

Domain-Driven, Generative Data Model for BigPetStore

Ronald J. Nowling
Red Hat, Inc.
Raleigh, NC 27601
Computer Science & Engineering,
University of Notre Dame,
Notre Dame, IN 46617
Email: rnowling@redhat.com

Jay Vyas
Red Hat, Inc.
Raleigh, NC 27601
Email: jvyas@redhat.com

Abstract—Abstract

I. INTRODUCTION

Big data applications process large, continuously-changing data sets with the general goal of information and knowledge extraction. New tools such as Hive, Pig and Crunch [reference](#) based around big data frameworks such as Hadoop, Storm, and Spark [reference](#) are being released every day. Development, testing and benchmarking of big data software, algorithms, and infrastructure pose unique challenges. Though software and algorithms will be deployed to large clusters for benchmarking and production usage, engineers develop and test these tools on desktops due to the speed and ease of development and debugging on a local machine and limited access to large clusters due to the high costs of purchasing, provisioning, and maintaining such clusters. Software and algorithms are being developed at incredible speed with new releases everyday. The high speed of development driven by the need to be the first to market in a competitive field often results in buggy software. There is great need for consistent and rigorous testing and benchmarking frameworks to aid in the identification of bugs while maintaining the high-speed of development. Access to high-quality, semantically-rich data sets and complex, realistic analytics workflow examples that can scale from local development environments to large clusters are key to increasing access to and development of consistent and rigorous testing and benchmarking frameworks.

Over time, more high-quality data sets are becoming publicly available, but these data sets often carry burdens that force developers to make undesirable trade offs. Since few data sets are scalable, allowing the user to pick how much data they want, users are forced to choose between small data sets which are useful for development and testing on a local machine, but not for production-scale tests or benchmarking, or large data sets which are prohibitively expensive to host, transfer, and store. Furthermore, access to large data sets can be unreliable due to dependence on third party benefactors for hosting and requires significant amounts of bandwidth to download. Several companies (such as Twitter [reference](#)) offer access to their data but the costs and usage restrictions can be undesirable.

Methods for generating synthetic data offer multiple advantages over pre-generated data sets. The code for generating the

data is often small, making it easy and cheap to host, transfer, and store. Licensing restrictions and data privacy issues are avoided [cite](#). And, the size of the data sets can be scaled arbitrarily, enabling seamless usage of the same data set across a range of situations.

Multiple synthetic data set generators exist. TeraGen [cite](#) and the Intel Hadoop Benchmark Suite [10] can generate data sets quickly and at any scale. Targetted towards specific benchmarks such as Map-Reduce jobs and Hive queries, the resulting data is semantically-void and not useful for other tasks. The PGDF (Parallel Data Generation Framework) can quickly generate data and queries that structurally satisfies a user-defined relational database schema [9], [15], [8], [14]. The Myriad Toolkit improves on PGDF through a more descriptive model specification language [6]. Myriad allows the user to specify the distributions for sampling values, simple arithmetic equations for setting the value of fields from other fields, and restrictions on fields domains. Due to the inability to specify underlying dynamic models with latent variables, the data generated by PGDF and Myriad lacks informational complexity and genuine semantics, making it inappropriate for analytics.

More advanced approaches such as Oligas [5] and SMOTE [7] attempt to generate or extend real data sets with synthetic data sets that satisfy characteristics learned from a real reference data sets. [privacy, accuracy without details of underlying process, SMOTE - feature space](#) Oligas appears promising but will need to overcome several difficult challenges in automatically characterizing complex reference data to determine relevant constraints and learn accurate models that capture the effect of latent variables.

Until such methods have reached maturity, the development of customized, domain-driven models is best suited for satisfying the needs of big data software and algorithms development. For example, the Information Discovery and Analysis System (IDAS) Data Set Generator (IDSG) [11], [13] was developed for the Department of Homeland Security to develop synthetic data for effectiveness of IDASs designed to identify potential terrorism threat scenarios.

[Huppler, Pavlo. standard benchmarks](#)

The BigPetStore application aims to be a big data appli-

cation blueprint built around a fictional chain of pet stores, which incorporates a data generator and complex examples of cleaning, aggregation, and analytics using well-known tools such as Hive and Mahout from across the Hadoop and Spark ecosystems [reference](#). BigPetStore’s usefulness has been demonstrated by its incorporation into the open-source Apache BigTop distribution [1] for testing and as a reference application.

In this work, we present a new domain-driven model and simulation for generating arbitrarily large, complex data for testing and benchmarking of algorithms, software, and infrastructure in the big data space. Compared with BigPetStore’s previous data generator, our model and implementation can generate data which contains geospatial, temporal, and quantitative features, similar to the type of data which businesses and organizations might typically encounter. Combining the features of the TeraGen (scalability) [reference](#) and MovieLens (semantically rich) [reference](#) input data sets commonly used for big data benchmarking, we present a scalable synthetic data set generation implementation which can be used to benchmark lower level tasks (such as sorting) as well as higher level tasks (such as clustering, recommending, and search indexing) at arbitrary scale. To validate this approach, we perform an example analysis designed to guide a fictional advertising campaign (a typical business use case) to demonstrate that the model generates data with desirable properties (global patterns and region specific transactions). We have made the model’s implementation available as open source software.

II. OVERVIEW OF MODEL & SIMULATION PROCEDURES

[need an intro paragraph](#)

A. Stores

Store records contain a unique identifier and a location in the form of a 5-digit zip code. The locations of the stores are modeled by a probability density function (PDF) that gives the probability that a store’s location is the given zip code. We designed the PDF to give to zip codes in high-population, high-income areas the highest probabilities. The PDF is composed of two individual PDFs. One PDF determines the probability of each zip code as the population of that zip code over the total population:

$$p(\text{location} = z | \text{population}(z)) = \frac{\text{population}(z)}{\sum_i \text{population}(i)} \quad (1)$$

The second PDF scales the probabilities of the zip codes by their incomes. The zip code with the highest-income has a probability s -times larger than that of the lowest-income zip code. The values in-between are interpolated using an exponential function:

$$p(\text{location} = z | \text{income}(z)) = s^{\left(\frac{\text{income} - \min_i \text{income}(i)}{\max_i \text{income}(i) - \min_i \text{income}(i)} - 1 \right)} \quad (2)$$

The combined PDF is given as:

$$p(\text{location} = z) = Z^{-1} p(\text{location} = z | \text{population}(z)) p(\text{location} = z | \text{income}(z)) \quad (3)$$

where Z is the normalization factor. Given the small size ($\approx 30,000$ zip codes) of the data set, Z can be found directly by iterating over all of the zip codes and summing the scores.

The PDFs are parameterized using external data sets. The zip codes’ latitudes and longitudes came from the Zip Code Database Project [4]. The population and household income data for the zip codes were taken from the U.S. Census American Community Survey [2].

The stores’ locations are generated by sampling zip codes from the PDF using a Roulette Wheel sampler.

B. Customers

Customer records consist of a unique identifier, a name, a location given as a 5-digit zip code, and the number and species of their pets.

[names](#)

For each proposed customer location zipcode z , we compute the distance between z and every store’s zipcode to find the closest store. We compute the distance d between z and the zipcode s_z of the closest store. The weight w_z for each zipcode z is modeled using an exponential distribution with the average distance β .

$$w_z = \beta^{-1} \exp(-\beta^{-1} d_m(z)) \quad (4)$$

$$d_m(z) = \min_{s_z} d(z, s_z)$$

To find the probabilities, the weights are normalized:

$$p(\text{location} = z) = W^{-1} w_z \quad (5)$$

$$W = \sum_z w(z)$$

The customer locations are sampled from the PDF using a Roulette Wheel sampler.

We determine the number of pets N_p each customer has by sampling from a discrete uniform distribution of integers. We then sample the number of dogs N_d as a discrete uniform distribution of integers between 0 and N_p . The remaining number of pets are assigned to be cats.

$$N_p \sim U(1, b) \quad (6)$$

$$N_d \sim U(0, N_p)$$

$$N_c = N_p - N_d$$

C. Transactions

A transaction model is created for each customer and integrated over the simulation's timeframe to generate transactions. Transaction records consist of a customer, a store, a date and time, and a list of items purchased and their prices. The store is set to the store located closest to the customer. The date and time is given as a floating-point number with units of days since the start of the simulation. (Timezones are not taken into account.) Transaction times are generated from a Monte Carlo process based on the customer will use up items such as dog food (Section II-C1). The purchased items are generated by a Hidden Markov Model (HMM) parameterized by the transaction time and amount of time remaining for the items in the customer's inventory (Section II-C2).

1) *Transaction Times*: Transaction times are simulated using a Monte Carlo method (Figure 1). The usage over time of each category is simulated, from which transaction times are found for each item category. The earliest transaction time is taken as the proposed transaction time. The probability of the transaction time is calculated using a PDF. If rejected, new transaction times are proposed. Otherwise, the purchased items are modeled using a separate process (discussed below). After the items are chosen, the process begins again with the simulation of the item category usages.

In the first stage of the process, the usage of items over time is modeled using a simple stochastic differential equation (SDE):

$$\frac{da_i}{dt} = -\mu_i + \sigma_i dW_i(t) \quad (7)$$

where a_i is the amount of item category i "stuff" remaining at time t . The variable u_i gives the average usage rate and σ_i^2 gives the variance of the usage rate for item category i . The SDE is numerically integrated using the Euler-Maruyama method[12] with timestep Δt_n sampled from a Poisson process:

$$\begin{aligned} a_{n+1,i} &= a_{n,i} - \mu_{i,c} \Delta t_n + \sigma_{i,c} \sqrt{\Delta t_n} R_n \\ \Delta t_n &\sim \text{Exp}(\lambda_i) \\ R_n &\sim N(0, 1) \end{aligned} \quad (8)$$

where λ_i is the item category's average usage rate, μ_i is the item category's average amount used per use, and σ_i^2 is the variance of the amount of the item category used per use. The parameter λ_i is defined for each item category. The parameters $\mu_{i,c}$ and $\sigma_{i,c}^2$ are found based on a base rate for each item category multiplied by the number of pets of the appropriate species the customer has. The exhaustion time $T_{E,i}$ for each item category is found by simulating the usage until $a_n = 0$.

The proposed transaction time is found from the exhaustion times. For each item category, the offset time $T_{O,i}$ between when a customer would visit the store and the exhaustion time is sampled from an exponential distribution parameterized separately for each customer c . A proposed transaction time $T_{P,i}$ for each item category is found by subtracting the offset time $T_{O,i}$ from the exhaustion time $T_{E,i}$. The earliest proposed

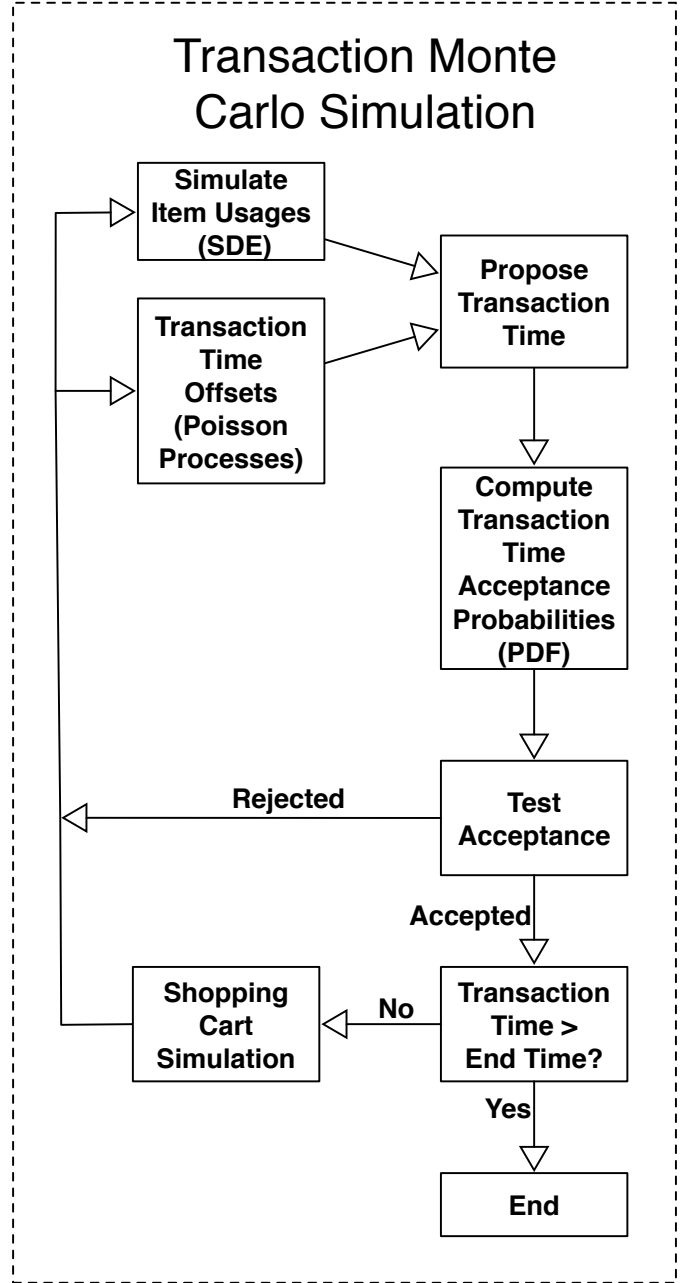


Fig. 1. Flowchart of Transaction Monte Carlo Simulation

transaction time is taken as the overall proposed transaction time T_T .

$$\beta_c \sim U(a, b) \quad (9)$$

$$T_{O,i} \sim \text{Exp}(\beta_c^{-1}) \quad (10)$$

$$T_{P,i} = T_{E,i} - T_{O,i} \quad (11)$$

$$T_T = \min_i \{T_{P,i}\} \quad (12)$$

The acceptance probability $p(t = T_T)$ of the proposed transaction time T_T is modeled using a PDF that incorporates factors such as store hours. The PDF could be extended to incorporate additional factors including time-dependent events

such as sales, weather, and the amount of money the customer has available at the time.

$$p(t = T_T) = p(t = T_T | \text{store hours}) \quad (13)$$

The proposed transaction time is accepted if $p(t = T_T) < r$, where $r \sim U(0, 1)$. Until the proposed transaction time is accepted, a new proposed transaction time is generated by sampling new offset times. If the proposed transaction time is accepted, the shopping cart of items is generated through a separate simulation discussed in Section II-C2 below. Transactions are generated until the accepted transaction time is later than the end time given as a simulation parameter.

2) *Shopping Cart Items*: A customer's purchases in each transaction are modeled using a layered Hidden Markov Model (HMM) (Figure 2). The HMM has three states: the start state, the purchase states, and the stop state. There are an infinite number of purchase states, identified by the exhaustion times of the customer's item categories, the transaction time, and the number of items purchased in the current transaction. The observables of the states correspond to the item categories. Each item category i is assigned a weight w_i according to the amount of time between the transaction time T_t and $T_{E,i}$ when each item category will be exhausted (based on the usage simulations):

$$w_i = -\lambda_i \exp(-\lambda_i(T_{E,i} - T_T)) \quad (14)$$

The probability $p(X = i)$ of choosing item category i is given by the normalizing the weight w_i by the sum of all of the weights W :

$$p(X = i) = W^{-1} w_i \quad (15)$$

$$W = \sum_i w_i$$

The item category for the purchase is chosen using a Roulette Wheel sampler.

Markov Models are used to model the customer's purchasing behavior for each item category. Each state corresponds to an item in that category. The transition probabilities are assigned according to the customer's buying habits. For example, if a customer tends to buy the same brand of dog food repeatedly, the outgoing edges to all products of the same brand as the product in the current state will be higher than those for other brands' products. If the customer does not care about flavor, the outgoing edges to all products from the current brand will be equal, otherwise the edge weights can be adjusted to prefer a specific combination of brand and flavor. Likewise, if a customer cares more about flavor or cost, the edge weights will be assigned appropriately. To simulate an item purchase, the Markov model is transitioned forward by one state. The current states of the Markov models are kept between purchases.

$$Pr(X_{n+1} = x | X_n) = \dots \quad (16)$$



Fig. 3. Distributions of Distances of Customers to Closest Stores

After an item is chosen, the item category's exhaustion time is updated by propagating the item category usage SDE described in Section II-C1. **state transition, stop state**

The HMM can be extended to consider additional factors such as the amount of money spent so far in the transaction and the customer's spending habits.

III. IMPLEMENTATION

The models and simulations were implemented using Python. The source code is available on GitHub [3] at <https://github.com/rnowling/bigpetstore-data-generator> under the Apache 2 license.

IV. EVALUATION OF THE MODEL

intro paragraph

A. Example Analysis of Generated Data

Using data for 10 stores, 10,000 customers, and five years of transactions generated from the model, we performed two example analyses in support of a fictional advertising campaign. The analyses were designed to be realistic and driven by real-world business concerns. Due to limited advertising budgets, businesses need to decide who to target and what sort of advertisements are likely to be effective for each type of customer. We analyzed data generated from the model to determine where the advertising campaigns should be targeted geographically and identify customer purchasing profiles.

First, we wanted to identify how close most customers live to a store since advertising to people who live or work too far away from a store is unlikely to be effective and will increase our costs. We analyzed the distribution of distances between customers' homes and their nearest stores. We found that most customers tend to live within 15 miles of their closest store, suggesting that we should limit our advertising to a 15-mile radius around each our store.

Secondly, we profiled our customers' purchasing habits to optimize the advertising campaign's effectiveness by customizing the advertised products for each customer. For each customer, we generated a feature vector by computing the

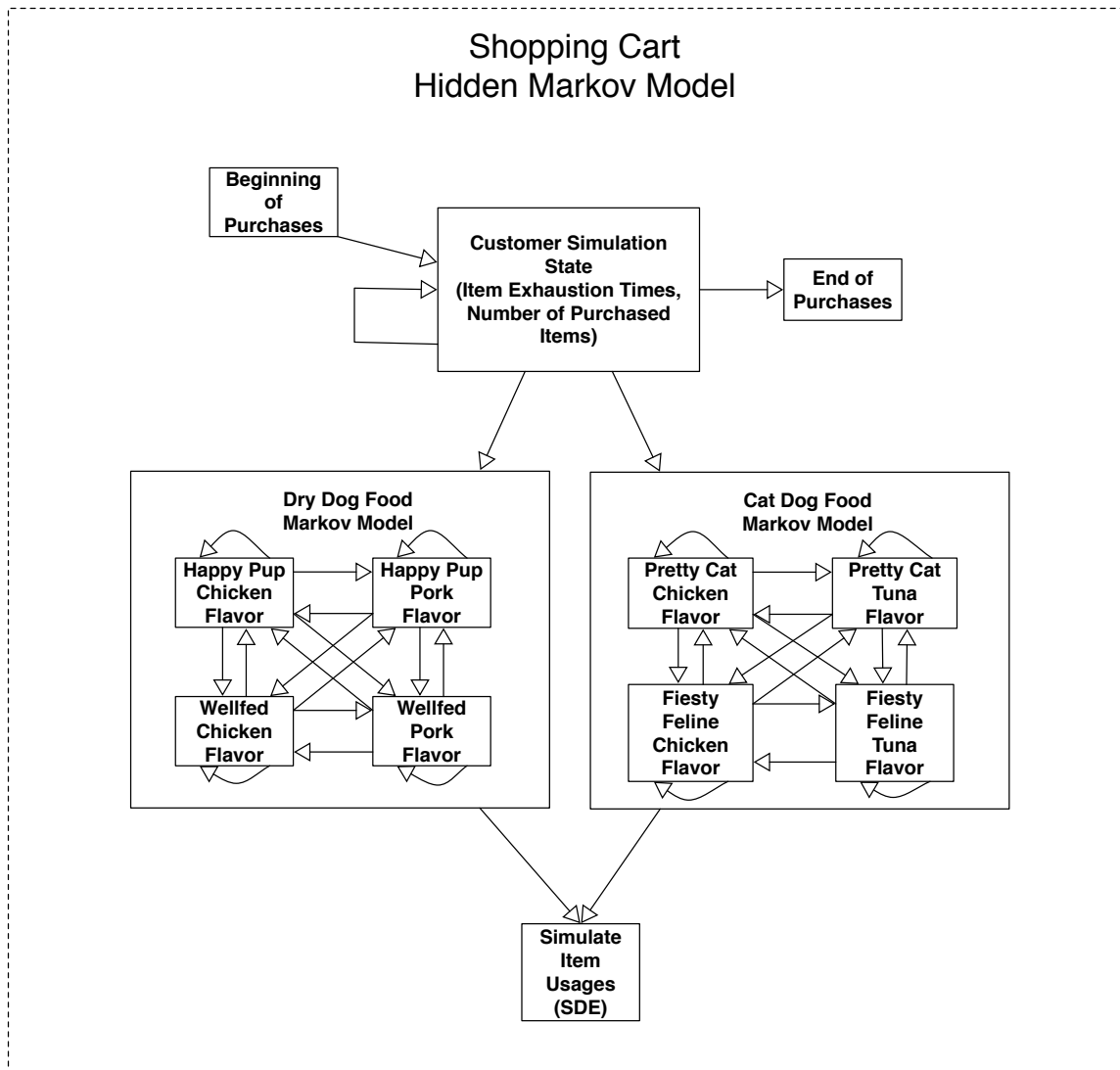


Fig. 2. Example Shopping Cart Hidden Markov Model

frequency with which they would purchase the same brand or flavor from one transaction to the next. We clustered the customers' feature vectors using the KMean algorithm with a range of cluster counts. Twenty clusters converged the error.

Analysis of the feature vectors and clusters revealed four predominant purchasing profiles (Figure 4). High frequencies of purchasing the same flavors repeatedly were tightly-correlated with high frequencies of purchasing the same brands – these customers were likely to be happy with a particular item and kept purchasing the same item repeatedly. For our advertising campaign, it would be unlikely to get these customers to purchase different items, so we should create incentives to purchase larger quantities, especially when inventory levels are high and needed to be depleted. Other customers had a tendency to purchase either the same flavor or brand repeatedly, but varied in their choice of the other. For customers who prefer a particular brand, we should target our advertising campaign to suggest other flavors sold by that brand. Likewise, for customers who prefer a particular flavor, we should target our advertising campaign to suggest other brands with that

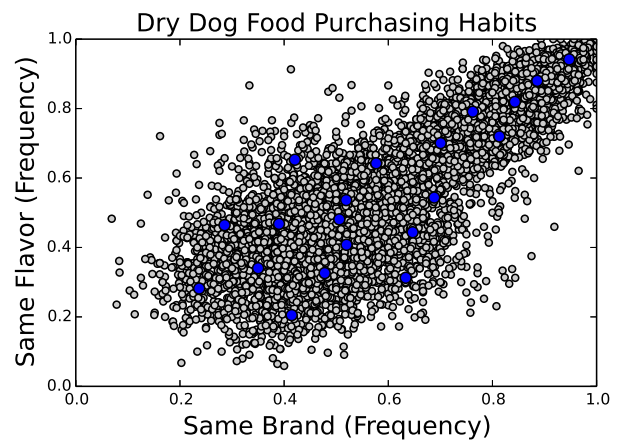


Fig. 4. Clustering of Brand and Flavor Purchasing Preferences

flavor. Lastly, some customers had very low frequencies of purchasing neither the same brand nr flavor. Other factors, such

TABLE I. BENCHMARKS OF DATA GENERATOR

Stores	Customers	Simulated Time (years)	Transactions	Data Size (MB)	Run Time (min)
10	10,000	1	279,870	18	3.5
100	10,000	1	279,586	18	3.5
10	1,000	5	123,309	8	1.1
100	1,000	5	127,064	8	1.2
10	10,000	5	1,275,542	84	11.0
100	10,000	5	1,268,403	84	10.5

as cost, not represented in this analysis may be driving these customers' purchasing habits and will need further study.

B. Scaling of Data Size and Run Time

BigPetStore aims to scale from a local desktop to a large cluster. As a result, we benchmarked the generation of data on a laptop with a 2 GHz Intel Core i7 CPU, 8 GB of RAM, and a 256 GB SSD using the Python implementation with a single thread. Using the test setup, between 1,500 and 2,000 transactions can be generated per second (Table I). As the customers' transaction simulations are independent of one another, the transaction generation can easily be parallelized so that 1,500-2,000 transactions can be generated per thread per second.

The number of transactions (and hence, data size) grows as $O(N_c T)$ where N_c is the number of customers and T is the amount of time to be simulated. The amount of data generated can be scaled as large as necessary simply by increasing N_c and T . Five years of transactions for 100,000 customers will generate approximately 1 GB of data. A terabyte of data can be generated by simulating 100 million customers over five years.

V. DISCUSSION AND CONCLUSION

implementation [want a modular framework with clearly-defined extension points, modularity so that users do not need to know about the machinations of the entire model, and so that modules are easy to develop and turn on/off dynamically]

[spark, hadoop ecosystem, parallelization of generating stores, customers, and transactions]

[model extensions]

ACKNOWLEDGMENT

The authors would like to thank...

REFERENCES

- [1] <http://bigtop.apache.org/>.
- [2] <http://www.census.gov/acs/www/>.
- [3] <https://github.com/>.
- [4] Zip code database project. <http://zips.sourceforge.net/>.
- [5] Alexander Alexandrov, Christoph Brücke, and Volker Markl. Issues in big data testing and benchmarking. In *Proceedings of the Sixth International Workshop on Testing Database Systems - DBTest '13*, page 1, New York, New York, USA, 2013. ACM Press.
- [6] Alexander Alexandrov, Kostas Tzoumas, and Volker Markl. Myriad: scalable and expressive data generation. *Proceedings of the VLDB ...*, 5(12):1890–1893, 2012.

- [7] Nitesh V Chawla, Kevin W Bowyer, Lawrence O Hall, and W Philip Kegelmeye. SMOTE: synthetic minority over-sampling technique. *Journal of Artificial Intelligence Research*, 16:321–357, 2002.
- [8] Michael Frank, Meikel Poess, and Tilmann Rabl. Efficient update data generation for DBMS benchmarks. In *Proceedings of the third joint WOSP/SIPEW international conference on Performance Engineering - ICPE '12*, Boston, Massachusetts, USA, 2012. ACM Press.
- [9] Ahmad Ghazal, Tilmann Rabl, Mingqing Hu, Francois Raab, Meikel Poess, Alain Crolette, and Hans-Arno Jacobsen. Bigbench: Towards an industry standard benchmark for big data analytics. In *Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data*, pages 1197–1208, 2013.
- [10] Shengsheng Huang, Jie Huang, Jinquan Dai, Tao Xie, and Bo Huang. The HiBench benchmark suite: Characterization of the MapReduce-based data analysis. In *2010 IEEE 26th International Conference on Data Engineering Workshops (ICDEW 2010)*, pages 41–51. IEEE, 2010.
- [11] Daniel R. Jeske, Ryan Rich, Behrokh Samadi, Pengyue J. Lin, Lan Ye, Sean Cox, Rui Xiao, Ted Younglove, Minh Ly, and Douglas Holt. Generation of synthetic data sets for evaluating the accuracy of knowledge discovery systems. In *Proceeding of the eleventh ACM SIGKDD international conference on Knowledge discovery in data mining - KDD '05*, page 756, New York, New York, USA, 2005. ACM Press.
- [12] Peter E. Kloeden and Eckhard Platen. *Numerical Solution of Stochastic Differential Equations*. Springer, 2013.
- [13] PJ Lin, B Samadi, and Alan Cipolone. Development of a synthetic data set generator for building and testing information discovery systems. In *Information Technology – New Generations, 2006*, pages 707–712, Las Vegas, NV, 2006. IEEE.
- [14] Tilmann Rabl, Michael Frank, Hatem Mousselly Sergieh, and Harald Kosch. A data generator for cloud-scale benchmarking. In Raghunath Nambiar and Meikel Poess, editors, *Performance Evaluation, Measurement and Characterization of Complex Systems*, pages 41–56. Springer, 2011.
- [15] Tilmann Rabl and Meikel Poess. Parallel Data Generation for Performance Analysis of Large, Complex RDBMS. In *Proceedings of the Fourth International Workshop on Testing Database Systems*, pages 1–6, 2011.