

Data Generator for SAP Solutions Using Benerator

Baba Pakruddin Tailor (baba.tailor@st.ovgu.de)
Julian Reddy Allam (julian.allam@st.ovgu.de)
Nagasai Krishna Ayinampudi (naga.ayinampudi@st.ovgu.de)
Sai Rajesh Vanimireddy (vanimireddy.sairajesh@st.ovgu.de)

*Otto-Von-Guericke University, Magdeburg
SAP UCC Magdeburg, Department of Informatics*

Abstract

Data analysis has proven useful in the present scenario for businesses when it comes to improving their business processes. Ideally, complete access to data is important for the production of high-quality and accurate analyses. However, real-world data is often subject to several privacy constraints, making them unavailable at full-granularity for some analyses. This can significantly impact the quality of the analysis based on it. Under these constraints, researchers often resort to generating realistic synthetic data to verify the efficacy of various methods. It is critical that the generated realistic data needs to be 'realistic' (i.e., follow the same distributions / trends found in the real world data that it attempts to approximate) for the analysis based on it to be useful. This project focuses on the use of 'Benerator', an open-source tool for the efficient generation of large amounts of realistic data for various SAP solutions. Different analyses are performed on this data to verify that it approximates distributions found in real-world data. .

Keywords: Data, privacy constraints, Realistic data, Analyses

1. Introduction

SAP is a German multinational software corporation that creates enterprise software that enables companies to manage their business processes. Since SAP solutions cater to a variety of business processes, the training of personnel to properly leverage the full functionality offered by these solutions is a critical step to realizing business outcomes. To ease this training process,

SAP UCC created 'Global Bike INC.', a fictitious multinational company that manufactures and sells bicycles. Along with the training material (Case studies, Data sheets, Slides, Sample Exercises, Lecture Notes, etc.) SAP also creates the master data required for the efficient use of this training material. However, to leverage the full capabilities of SAP systems, this data needs to be generated efficiently.

Generating a large amount of data is useful for research, analysis and pattern recognition. Data generators are usually used when actual data is challenging to obtain due to security and / or privacy concerns[1]. In these situations, realistic synthetic data from data generators are used as a substitute for actual data. In addition to availability at any desired granularity, realistic data has the advantage that it can be generated to test events or conditions (supply chain issues, policy changes, etc.) that may not occur in real data[1]. It is obvious that the generated data should approximate the trends and patterns in the actual data as closely as possible. Benerator was chosen as the tool because of its support for realistic data generation and for connecting with various databases.

It is the aim of this project to customize 'Benerator' to generate realistic master data for SAP solutions. Benerator is an open source tool capable of efficiently generating large data. In addition, Benerator is domain-independent, allowing the creation of data for use in diverse fields. The resultant data from Benerator tool was analyzed for seasonal and regional sales trends to ensure that the generated data is indeed realistic. It has already been shown in that properly tuning the distributions of the functions generating the data ensures that the data is realistic[1]. It is seen that this method can accurately generate data that shows realistic variations in sales volumes in Germany and the USA. In addition, the generated data is also shown to capture variations in sales due to seasonal and regional and geographic effects.

The following section discusses the details of the tools used, their configuration and the implementational details of generating large, realistic datasets. Section 3 details some graphical analyses performed on the generated data sets to confirm that it is realistic. The last section summarizes our findings and briefly discusses possible improvements.

2. Tools and Methodology

This section begins by briefly describing about each tool and its function within the project, followed by their set-up and/or configuration.

Generater is an open source tool used for fast generation of large amounts of realistic data. It provides a large set of plug-in interfaces for mappings and customizable extensions and configuration options. This tool supports CSV, Flat Files, XML and also helps in extraction and anonymization of real production data. In this project, we have generated realistic datasets using weighted distribution function and further used in various analysis.

MySQL is a relational database management system (RDBMS) based on Structured Query language. It is used for accessing, managing and adding content in a database. MySQL connections panel enables data architects to easily manage database connections. This database is used to define a user and subsequently create customer, product, vendor, material and master data tables for USA and Germany and then insert the values into them.

Apache Maven is a comprehensive tool and software project management system. Based on the project object model (POM) file Maven tool can manage to report, makes building projects easy, transparent migration to new features and documentation from a central part of information. The POM file manages dependencies and provides the entire configuration for a project. You can also compose individual phases of the build process which are further implemented as plugins.

2.1. Data Generation overview

As already discussed, access to real data may be restricted in some situations. To generate realistic sales data, the generated data must approximate seasonal, regional and product preference trends in the real world. To find these trends, real data was collected related to

- Sales of bikes and accessories
- Number of bike commuters in different cities / countries,
- Monthly bike accidents and
- Commuter transportation preferences

The data from these sources were analyzed for patterns. By considering these patterns and distributions, we have customized the Benerator to generate sales data for the period 2009-2011.

In this paper, we have followed four rules to generate realistic synthetic data.

- Region-wise customer distribution
- Month-wise customer distribution
- Year-wise distribution
- Product preferences

It was seen from [2] that the sales of bicycles are subject to regional effects. Regional factors like population, geography, general differences in population percentage that uses bikes to commute, and infrastructure contribute to the differences in bikes usage and sales figures.

In order to approximate this, the customers that belong to specific regions in Germany and USA were given higher weights, as can be seen in table 2 and table 4. The weights file is designed as a CSV, with the first column identifying the customer and the second column assigning the weight. Benerator automatically normalizes all weights.

Since accurate sales data was not readily available at a monthly level, the percentage of bicycle accidents per month has used a proxy for total bicycle usage [2]. It is assumed that the usage data and sales data are well correlated. The data follows intuitive bicycle usage patterns that would be expected from prevailing weather conditions. In order to generate sales data that follows the same trends (max. sales in summer months), the Benerator descriptor file was modified to assign a fixed proportion of total yearly orders to each month.

Sales data was generated with peak sales in 2011, to accurately reflect the sales as described in [3]. The total sales for the years are controlled by controlling the sales for the individual months comprising it. It is also to be noted that the total sales for Germany are much higher than that for USA.

In addition to the three rules described above, the sales data is also created to reflect real-world customer preferences. The following rules were applied to the data to control product-level purchase preferences:

- Gender distribution: Higher sales products aimed at men [4]

- Color Preferences: Highest sales for black products, followed by silver and red [5][6]
- Product Type: Higher sales for bikes than for accessories shown in figures 13 and 14

To implement these rules, the benerator is customized to generate sales orders where some products are picked with a higher probability ('weight') than others.

Figure 1 shows an overview of the components involved in the Data generation process. The main steps that need to be taken are: configuration of descriptor file, defining and populating the database schema and defining the CSV files for customizing the output. The customer, product and vendor data are stored as tables in the database. All the customers and products are then assigned with their respective weights and stored as additional data in CSV files. The Benerator is configured to populate two tables, *sales – order* and *order – details*.

We Configure the descriptor file to create two entities with their attributes, specifying number of rows to be generated in the entities. The Configuration of descriptor file includes querying the stored tables (for customer, vendor, and product details) and CSV files to access data and assign a distribution to be followed in the dataset.

2.2. Realistic data Generation process

The Benerator is customized to generate data that follows the rules described earlier in the section. This is done through modifying the *benerator.xml* file. The benerator is configured to use weights for the customers and products from the predefined in CSV file. The output of the Benerator is directed to the *SalesOrder* and *OrderDetails*. The attributes for which data is generated are *orderId*, *CustomerId*, *Sale – Time* (*SalesOrder* table), *orderdetailsId*, *productId* and *count* (*OrderDetailsId* table, along with foreign key *orderId*). The benerator is configured to automatically fetch customer, product and vendor metadata from their respective tables in the database. It is to be noted that supporting data for generation of purchase orders already exists in the database and that the benerator can be easily configured to support its generation. The number of sales orders to create is also defined in the file. The output of the benerator is directed to the

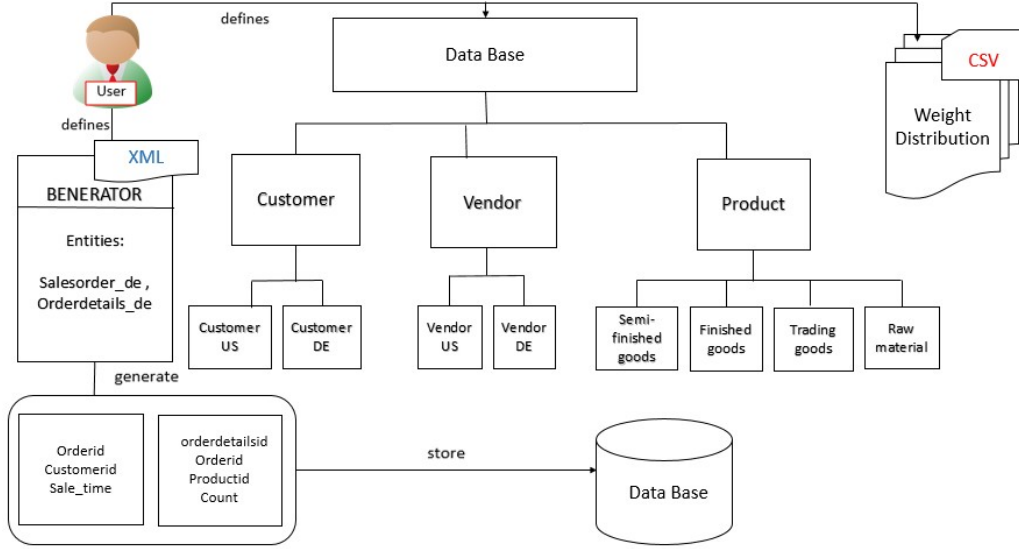


Figure 1: Data Generation Process Overview

database tables discussed above and is extracted as CSV from MySQL for analyses described in the following section.

2.3. Data Schema

The data model follows the star schema. Figure 2 shows customer, Product, Sales order and Time entities and their attributes. The primary and foreign keys are marked 'PK' and 'FK'.

2.3.1. Descriptor file

As discussed already, the *benerator.xml* file or 'descriptor file' is the core component of the project. It defines the database URL, the entities and attributes that are used, and the tables that receive the output. The file also contains customization that ensures that the generated data follows the desired distributions (see section 2.1). It is to be noted that each sales order also consists of a *count* attribute which specifies how many items of a particular product is purchased. The count attribute takes a random value arbitrarily fixed between 1 and 20.

2.3.2. Data Elements

We used 'attribute', 'variable', 'beans', 'reference' and 'generate' tags to generate the data.

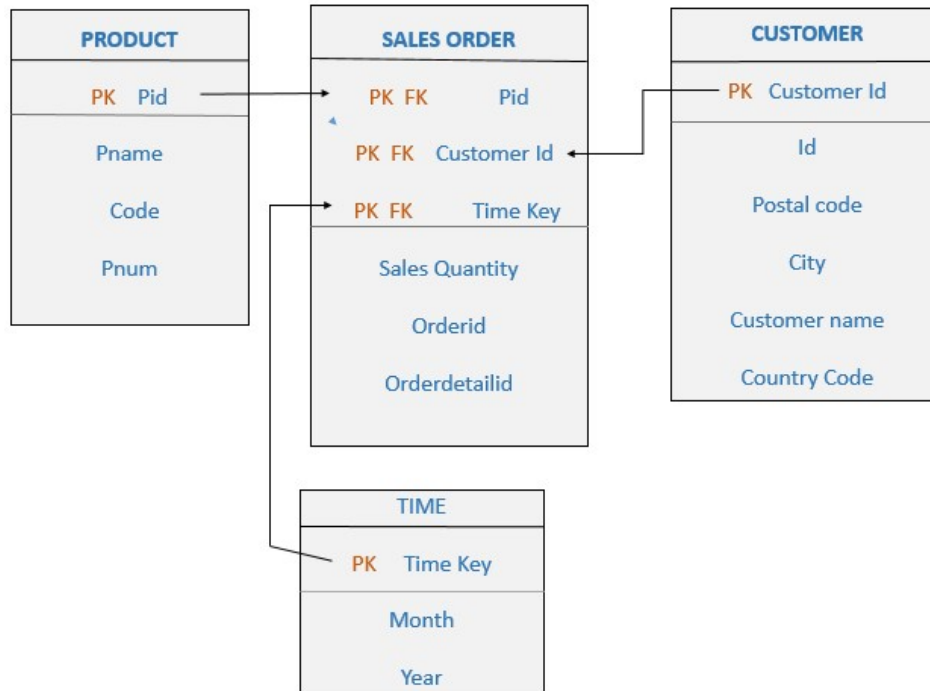


Figure 2: Database Design

- Attribute Name: Refers to columns inside an entity.
Syntax: `< attributename = "sale-time" type = "datetime" generator = "dtGen" / >`
- Variable: Used to assign fixed weights to different products and customers.
Syntax: `< variablename = "weightings" source = "weightings.wgt.csv" distribution = "weighted" / >`
- Beans: Classes in Benerator that provide existing functionality
Syntax: `< beanid = "dtGen" class = "DateTimeGenerator" / >`
- Reference: Used to refer to existing entities in database. (Uses SQL code)
Syntax: `< referencename = "customerid" type = "int" targetType =`

"salesorders-de" source = "db" selector = "selectidfromcustomer-de" nullable = "false" cyclic = "true" script = "weightings" / >

- Generate: The main tag that initiates the data generation process. It encloses all the entities and attributes for which data is generated.
Syntax: < generatename = "salesorders - de" type = "salesorders-de" count = "100" consumer = "db, ConsoleExporter" / >

2.3.3. Weight distribution function

The weight function allows fine-grained control over the selection proportions of the available values for a particular attribute. For instance, if an attribute A has five values p, q, r, s and t , and each of these values has is configured with weights (1, 1, 1, 1 and 10), then item t will be selected 10 times more often than $p-s$. Weight functions are used in this project to control the distribution of sales for products and customers in such a way that the rules in section 2.1 are satisfied. The weights are provided to the tool in a CSV file that contains the (*Identifier, weight*) pair. Benerator normalizes all the weights in the CSV file during data generation. < variablename = "weightings" source = "weightings.csv" distribution = "weighted" / >

The weight percentages given for products for Germany are shown in below table 1 and Customers is table 2. Similarly weight percentages given to products and customers in the USA are shown in tables 3 and 4

2.4. Configuring CSV Files

The CSV is used to weight the frequency of occurrence of an attribute value in the generated data. For each attribute, a single file is configured. Dependent attributes like the customer name for a particular customer ID are generated automatically by Benerator using the database schema.

2.5. Generated Datasets (CSV file)

In the figures 3 and 5 shows a generated datasets for sales order of Germany and also in figures 6 and 7 shows a generated dataset of the USA during the period 2009-2011.

Product	Weight percentages
1	1.9
2	2.5
3	3.3
4	2.8
5	4.2
6	3.9
7	4.1
8	1.7
9	4.5
10	2.2
11	10.5
12	6.8
13	8.9
14	8.6
15	9.8
16	6.8
17	7.4
18	9.1

Table 1: Product Weights table - Germany

3. Analyses

This section focuses on the analysis of the generated data to ensure that all the rules specified in Section 2.1 are satisfied.

We have analyzed the bike sales for two countries that data was generated for, Germany and the USA. There are 12 customers for Germany, from 11 cities, and 13 customers for the USA from 13 cities. The following sections describe Region-wise sales distribution, Monthly sales distribution, Yearly trends and Product preferences to prove that the data generated by benenerator tool is realistic.

3.1. Regional Sales

Regional sales analysis yielded the following patterns:

- Sales for USA is lower than for Germany

Customer	Weight percentages
1	5.5
2	5
3	10.5
4	8.5
5	10.5
6	6
7	11.5
8	8
9	9.5
10	6.5
11	9
12	7.5

Table 2: Customer Weights table - Germany

Product	Weight percentages
1	1.9
2	2.3
3	3.1
4	2.7
5	5.2
6	6.5
7	5.1
8	1.9
9	4.7
10	3.2
11	8.3
12	6.4
13	7.6
14	10
15	8.2
16	6.3
17	7.7
18	9.5

Table 3: Product Weights table - USA

Customer	Weight percentages
1	30.48
2	5.04
3	10.59
4	10.41
5	2.05
6	0.35
7	3.24
8	0.66
9	14.78
10	7.26
11	1.28
12	2.21
13	11.58

Table 4: Customer Weights table - USA

orderid	productid	orderid	count	orderid	customerid	sale_time	pid	pname
1	15	1	17	1	1	22-01-2009	15	PROFESSIONAL TOURING BIKE (BLACK)
2	11	1	6	1	1	22-01-2009	11	DELUXE TOURING BIKE (BLACK)
3	6	1	20	1	1	22-01-2009	6	REPAIR KIT
4	9	1	20	1	1	22-01-2009	9	WATER BOTTLE
5	12	1	5	1	1	22-01-2009	12	DELUXE TOURING BIKE (RED)
6	18	1	12	1	1	22-01-2009	18	WOMEN'S OFF ROAD BIKE EN
7	13	1	20	1	1	22-01-2009	13	DELUXE TOURING BIKE (SILVER)
8	11	1	1	1	1	22-01-2009	11	DELUXE TOURING BIKE (BLACK)
9	11	2	20	2	4	31-01-2009	11	DELUXE TOURING BIKE (BLACK)
10	3	2	6	2	4	31-01-2009	3	FIRST AID KIT
11	5	2	1	2	4	31-01-2009	5	OFF ROAD HELMET
12	14	2	4	2	4	31-01-2009	14	MEN'S OFF ROAD BIKE
13	9	2	18	2	4	31-01-2009	9	WATER BOTTLE
14	16	2	10	2	4	31-01-2009	16	PROFESSIONAL TOURING BIKE (RED)
15	3	2	2	2	4	31-01-2009	3	FIRST AID KIT
16	12	2	17	2	4	31-01-2009	12	DELUXE TOURING BIKE (RED)

Figure 3: Sales data for Germany

- USA shows low sales figures for the year 2011 figure 9

It can be seen that sales trends in the generated data follow real world patterns because they capture the effect of population, personal preferences (higher percentage of bike users in Germany), etc. It can be concluded that the data is realistic at the regional level.

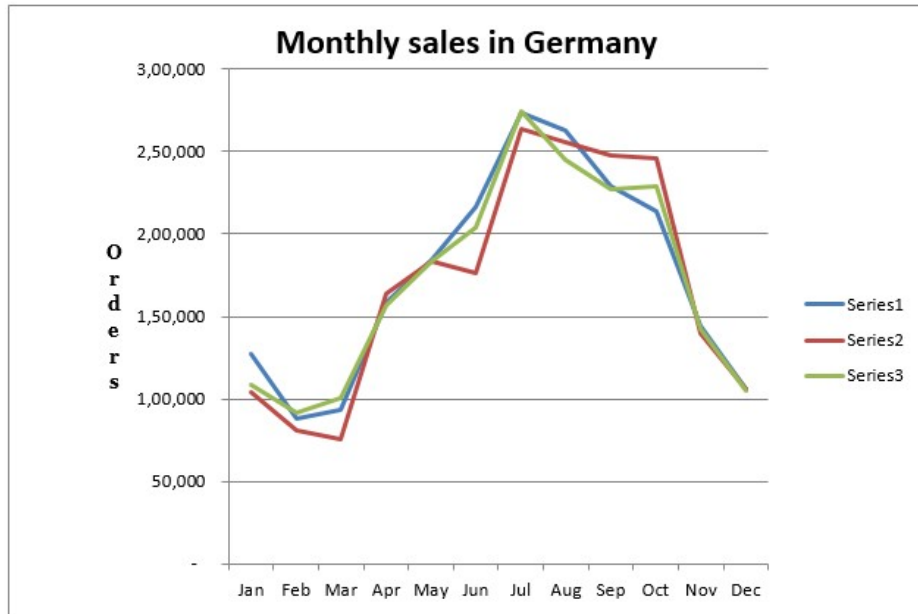


Figure 4: DE sales

code	pnum	type	price	id	customerid	postalcode	city	customername	countrycode
EN	PRTR1000	FG	510	1	13000	60549	FRANKFUI	AIRPORT BIKES	DE00
EN	DXTR1000	FG	380	1	13000	60549	FRANKFUI	AIRPORT BIKES	DE00
EN	RKIT1000	TG	45	1	13000	60549	FRANKFUI	AIRPORT BIKES	DE00
EN	BOTL1000	TG	5	1	13000	60549	FRANKFUI	AIRPORT BIKES	DE00
EN	DXTR3000	FG	400	1	13000	60549	FRANKFUI	AIRPORT BIKES	DE00
EN	ORWN1000	FG	320	1	13000	60549	FRANKFUI	AIRPORT BIKES	DE00
EN	DXTR2000	FG	410	1	13000	60549	FRANKFUI	AIRPORT BIKES	DE00
EN	DXTR1000	FG	380	1	13000	60549	FRANKFUI	AIRPORT BIKES	DE00
EN	DXTR1000	FG	380	4	16000	16341	BERLIN	CAPITAL BIKES	DE00
EN	FAID1000	TG	50	4	16000	16341	BERLIN	CAPITAL BIKES	DE00
EN	OHMT1000	TG	20	4	16000	16341	BERLIN	CAPITAL BIKES	DE00
EN	ORMN1000	FG	310	4	16000	16341	BERLIN	CAPITAL BIKES	DE00
EN	BOTL1000	TG	5	4	16000	16341	BERLIN	CAPITAL BIKES	DE00
EN	PRTR3000	FG	510	4	16000	16341	BERLIN	CAPITAL BIKES	DE00
EN	FAID1000	TG	50	4	16000	16341	BERLIN	CAPITAL BIKES	DE00

Figure 5: sales data for Germany1

3.2. Seasonal Analysis in Germany and the USA

Seasonal analysis of the data shows the following:

orderdetailid	productid	orderid	count	orderid	customerid	sale_time	pid	pname	code
1	15	1	2	1	2	2009-01-08	15	PROFESSIONAL TOURIN	EN
2	6	1	9	1	2	2009-01-08	6	REPAIR KIT	EN
3	7	1	14	1	2	2009-01-08	7	ROAD HELMET	EN
4	14	1	18	1	2	2009-01-08	14	MEN'S OFF ROAD BIKE	EN
5	15	1	1	1	2	2009-01-08	15	PROFESSIONAL TOURIN	EN
6	13	1	10	1	2	2009-01-08	13	DELUXE TOURING BIKE (EN
7	10	1	8	1	2	2009-01-08	10	WATER BOTTLE CAGE	EN
8	14	1	6	1	2	2009-01-08	14	MEN'S OFF ROAD BIKE	EN
9	10	1	7	1	2	2009-01-08	10	WATER BOTTLE CAGE	EN
10	13	1	18	1	2	2009-01-08	13	DELUXE TOURING BIKE (EN
11	15	1	18	1	2	2009-01-08	15	PROFESSIONAL TOURIN	EN
12	13	1	8	1	2	2009-01-08	13	DELUXE TOURING BIKE (EN
13	17	1	1	1	2	2009-01-08	17	PROFESSIONAL TOURIN	EN
14	14	1	12	1	2	2009-01-08	14	MEN'S OFF ROAD BIKE	EN
15	2	1	7	1	2	2009-01-08	2	ELBOW PADS	EN
16	13	1	7	1	2	2009-01-08	13	DELUXE TOURING BIKE (EN

Figure 6: sales data for US

pnum	type	price	id	customerid	postalcode	city	customername	countrycode
PRTR1000	FG	550	2	0000002000	10014	NEW YORK CITY	BIG APPLE BIKES	US00
RKIT1000	TG	45	2	0000002000	10014	NEW YORK CITY	BIG APPLE BIKES	US00
RHMT1000	TG	20	2	0000002000	10014	NEW YORK CITY	BIG APPLE BIKES	US00
ORMN1000	FG	350	2	0000002000	10014	NEW YORK CITY	BIG APPLE BIKES	US00
PRTR1000	FG	550	2	0000002000	10014	NEW YORK CITY	BIG APPLE BIKES	US00
DXTR2000	FG	460	2	0000002000	10014	NEW YORK CITY	BIG APPLE BIKES	US00
CAGE1000	TG	3	2	0000002000	10014	NEW YORK CITY	BIG APPLE BIKES	US00
ORMN1000	FG	350	2	0000002000	10014	NEW YORK CITY	BIG APPLE BIKES	US00
CAGE1000	TG	3	2	0000002000	10014	NEW YORK CITY	BIG APPLE BIKES	US00
DXTR2000	FG	460	2	0000002000	10014	NEW YORK CITY	BIG APPLE BIKES	US00
PRTR1000	FG	550	2	0000002000	10014	NEW YORK CITY	BIG APPLE BIKES	US00
DXTR2000	FG	460	2	0000002000	10014	NEW YORK CITY	BIG APPLE BIKES	US00
PRTR2000	FG	540	2	0000002000	10014	NEW YORK CITY	BIG APPLE BIKES	US00
ORMN1000	FG	350	2	0000002000	10014	NEW YORK CITY	BIG APPLE BIKES	US00
EPAD1000	TG	20	2	0000002000	10014	NEW YORK CITY	BIG APPLE BIKES	US00
DXTR2000	FG	460	2	0000002000	10014	NEW YORK CITY	BIG APPLE BIKES	US00

Figure 7: sales data for US

- Sales trends are affected by weather patterns (highest sales in summer) as in figure
- The sales in 2009, 2010 and 2011 found to be increased in figure 10 in the generated data for Germany . While in US, 2009,2010 sales were increased. But in 2011 its decreased figure 9.

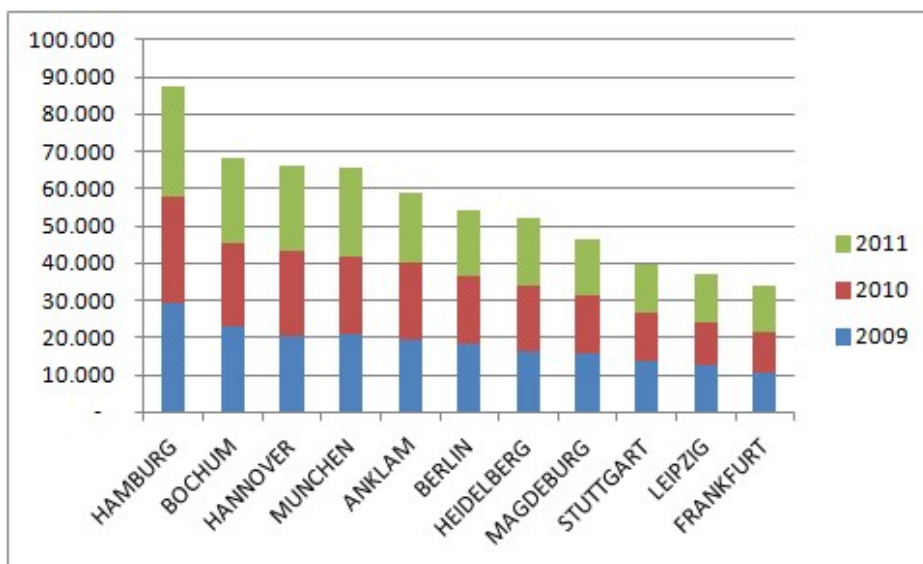


Figure 8: Region-wise sales in DE from 2009-2011



Figure 9: US sales

It can be concluded that the data shows both short and long-term trends found in real data. We realized that seasonal analysis of bike sales are low in

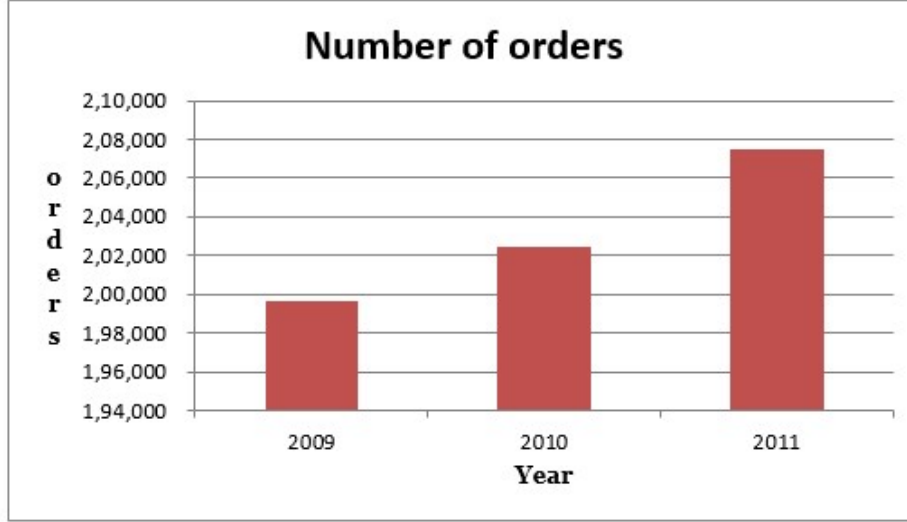


Figure 10: DE sales

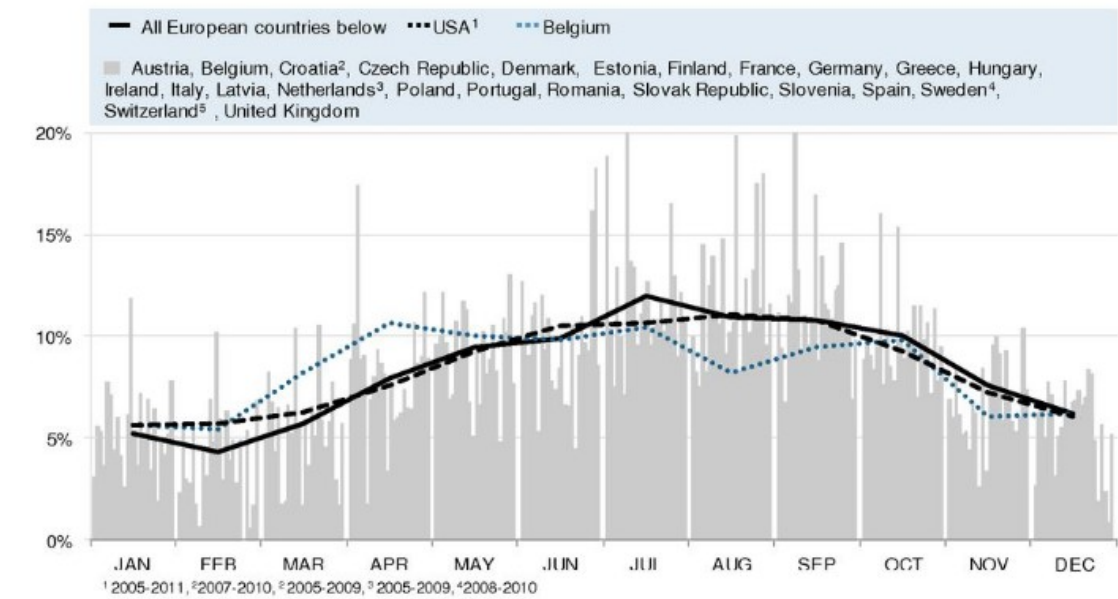
the winter season in Germany and the USA. During May to October months, the number of bike sales are increased in both countries figures 12 and 4. But, the highest number of sales are happening in Germany and the USA during the summer time from figures 10 and 12.

Figures 11, 12 and 10 shows the comparison between real data [7] and realistic synthetic data sales (we generated data for Germany and the USA using Benerator). We can see that the sales data of real and realistic synthetic data are almost similar when compared with the summer and winter months in a year. This states that the generation process we follow to generate a data also provides accuracy from the real data.

3.3. Product preference

3.3.1. Color

Colors play an important role in designing for marketing materials. Usage and Meanings can vary from region to region. We should know about culture and behavior of customers before selling the finished products. Therefore it can increase in the attention of customers which leads to more sales in bikes.



Source: EU CARE database, 2005-2010 and USA FARS database 2005-2011

Figure 11: Percentage of Fatal Bicycle Accidents

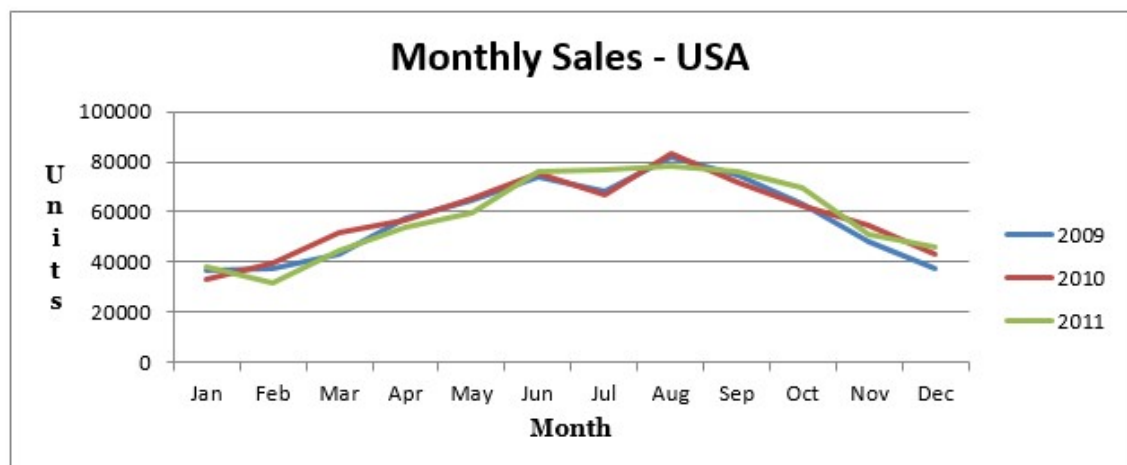


Figure 12: US sales

After studying about given colors black, red and silver. We made an analysis for touring bikes (deluxe, professional) and off-road bikes(men, women) of GBI company. We analyzed that bikes sells more in black colored in figure 13, 14 among all categories and then followed by silver and red will be the least preference. We also found men are using more bikes than a women [4]

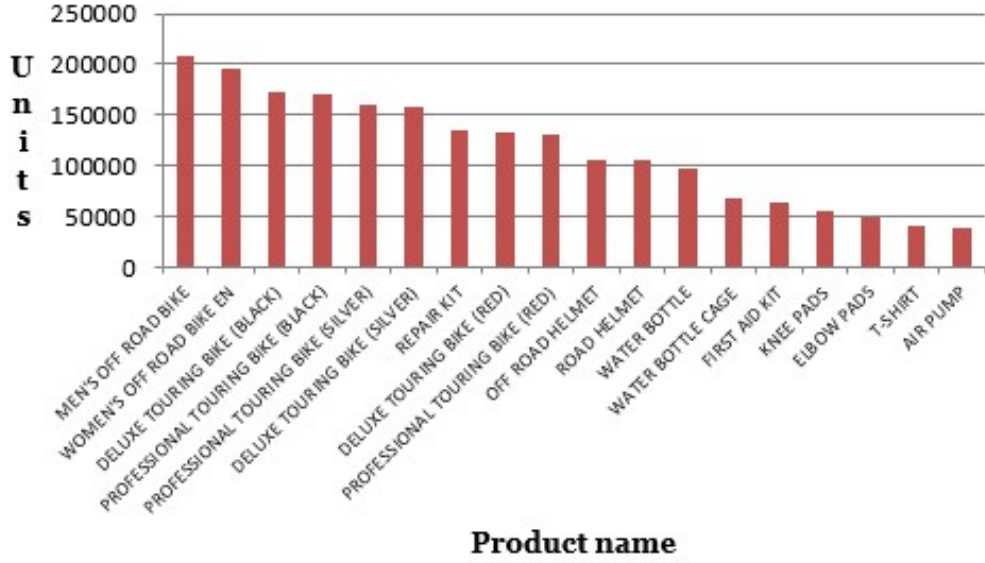


Figure 13: product sales in DE

3.3.2. Accessories

Figure 13 shows sales in parts and accessories in the year 2009, 2010 and 2011 respectively in Germany.

Therefore, we concluded that the generated data meets all the requirements discussed in section 2.1.

4. Related Work

In present market we have many other similar tools to generate realistic datasets having specific features used for testing purposes like GS Data Generator, Turbo Data, GT DataMaker etc. Many of these tools only generate

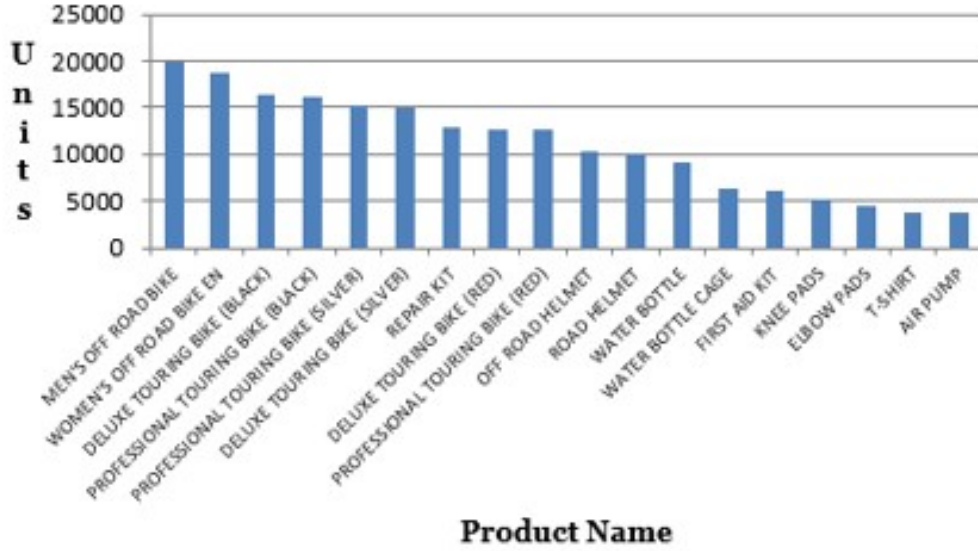


Figure 14: Product sales in US

a large volume of data but hard to create anything that stands for workflow validation and also not much useful in real-world enterprise applications. Most of these tools are compatible only with limited operating systems and flexible for few database management systems. When these tools compared with Benerator, an open source tool it has many better functionalities and Performance with built-in configuration and supports all types of formats. Benerator tool is easy to understand and implement. It helps in processing fast data generation and can create data in large volumes.

5. Conclusion

In this paper, we have presented Benerator tool used for automatic generation of a large amount of data. We have studied the behavior of Benerator tool using different techniques, which leads to extend ways of Benerator tool to add a new area for data generation. We explained how we extended the given Benerator tool by using different functions. Benerator tool controlled by the descriptor file combines to produce an effective tool for data generation. MYSQL database is also employed in the implementation of data generation in addition to storing the attributes of generated data.

5.0.1. *Learning from Project*

The chief learnings from this project are:

- Benerator can be used to generate realistic data for SAP solutions.
- Support and development for Benerator have ended, therefore long-term development should consider other tools.
- Benerator does not meet the performance requirements for generating datasets that are more than one million rows figure 15

5.0.2. *Problem*

The problems faced during the project using this tool were

- Development for the tool ended in 2009, with release v 0.98
- The online forum is inactive and no longer accepts registrations
- The documentation is not exhaustive enough to cover all use cases.

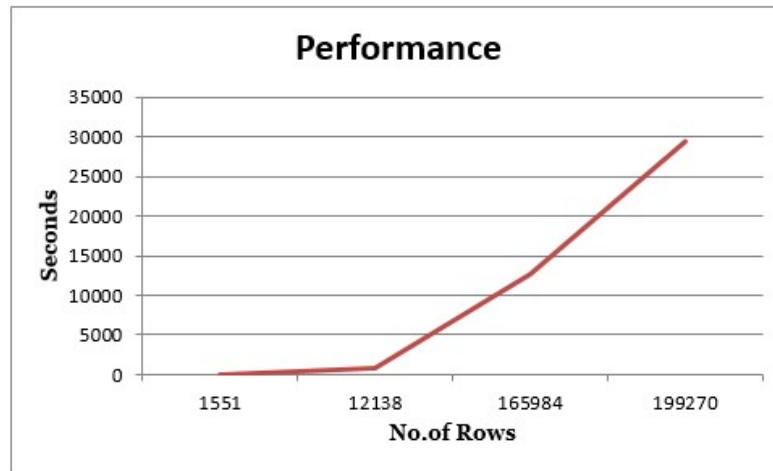


Figure 15: Tool performance per data generation

6. Acknowledgement

This project was assisted by SAP UCC Magdeburg, we are really thankful to Mr.Tim Bottcher, Mr.Chris Bernhardt our project guides and SAP UCC Team for encouraging to complete our project in a successful manner. This task is also successfully finished with the effort and coordination of our Team members.

7. References

- [1] V. Ayala-Rivera, P. McDonagh, T. Cerqueus, L. Murphy, Synthetic data generation using benenerator tool, arXiv preprint arXiv:1311.3312 (2013).
- [2] K. McLeod, D. Flusche, A. Clarke, Where we ride: Analysis of bicycling in american cities (2013).
- [3] B. G. G. B. F. F. H. I. T. N. P. S. S. Members of national Bicycle Industry Associations in 14 different countries: Austria, Belgium, Turkey., European bicycle market (2015) 10–75.
- [4] C. Emond, W. Tang, S. Handy, Explaining gender differences in bicycle behavior, in: Active Living Research Conference February, volume 19.
- [5] M. Nezhad, K. Kavehnezhad, et al., Choosing the right color: a way to increase sales, International Journal of Asian Social Science 3 (2013) 1442–1457.
- [6] J. A. Bellizzi, R. E. Hite, Environmental color, consumer feelings, and purchase likelihood, Psychology & marketing 9 (1992) 347–363.
- [7] N. Torslov, Cycling,health and safety, OECD/International Transport Forum(2013),Cycling, Health and Safety OECD Publishing / ITF (2013) 24–25.