmic program, implementations of the algorithmic program in programming languages apart from Java or different information solutions than SQL don't seem to be thought-about.

## **Research Outline**

The rest of this research paper is split into four stages - Simulation, Method, Results and Discussion. The Simulation chapter begins with a close description of however discount rate simulations work and also the issues that this implementation has introduced. The second half contains a specification of the scaling simulation practicality and a proof of the technical difficulties that are introduced by this extension. The Method chapter describes the models and algorithms that are developed during this project. It additionally contains a theoretical analysis of those and comparisons between this implementation in percentage and our answer. In the Results chapter, the performance of percentage also as of our solutions for each the optimized client discount model and also the scaling extension are given. This is separate into a collection of take a look at cases, with motivations of their connotation for actual usage eventualities. The research paper provides an elaborate discussion on the results. This can be wherever the results are mentioned and conclusions and concepts for future work are given.

## **Simulation**

When a sales representative negotiates with a client, one will consider it as a kind of equalization downside. The sales representative desires to maximize the professional t gained by keeping discounts at a minimum, whereas the client desires to reduce his or her prices by maximizing the discounts. This is often wherever the simulation method comes in handy - by simulating the results of latest discounts, it's potential to come to a decision whether or not they area unit professional table enough or not.[2] When each the sales representative and also the client area unit happy with the results, they will save the discounts as conditions within the system's information. Discount rates from such conditions can then be applied to the customer's future orders. This chapter is split into 2 main elements. The first one describes however client discount simulations work and the way these area unit presently enforced in per centum. The second half focuses on a scaling extension to the simulation method that makes it potential to use completely different discount rates looking on individual order volumes. The models and algorithms bestowed during this chapter aren't essentially a twin of those enforced in per centum or the optimized system. They're alleged to be browse as explanations of the expected practicality of Associate in nursing discretionary implementation, unless the rest is expressly declared. 2.1 Customer discount simulation Customer discount simulations area unit presently absolutely enforced in per centum.

By running a simulation over the information represented in this research paper, a sales representative can conclude that professional t would be gained if the client bought a similar articles as within the period however victimization current evaluation conditions. Even a lot of significantly, new discount rates is applied to the simulation that means that the sales representative will see that effects they'll provide and whether or not they appear professional table enough or not. The details of the simulation method area unit described first in this research paper, however reading the chapter within the bestowed order is extremely suggested. Understanding of the under-lying ideas could be a nice advantage once making an attempt to achieve insight into the workings of the simulation method. Data required for a client discount simulation A simulation is predicated on information from the subsequent sources: Article tree - A tree structure wherever branch nodes represent article classes and leaf nodes represent articles Sales history - a collection of mass order rows, containing data concerning previous sales history Existing client conditions - in agreement discount rates from existing contracts, that set a definite discount rate to a specific node within the article tree User input - numerous parameters that specify that historical information and discount rates to use within the simulation Since the contents of those information sources area unit terribly central to the simulation method, a fast description of every one amongst them is bestowed below. The article tree the article tree categorizes all of the company's articles into article teams. These area unit successively classified along into a lot of general classes in 3 price levels", wherever level three is that the most specific and level one is that the most general class. Associate in nursing example tree victimization this structure is bestowed in figure 1.

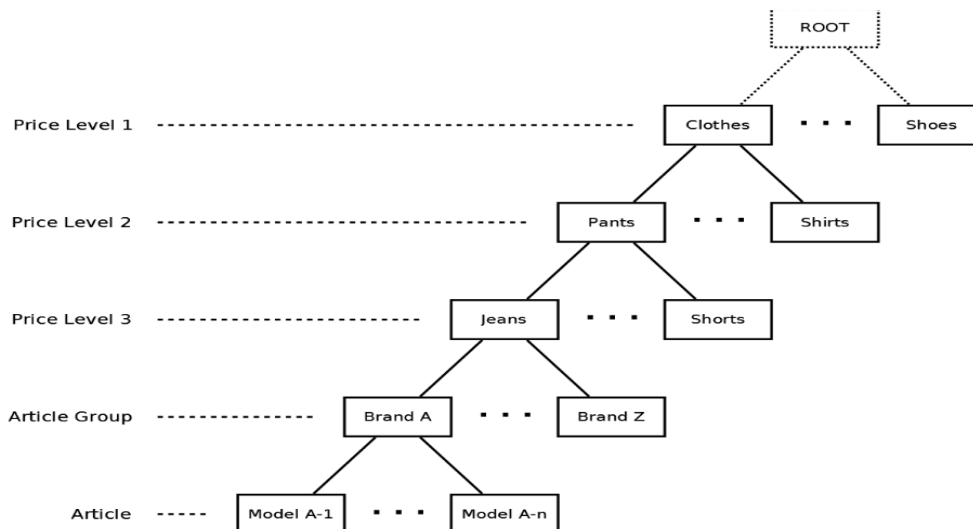


Figure 1: An example article tree containing clothes and accessories

As seen during this figure, leaf nodes contain articles whereas branch nodes represent article classes. The foremost general classes square measure hold on in price index one (Clothes and Shoes), subcategories of those in price index a pair of (Pants and Shirts square measure each subcategories of Clothes) then on. Within the current article information, every price index three node contains specifically one article cluster that means that these 2 levels square measure equally specific. While the tree within the figure is simply associate degree example (using created up names of garments and accessories rather than the lot of cryptically class codes from the important system), it ought to be enough to elucidate the conception of the article tree structure utilized in percentage. The number of nodes for every level within the reduced tree provided for this project is shown in table a pair of.1. The corresponding numbers for the particular systems tree square measure even larger. Since internal company policies don't permit sharing of the complete article tree, this reduced tree has been used throughout the complete project.

Node level	#Nodes
Price level 1	8
Price level 2	64
Price level 3	802
Article Group	802
Article	9,706

Table 1: The amount of nodes for each level of the article tree provided for this project Sales history

The sales history consists of a database containing a large set of historical order information, aggregated on a per month basis. An example aggregated order row is shown in figure 2.

Period	Customer ID	Article no	Actual discount
201207	123456	88084686	257.87
Agreed discount	Avg. target discount	c0	Currency
0	184.76	56.25	EUR
Customer level	Price	Organization	Target discount
8	422.8	1,000	184.76
Value	Volume	Weight unit	Market code
82.465	55.44	KG	DE

Table 2: An example order row from the aggregated sales history database

This example row shows that customer 123456 in the German market bought a total of 55.44 kg of article 88084686 during July 2012. One can also see what the total value of the sold articles was, which discount the customer received and so on. The fields which are relevant for the simulation process are described in greater detail in this research paper. To provide the reader with a perspective of the amount of data stored here, the sales history database for the German market alone stores around 750,000 such aggregated order rows during a single year.

### Existing customer conditions

When a sales representative and a customer agree on a discount rate for a certain article or article category, this is added to a database of customer conditions. An example condition is shown in figure 3.

ID	Opt lock	Aggvalorvol	Channel
abcdef0123456789	0		01
Command	Dirty	Discount	E . stop date
1	FALSE	10.5	
Freeze enddate	Freeze pl date	Freeze start date	New freeze
			FALSE
Note	End date	Start date	Status
This is just an example	201310	201211	
Contract ID	Created by	Customer ID	Price level ID
fedcba0987654321	SalesRep01	123456	DEPL3
Unfreeze cond. ID			

Table 3: An example customer condition

In this example, we will see that the sales representative SalesRep01 has in agreement to grant client 123456 a ten:5% discount on all articles within the class DEPL3 10. alternative fields indicate the ID of the condition, the ID of the contract that the condition belongs to, whether or not or not the condition is briefly disabled (frozen"), AN nonobligatory note specific erectile dysfunction by the sales representative so on. Once again, the fields that square measure relevant for the simulation method square measure delineated in bigger detail in this research paper. User input The final information required for a simulation is provided by the user. This information consists of a client, a path leading from the basis all the way down to discretionary node within the article tree, a period and a collection of discount rates for the nodes within the path.[3]The client is specific erectile dysfunction as a regard to the ID of a client within the client information. Every sales representative includes a set of appointed customers whom he or she will be able to make a choice from. The path is as a collection of hand-picked nodes within the article tree, wherever the first hand-picked node lies on indicant one and any node additional afterward should be a baby of the last hand-picked node. this implies that there's continuously a indicant one node within the path which the sales representative might opt to add a indicant two node further. If a indicant two node was additional, the user will opt to proceed by adding a indicant three node so on. The shortest potential path has the length one (meaning that the trail consists of a indicant one node) and therefore the longest potential path has the length five (containing one node every from worth levels 1-3, a piece of writing cluster node and finally a piece of writing node), that is capable the peak of the article tree. The period is portrayed by a begin date and an finish date, every portrayed as a mixture of a year and a month. Once a simulation is run, historical order information whose amount parameter lies within this interval are going to be used and any information outside of the interval are going to be unnoticed. the top date should after all lie once {the begin the beginning} date and therefore the start date should lie among the last thirteen months. This limit ensures that a full year's history will continuously be used, since the sales history information might not contain the present month's full history nonetheless. Finally, the user can specify a reduction rate for every node within the path. a reduction rate may be a decimal variety between 0:0 and 99:9 with one decimal worth, representing the discount share.



It's additionally potential to let a node inherit its parent's worth by not assignment a reduction rate thereto. Since the value level one node doesn't have a parent node to inherit from, its discount is about to 0:0% if no discount rate is entered on this level. The discount rates square measure generally modified erectile dysfunction multiple times throughout the simulation method, since the sales representative should simulate over multiple configurations so as to find an acceptable set of discount rates to feature to a contract. The simulation method The sales representative starts by coming into that client he's negotiating with and choosing a path within the article tree that discounts are going to be entered. Next up, a begin and stop month is specific erectile dysfunction and currently the system is prepared to run the first simulation. Since no discount rates are entered at this time, all nodes within the path can use their existing discount rates if any such exist within the active conditions and 0:0% otherwise.

All indicant one nodes that don't seem to be affected by the prevailing conditions also will have their discounts set to 0:0%. because of the conception of discount inheritance, all alternative nodes can inherit their parent's discount rate top-down if they are doing not have AN existing condition. this implies that the results of the first run can continuously show the economical results which will follow if identical item quantities square measure sold as within the historical information used for the simulation, taking solely presently active conditions into consideration.

Conditions might are additional or removed since the historical orders were handled, therefore it's not enough to simply mixture the values and professional from the history information. Instead, the "base value" (which one will consider because the worth for the order rows if no discounts had been applied) should be calculated for every article. By applying discount rates from existing conditions to those base values, the system finds out what proportion the client would get to get hold of identical orders if that they had been placed mistreatment current conditions.

In the next step, the sales representative sets discounts for the nodes within the selected path and runs another simulation over a similar information. Any conditions affecting discount rates for the trail nodes are overrun by the discount rates set by the sales representative, whereas conditions affecting alternative nodes can still be taken into thought. The user specific impotence discount rates can then be hereditary down through the article tree rather like those from the conditions. The result can thereby correspond to the professional t which might be achieved if these new rates were other to the conditions information and also the same orders as within the historical information were then placed once more by the client. This simulation step can usually be run multiple times with different discount rates for the nodes within the path, till square measure they're} balanced in such how that each the client and also the sales representative are satisfaction impotence with the results. Running multiple simulations with different discount rates for a similar fundamental quantity and historical information till one gets satisfying results is brought up as inquiring a simulation method. Simulation output So far, the output of simulations has been represented in terms of "pro t" and "value". the particular values computed throughout a simulation square measure after all additional specific than that and per se, the specification of necessities presents pointers for the output information layout. The specification indicates that the output ought to be given as a table, wherever every node within the selected path is delineate as a row. There's additionally a prime row tagged "Total", that shows the overall simulation values of all articles within the whole article tree. A print screen showing however this appearance within the current version of percent is shown in figure 2. The columns of every row square measure represented in table 4.

Customer Total	Volume (kg)	Value (EURO)	C0	C0%	Actual Discount	Agreed Discount				
Total	128,167	471,233	365,257	77.5	86.4					
Price Level 1	Volume (kg)	Value	C0	C0%	Actual Discount	Agreed Discount		Avg. Agreed	Target	Avg. Target
PL1_10	114,635	446,862	355,066	79.5	87.1	44.0	<input type="checkbox"/> * 0.0	0.0	58.5	66.2
Price Level 2	Volume (kg)	Value	C0	C0%	Actual Discount	Agreed Discount		Avg. Agreed	Target	Avg. Target
PL2_01	10,349	24,014	15,670	65.3	87.6	67.1	<input type="checkbox"/> * 0.0	0.0	66.9	66.9
Price Level 3	Volume (kg)	Value	C0	C0%	Actual Discount	Agreed Discount		Avg. Agreed	Target	Avg. Target
PL3_5751F1	26	187	161	86.0	78.1	12.4	<input type="checkbox"/> * 0.0	0.0	54.7	54.7

Figure 2: A print screen from PCT showing how simulation output is presented in the current system

Field name	Unit	Type	Description
Node name	n/a	String	The name of the row's node
Volume	kg	Integer	The total volume of all orders for articles under the row's node in the article tree
Value	Euro	Integer	The total amount of money which the customer would have to pay if all historical orders for articles under the row's node were placed again, with the new discounts applied
CO	Euro	Integer	The profit which the company would gain if all historical orders for articles under the row's node were placed again, with the new discounts applied
CO%	Percent	Decimal number	Shows how many percent of the row's value CO corresponds to, i.e. $(Value/CO)*100$
Actual discount	Percent	Decimal number	The average historical discount for articles under the node in the simulation period
Above target	n/a	Boolean	A warning flag which shows whether the agreed discount is higher than the node's target discount
Agreed discount	Percent	Decimal number	The discount used for the row's node in the current simulation
Avg Agreed	Percent	Decimal number	The average agreed discount for articles under the row's node in the article tree
Target discount	Percent	Decimal number	A recommended target discount for the node, based on the customer's pricing level
Avg Target	Percent	Decimal number	The average historical target discount of articles under the row's node

Table 4: The columns which are used to structure the output from a simulation

The five last columns are empty for the Total" row, since these values are considered irrelevant to display for the whole article tree.

### Discount inheritance

Discounts can be applied to nodes on any level of the article tree - from price level 1 down to specific articles. It is intuitive that a discount which is set for a single article will only affect the price of that specific article. When it comes to discounts set on article groups or price level nodes, the system uses a concept called "discount inheritance" to let this affect underlying nodes. In order to determine which discount rate to apply to a given node, the method presented in algorithm 2.1.1 is used.

---

#### Algorithm 2.1.1: Find Discount Rate(Node n)

---

```

Input: A node n from the article tree
Result: The discount rate which should be applied to n
if n is a node in the path for which a discount rate d is set then
    return d
else if n is not a node in the path AND n has an active condition c then
    return the discount rate from condition c
else if n is a price level 1 node then
    return 0:0%
else
    parent := n's parent node in the article tree
    return find Discount Rate(parent)
end
    
```

The concept of discount inheritance is easy to visualize due to the tree structure of the article database. An example tree with some existing discount rates is shown in figure 3. Existing discount rates are written directly onto the grey nodes to which they belong, while nodes without such rates are white. The final result of the discount rate inheritance in the same tree can be seen in figure 4, where arrows show how discount rates are passed down through the tree.

### Current implementation

As mentioned within the project motivation, this implementation of percentage suffers from vital performance problems. Since the ASCII text file of this technique isn't allowed to be enclosed during this report, the issues of its formula ought to be explained in terms of unhealthy structure decisions and complexness instead of examples and excerpts from the particular code.

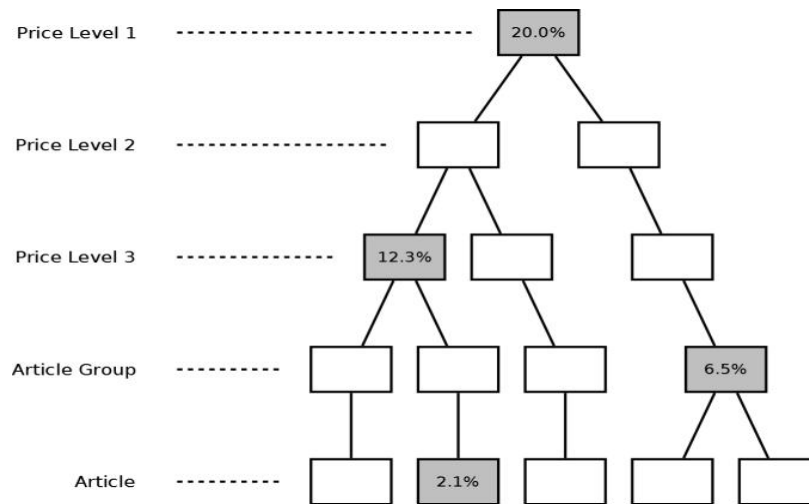


Figure 3: An example article tree where discount rates have been set for four nodes

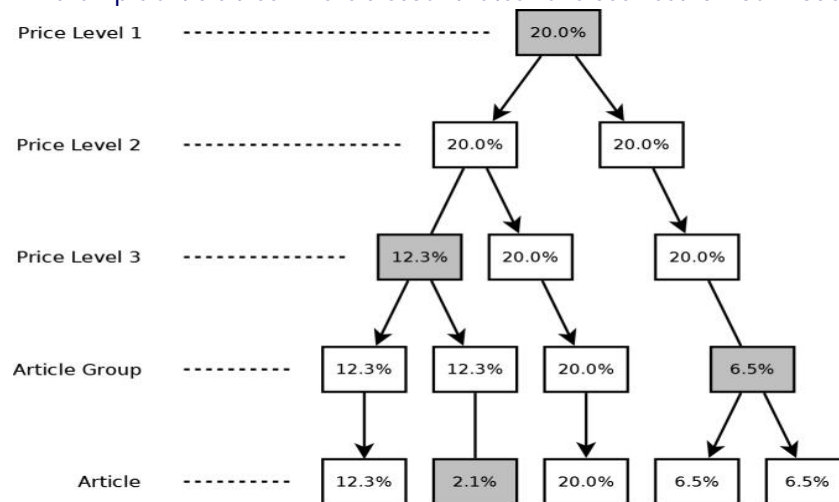


Figure 4: Discounted inheritance

A (very) rough define of the formula structure wont to perform simulations in percentage is conferred in formula two.1.2. whereas it doesn't inspire or make a case for the small print of every step, it will give enough info to analyze its complexness. to provide the reader some type of plan of the particular magnitude of the implementation of this formula, its Java ASCII text file takes up many hundred kilobytes (not together with GUI, server connections, info handling and different elements that don't seem to be directly associated with the algorithm). In different words, a line describing e.g. criteria matching means that running a separate algorithm that successively includes a complexness value mentioning.

Algorithm 2.1.2: Structure of the simulation process in PCT

```

if this is the rst run of the simulation process then
initialize connection to each input data element in the GUI [O(k)]
end
foreach price level in the article tree [O(k)] do
match condition level [O(k)]
match price level [O(k)]
foreach item in the customer's cache [O(n)] do
match criteria [O(k)]
end
retrieve target discount [O(k)]
foreach article in the article tree [O(a)] do
foreach article in the customer's cache [O(n)] do
match criteria [O(k)]
end
end
foreach price level in the article tree [O(k)] do
retrieve data and calculate results

```

```

end
end
retrieve agreed discounts [O(k)]
compare discounts to target discounts [O(k)]
end
end
foreach article in the customer's cache [O(n)] do
  calculate results for articles under price level 1 nodes 2= path
end
  
```

In the pseudo code above, the complexity has been included on each line where O notation is applicable. The meaning of each occurring variable in the O notation is presented in table 5.

Variable	Magnitude	Description
a	30,000	The amount of articles in the article tree (exact numbers for the reduced tree used in this project can be seen in table 2.1)
k	5	Height of the article tree (constant in this program, but may vary between implementations)
n	1 n 13a	Distinct (month, article) tuples in the selected customer's sales history for the last 13 months

Table 5: The meaning of different complexity variables in algorithm 2.1.2

The total complexity of the implementation of the current simulation algorithm is  $O(k+k(k+k+nk+k+a(n(k+k)))+k+k)+n) = O(k+5k^2 + nk^2 + 2ank^2 + n) = O(ank^2)$ . It should also be noted that the complexity of repeated runs of the algorithm is  $O(k(k+k+nk+k+a(n(k+k)))+k+k)+n) = O(5k^2 + nk^2 + 2ank^2 + n) = O(ank^2)$ . This is barely associate improvement from the first run in any respect - one single k term is removed since the format tread line a pair of doesn't got to be run once more. The entire complexness of  $O(ank^2)$  is high in itself, since each a and n will hold quite giant numbers and also the k term is employed at multiple places. However, this can be not the sole reason behind the high running times of the algorithmic program. Another massive downside is that the on-demand usage of information resources. Whenever a group of values from the information is required, a brand new association to the information is opened. [4] The sought-after values square measure then retrieved by associate SQL question and later on the information association is closed once more. Repeatedly gap and shutting information connections takes time and this can be drained several elements of the algorithmic program, together with the info retrieval strategies mentioned in line fifteen. this implies that  $O(ank^2)$  information connections and SQL queries might need to be opened and run within the worst case. Some values square measure even retrieved from the information multiple times throughout one execution, since they're utilized in multiple places within the code however aren't saved once being retrieved the first time. It ought to but be noted that some actions are taken so as to cut back the number of knowledge retrieved from the information per simulation. each client's order history for the last year is hold on as an inventory in a very hash map indexed by the customer ID, that thereby works as associate in-memory information. This makes retrieval of a customer's historical data (without direct information base access) in  $O(1)$  time doable. Of course, iterating over the ensuing list can still take  $O(n)$  time. This cache is formed on server startup and updated frequently, thus creation and updates of the cache don't a shock therapy the time period of the simulation algorithmic program. Some loops in percentage square measure still performed over all distinct values in sure information tables once data from this cache may are used instead, resulting in even a lot of extra information lookups. A third reason for the high time period is that the redundant calculation of sure values. The for each loop on lines 4-21 within the algorithmic program on top of runs once for every index number and calculates the results for all articles underneath the node on it level. this implies that the results for all articles underneath the value level one node are calculated first, followed by a computation of all articles underneath the value level a pair of node then on for a complete of up to k calculations of an equivalent values for a few articles. Let the chosen path within the simulation be known as P and also the path from associate arbitrary T article a up to its index number one root be known as Pa. Then,  $jP Paj$  (the quantity of nodes that square measure each in P and in Pa) equals the amount of times article ad's results needs to be calculated throughout every simulation. shrewd an equivalent price quite once is in fact redundant and adds extra time period. this can be visualized in figure a pair of.5. Nodes within the path P square measure marked with thick outlines within the figure. Values for articles inside the blue box (A1 A5) square measure calculated once, whereas articles inside the inexperienced box (A1 A4) square measure calculated another time and articles inside the red box (A2:A3) yet one more time. The purple box (A6) marks articles that lie underneath a index number one node  $2= P$ , whose values square measure continuously calculated just the once.



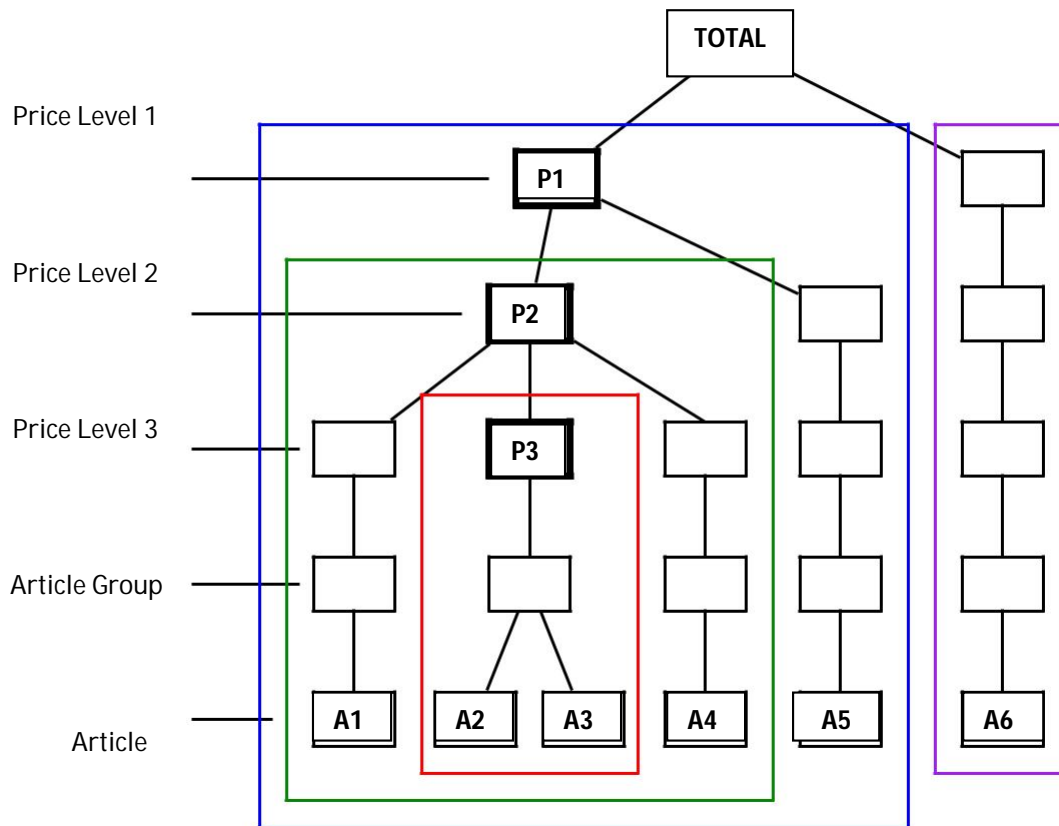


Figure 5: Redundant calculation of values in PCT

The combination of a high time complexity, inefficient database usage and redundant calculations cause the running time of each simulation to grow rapidly for increasing values of n.

### Scaling simulation

As mentioned in section 1.1, the company wants a scaling extension of the simulation process and conditions in PCT. The purpose of this extension is to make it possible to apply different discount rates depending on the volume of each individual order. If larger order volumes are rewarded with higher discounts, customers will be more likely to place large orders a few times per year instead of small orders every week or month, thus decreasing shipping and warehouse charges. The extension's specification is built around a concept called discount stairs". These are set by sales representatives on a per customer and node basis, in order to define which discount rates will be applied to orders of certain volumes. This concept is described in greater detail in this research paper. The data sources which are needed for scaling simulations are in turn described in this research paper.

### Data needed for a scaling simulation

A scaling simulation is based on data from the following sources:

**Article tree** - A tree structure where branch nodes represent "price levels" (article categories) and leaf nodes represent articles

**Sales history** - A set of order rows, containing information about previous sales history. The aggregated sales history from the customer discount model is also required.

**Existing scaling conditions** - Agreed scaling conditions from existing contracts, which set a certain discount stair to an article, article group or price level 3 node in the article tree

**User input** - Various parameters that specify which historical data to use, which discount stair to use and various other simulation settings

### The article tree

This is exactly the same tree as the one used for customer discount simulations, which is described in this research paper.

### Sales history

The sales history used for scaling simulations consists of a database containing a large set of historical orders. Note that these order rows are not aggregated, as opposed to the ones used for customer discount simulations. An example order row is shown in table 6.

Invoice date	Invoice no	Invoice value	Invoice volume
2012-07-03	a138215	40.345	35
Article ID	Customer ID	Unit ID	ID
DE88001104	123456	kg	379311

Table 6: An example order row from the order history database

However, in some situations data from the aggregated database table is used as well. This means that scaling simulations require both of these database tables to be present.

**Existing scaling conditions**

Agreed discount stairs for specific article tree nodes and customers are added to a database of scaling conditions. An example scaling condition is show in table 7.

Price level ID	Customer ID	v0	v1	v2	v3	v4
DEPL35690F 1	123456	5	10	20	30	50
v5	d0	d1	d2	d3	d4	d5
	12.1	15.0	17.5	20.0	25.3	

Table 7: An example scaling condition

In this example, the scaling condition covers the worth level three node DEP L3 5690F one for client 123456. The step has five thresholds, whose volume limits square measure shown in columns v0 v4 and their several discounts in columns d0 d4. Since there's no sixth threshold during this example, v5 and d5 square measure left empty. an outline of however these values square measure utilized in a scaling simulation is conferred in section a pair of.2.3. According to the scaling extension specification, all customers may be expected to own a complete of at the most 10 active scaling conditions. Most discounts square measure in different words still expected to be handled through client discount conditions once the scaling extension has been enforced. User input The final knowledge required for a scaling simulation is provided by the user.[5] This knowledge consists of a client, a period, a node and a reduction step. The client could be a respect to the ID of a client within the client information, similar to in client discount simulations. The period is depicted by a begin date associate degreed an finish date, that is slightly different from the period in client discount simulations. Since the historical order rows used for scaling simulations don't seem to be mass per month just like the standard order history, these dates specify every day of the month similarly.

The node could be a respect to either a commentary, a commentary cluster or a price index three node within the article tree. this can be quite different from the trail utilized in client discount simulations - not solely as a result of solely one node is chosen, however additionally as a result of price index one and price index a pair of nodes cannot be used. A discount step could be a approach of Diamond State fining different discount rates for an equivalent node, looking on individual order volumes. this can be delineated in bigger detail in this research paper. Discount stairs As mentioned in section a pair of.2.1, discount stairs create it attainable to use different discount rates for orders looking on their volumes. similar to in client discount simulations, discount inheritance is applied. However, scaling conditions cannot be set for price index one or price index a pair of nodes within the article tree. the rationale behind this can be that articles United Nations agency don't share an equivalent relation on price index three square measure usually thought of too various to share volume limits. In different words, it might not perpetually be to use an equivalent discount rate for i.e. orders between twenty and twenty five kilogram on 2 articles of terribly different varieties.[6] A discount step consists of between one and 6 volume thresholds and a reduction rate for every one among these. The thresholds indicate boundaries between weight intervals, that means that i thresholds Diamond State ne i+1 intervals. a group of such thresholds is shown in table a pair of.8 and its ensuing interval limits square measure shown in table a pair of.9. Since there square measure six thresholds during this example, there square measure seven discount intervals

Threshold volume	Threshold discount (%)
5	12.1
10	15.0
20	17.5
30	20.0
50	25.3
70	30.1

Table 8: An example set of discount thresholds

Volume v	Discount (%)
v < 5	0.0
5 v < 10	12.1

10 v < 20	15.0
20 v < 30	17.5
30 v < 50	20.0
50 v < 70	25.3
70 V	30.1

Table 9: The resulting discount intervals from the thresholds in table 8

From the last table we can easily see that an order with e.g. volume  $v = 4$  would get 0:0% discount if this stair is used, while an order with volume  $v = 27$  would receive a 17:5% discount.

### The scaling simulation process

The general workings of scaling simulations area unit the same as those of client discount simulations, however some differences area unit in fact gift. A scaling simulation begins with a sales representative choosing a client. In more than this, a indicant three node, article cluster node or article node within the article tree is chosen and a amount specification impotence.[7] Moreover, a group of 1 to 6 volume thresholds is specification impotence. Note that solely volumes area unit entered before the first run - actual discount rates for these don't seem to be entered till later. The system is currently able to run the first scaling simulation. During the first run, the chosen node  $n$  can use the discount rate 0:0% for every volume interval specification impotence within the input step. when every run, the user will modify the discount rates for every interval of no's discount support, excluding very cheap one that is often barred to 0:0%. The values for all articles area unit calculated in line with the strategy represented in formula.

---

#### Algorithm 2.2.1: getValues(Node $n_a$ )

---

```

Input: An article node  $n_a$  from the article tree
Result: Total simulated cost, value and volume for article  $n_a$ 
current :=  $n_a$ 
repeat
  if current is the selected simulation node with discount stair  $d_s$  then
    return scalingSimulation( $n_a$ ,  $d_s$ )
  else if current has an existing discount stair  $d_a$  in a scaling condition then
    return scalingSimulation( $n_a$ ,  $d_a$ )
  end
  current := current's parent in the article tree
until current is a price level 3 node
return Total cost, value and volume from aggregated (i.e. not scaling) history database for  $n_a$ 
during selected time period
  
```

As line ten within the algorithmic program on top of shows, values for articles that don't seem to be affected by the scaling node or scaling conditions are retrieved directly from the aggregative information table (which is additionally used for client discount simulations). This works since the aggregated knowledge for a month per Delaware definition equals the total of all individual orders from identical month. Even if solely a neighborhood of the month is roofed by the chosen interval for the simulation, the complete month's history can still be employed in this case. It ought to even be noted that the calculation of base costs and application of client discounts are unnoticed - the aggregative values are used directly so as to lower the scaling simulation's complexity.[8] Algorithm 2.2.2 shows however the scaling Simulation () perform that is named in algorithm two.2.1 works. Since the specification of historical order rows doesn't embody any columns for worth and value, these values have to be compelled to be calculated from the aggregate historical knowledge. First o, the article's \list price" and \list cost" square measure calculated as a kind of base values on the shape currency unit/volume unit (e.g. Euro/kg) for the specified article and month.

---

#### Algorithm 2.2.2: scalingSimulation(Node $n_a$ , DiscountStair $d_s$ )

---

```

Input: An article node  $n_a$  from the article tree and a discount stair  $d_s$ 
Result: Simulated values (cost, value and volume) for article  $n_a$ 
totalCost := 0
totalValue := 0
totalVolume := 0
orderRows := all order rows from the scaling sales history (see section 2.2.1) for article  $n_a$  within
selected time period
foreach row  $r$  in orderRows do
  rowVolume :=  $r$ 's volume
  agrPrice := price for  $n_a$  in  $r$ 's month in the aggregated history database
  
```

```

aggrCost := cost for  $n_a$  in  $r$ 's month in the aggregated history database
aggrVolume := volume for  $n_a$  in  $r$ 's month in the aggregated history database
listPrice := aggrPrice / aggrVolume
listCost := aggrCost / aggrVolume
rowDiscount := The discount from  $d_s$  whose volume interval covers rowVolume
rowCost := rowVolume * listCost
rowValue := rowVolume * listPrice * (1 - rowDiscount*0.01)
totalCost := totalCost + rowCost
totalValue := totalValue + rowValue
totalVolume := totalVolume + rowVolume
end
return (totalCost, totalValue, totalVolume)

```

These values square measure then increased by this order row's volume so as to urge its worth and value.[9] Next up, the row's simulated worth is calculated. The row's volume is matched to a volume interval within the discount support and also the corresponding discount is applied to the row's worth so as to find its worth. Finally, the total of all row's prices, values and volumes square measure came. The simulation results square measure obtained by aggregating the ensuing values for all articles, together with those wherever values square measure retrieved from the aggregate historical information. Scaling simulation output the output of a scaling simulation ought to be conferred a bit like the output of a client discount simulation that is delineated in section a pair of.1.2. All-time low row covers the scaling node, all of its ascendant nodes have a row every within the middle and also the prime row contains the full values for the entire article. Technical difficulties Scaling simulations haven't however been enforced in percentage. Adding this practicality has been thought of impractical, since scaling simulations run over way larger knowledge sets than client discount simulations (which have already got issues with high running times). The non-aggregated knowledge is even large to be control in associate in-memory information cache, which suggests that each historical order row can have to be compelled to be retrieved from normal information. This slows down the information handling even additional. Customers square measure calculable to put orders for an equivalent article up to once per week and also the fundamental measure used for a scaling simulation are often at the most one year. this suggests that one will assume a most of fifty two order rows per article in a very single scaling simulation. As shown in table a pair of.1, the reduced article tree contains 9;706 distinct articles whereas the complete article tree has around 30;000 nodes. it's deemed potential that one client buys up to 1000 different articles frequently. Scaling simulations will solely be run over article nodes or price index three nodes within the article tree, however if scaling conditions square measure specified for different such nodes than the chosen scaling node, these need separate scaling simulations on their own. The time period of a scaling simulation over associate discretional node for a client can thereby increase for each scaling condition other for an equivalent client. price index three nodes have a mean of twelve.3 and a most of 172 underlying articles, thus adding one condition may increase the quantity of needed historical order rows by virtually 9;000. This means that running a scaling simulation will need multiple separate scaling simulations for articles with existing scaling conditions. In total, these may need computations over as several as 52;000 historical order rows.[10]

## RESULTS & CONCLUSION

The purpose of this research work is to analyze the project's results and to debate however our resolution will improve the simulation method in proportionality. The first 2 sections discuss the results from chapter four and the way these correlate with our optimizations. this can be followed by some ideas regarding doable future work and a outline of the conclusions from this project. 5.1 Customer discount simulation As specific function in section two.1.4, the long running times of client discount simulations in proportionality square measure caused by 3 major issues. These square measure this algorithm's time complexness, its inefficient usage of info resources and also the redundant calculations caused by PCT's inefficient model. This section aims to elucidate however these issues square measure handled in our model and the way this will be seen within the results. Time complexness The enhancements of the time complexness square measure simple to spot thanks to the complexness analysis of proportionality in section two.1.4 and of our model in section three.1.3. PCT's complexness of  $O(nk^2)$  ought to be compared to our model's  $O(nk)$ , wherever  $n$  is proscribed by  $a$  and  $u$  in any realistic situation. this means that the calculations in our model square measure performed in an exceedingly much more efficient method. The actual decrease of the period caused by the reduced complexness is after all gift all told simulations, however it's terribly simple to check in check cases one and 3. Within the latter of those, the period for the most important information set victimization our model was roughly 99:4% less than the one for proportionality. If we tend to compare the running times for the tiniest information sets in test suit one instead, we are able to see that the advance remains 82:0%. It is thus affordable to assume that simulations run victimization our model can generally be > eightieth quicker than a similar simulations run victimization proportionality which the relative improvement can get even larger because the underlying simulation information grows.

Database usage As delineated in section two.1.4, proportionality suffers from AN ineffective technique of retrieving information from the underlying info. a similar information is usually retrieved multiple times throughout one simulation and loops square measure typically run over info rows rather than the cache entry for the client. within the worst case situation, proportionality should run  $O(nk^2)$  SQL queries. When developing and implementing our model, we've created absolute to avoid these issues by saving information that should be used multiple times and by victimization values from the in-memory cache as usually as doable. All active conditions for the chosen client square measure retrieved employing a single question and hold on in an exceedingly hash map victimization the condition nodes as keys. Any succeeding checks against the user's conditions will then be performed in  $O(1)$  time, while not accessing the info in any respect. This implies that the exaggerated simulation time caused by increasing numbers of conditions square measure barely noticeable - one thing that is so reflected within the results of test. Test case five additionally shows that the quantity of conditions in proportionality create a negligible impact on the running time; a minimum of in realistic situations wherever customers have so much but 100 conditions. The rationale behind this can be that alternative (about equally slow) calculations square measure performed instead for articles while not existing conditions. The removal of the article cluster nodes from the first model weakened each the article tree's height and also the most path length from five rights down to four. As a result of this, iterations over these parameters run quicker victimization our model notwithstanding the particular code would otherwise look identical. Furthermore, our model ne'er iterates over a listing of articles retrieved directly from the info. Such loops square measure instead run over the key set of another hash map, whose keys represent solely the articles that the client has bought throughout the specification function period of time. This implies that fewer info queries square measure required, and that supererogatory iterations square measure unheeded that successively improves our solution's complexness moreover. The total quantity of info queries per simulation in our model is one for the conditions and an extra question for each node within the path to retrieve their various target discounts, for a complete of two five queries. The difference between this low range and also the  $O(nk^2)$  queries for the worst case situation in proportionality mentioned higher than is clearly terribly huge.

#### REFERENCES

1. Informatica and Capgemini, The Big Data Payoff: Turning Big Data into Business Value, 2016.
2. V. Kayser, B. Nehrke, and D. Zubovic, Data Science as an Innovation Challenge: From Big Data to Value Proposition, Technology Innovation Management Review, 2018.
3. M. S. Hopkins and R. Shockley, Big Data, Analytics and the Path from Insights to Value, MITS Ioan Management Review, 2011.
4. Harvard Business Review, The Enterprise Lacks a Big Data Strategy for IoT Transformation, 2017, pp.1-12.
5. S. Lavallo, M. S. Hopkins, E. Lesser, R. Schockley, and N. Krushchwitz, Analytics: The New Path to Value, MITS Ioan Management Rev., 2010.
6. S. Viaene and A. Van den Bunder, The secrets to managing business analytics projects, MITS Ioan Manag. Rev., 2011, pp. 65-69.
7. A. Chebotko, A. Kashlev, and S. Lu, "A Big Data Modeling Methodology for Apache Cassandra," Proc. of IEEE Int. Congress on Big Data, 2015, pp.238- p.245.
8. A. Fink, R. Guzzo, and S. Roberts, "Big Data at Work: Lessons from the Field," Society for Industrial and Oranizational Pshchology, 2017.
9. S. Nalchigar and E. Yu, "Business-driven data analytics: a conceptual modeling framework", Data & Knowledge Engineering, 2018, pp. 1-14.
10. M. A. Berry and G. S. Linoff, Mastering data mining: the art and science of customer relationship management, Industrial management data system, 2000.