

Detecting Changes in Large Data Sets of Payment Card Data: A Case Study

Chris Curry, Robert Grossman*,
David Locke and Steve Vejckik
Open Data Group

Joseph Bugajski
Visa International

Abstract

An important problem in data mining is detecting changes in large data sets. Although there are a variety of change detection algorithms that have been developed, in practice it can be a problem to scale these algorithms to large data sets due to the heterogeneity of the data. In this paper, we describe a case study involving payment card data in which we built and monitored a separate change detection model for each cell in a multi-dimensional data cube. We describe a system that has been in operation for the past two years that builds and monitors over 15,000 separate baseline models and the process that is used for generating and investigating alerts using these baselines.

Categories and Subject Descriptors: G.3 [Probability and Statistics]: Statistical computing, statistical software; I.5.1 [Models]: Statistical Models

General Terms: Algorithms

Keywords: baselines, data quality, change detection, cubes of models

1 Introduction

It is an open and fundamental research problem to detect interesting and actionable changes in large, complex data sets. In this paper, we describe our experiences and the lessons learned over the past three years developing and operating a change detection system designed to identify data quality and interoperability problems for Visa International Service Association (“Visa”). The change detection system produces alerts that are further investigated by analysts.

Robert Grossman is also a faculty member at the University of Illinois at Chicago.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.
KDD’07, August 12–15, 2007, San Jose, California, USA.
Copyright 2007 ACM 978-1-59593-609-7/07/0008 ...\$5.00

The problem is difficult because of the following challenges:

1. Visa’s data is high volume, heterogeneous and time varying. There are 6,800 payment transactions per second that must be monitored from millions of merchants located around the world that are processed over a payment network that connect over 20,000 member banks. There are significantly different patterns across regions, across merchants, during holidays and weekends, and for different types of cardholders. See Figure 1 for example.
2. Alerts arising from the change detection system generally require human examination. Because of this it is necessary to balance generating a meaningful number of alerts versus generating a manageable number of alerts. If too many alerts are generated, it is not practical to manage them. If too few alerts are generated, they are generally not meaningful.

In this paper, we describe our experiences using a methodology for addressing these challenges. The methodology we use is to build a very large number of very fine grained baselines, one for each cell in a multi-dimensional data cube. We call this approach Change Detection using Cubes of Models or CDCM.

For example, we built separate baseline models for each different type of merchant, for each different member bank, for each different field value, etc. In total, over 15,000 different baseline statistical models are monitored each month and used to generate alerts that are then investigated by analysts.

We believe that this paper makes the following research contributions:

1. First, we have highlighted an important category of data mining problems that has not received adequate coverage within the data mining community and whose importance will continue to grow over time.
2. Second, we have described our experiences implementing and using a new change detection algorithm called CDCM that is designed to scale to large, complex data sets by building a separate change detection model for each cell in a data cube.
3. Third, we have introduced a software architecture that can reliably scale with tens of thousands of different individual statistical or data mining models.

Section 2 contains some background on payment card transactions. Section 3 describes the Change Detection using Cubes of Models

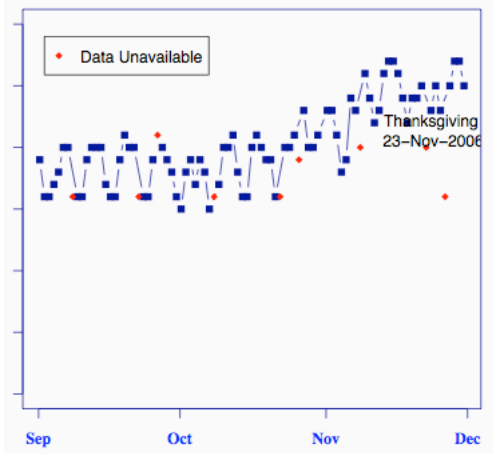


Figure 1. This is an example of one of the measures monitored in this project. This graph shows how the ratio of declined transactions varies for one of the Merchant Category Codes (MCC) monitored. Note the daily, weekly and monthly variation in the data. This variation, which is typical of the measures tracked, is one of the reasons detecting changes in this data is challenging. The vertical axis scale has been omitted due to confidentiality reasons.

(CDCM) algorithm that we introduce. Section 4 describes some typical alerts detected by the system we developed and deployed. Section 5 describes the architecture of the system we developed. Section 6 describes the implementation. Section 7 describes some of the lessons learned. Section 8 describes related work. Section 9 is the summary and conclusion. There is also an appendix that describes the XML representation of the change detection models we used.

2 Background on Payment Card Transactions

2.1 Processing a Transaction

We begin by providing some background on payment card transactions that will make this paper more self-contained. In this section, we define cardholder, merchant, acquiring bank, and issuing bank and describe the major steps involved when using a payment card.

1. A transaction begins when *cardholder* purchases an item at a *merchant* using a payment card. The payment card contains an account number which identifies the cardholder.
2. The merchant has a relationship with a bank called the *acquiring bank*, which agrees to process the payment card transactions for the merchant. The acquiring bank provides the merchant with a terminal or other system to accept the transaction and to process it.
3. The acquiring bank has a relation with a financial payment system, such as those operated by Visa and MasterCard. The transaction is processed by the acquiring bank and passed to the payment system. Visa operates a payment system called VisaNet.
4. The payment system processes the transaction and passes the transaction to the bank (the *issuing bank*) that issued the payment card to the cardholder. In other words, one of the essential roles of the payment system is to act as a hub or interme-

diary between the acquirer and the issuer.

5. The issuing bank processes the transaction and determines if there are sufficient funds for the purchase, if the card is valid, etc. If so, the transaction is authorized; the transaction can also be declined, or a message returned asking for additional information. The issuing bank also has a relationship with the cardholder or account holder. For example, with a credit card, the cardholder is periodically billed and with a debit card the appropriate account is debited. For the year ending September 30, 2006, there were over 1.5 billion Visa cards in circulation.
6. For each of these cases, the path is then reversed and the transaction is passed from the issuing bank to the payment system, from the payment system to the acquiring bank, and from the acquiring bank to the merchant.

Our problem was to use baselines and change detection algorithms to help detect data and interoperability problems at Visa [5]. Payment data arrives at Visa from millions of merchant locations worldwide. For the year ending September 30, 2006, total annual global card sales volume was over USD \$4.45 trillion¹ Payment data is processed through risk management rules set by over 20,000 individual member banks (issuing and acquiring banks). These rules determine if a payment authorization request from a merchant either is approved or rejected by the paying bank.

For this problem, we built separate baselines for a variety of data fields, for each member bank, and for thousands of merchants. Overall, over 15,000 separate baselines are currently used each month to monitor payment card transactions at Visa.

2.2 Baselines for Field Values

Note: The examples in this section are hypothetical and only used for the purposes of illustrating how to define baselines.

A payment card transaction typically includes a number of fields, such as information about the point of services (POS) environment, the merchant's type of business and location, the cardholder's identity, the transaction currency, the transaction amount, and bank routing information.

We begin with an informal description of baselines based upon a simplified example. In this simplified example, assume that one of the fields of interest describes characteristics of the point of service (POS). Specifically, we assume that this field can take the following (hypothetical) values: 00, 01, 02, 03, and 04.

For an observation period of a week, assume that the frequency of these values for a certain acquirer is given by the first table in Figure 2. Later, during the monitoring, assume that distribution is instead given by the right hand table in this figure. The observed distribution in Figure 2 is similar, except the value 04 is six times more likely in the observed distribution compared to the baseline distribution, although in both cases the values 02, 03 and 04 still as a whole contribute less than 3% of the distribution.

The challenge for detecting significant changes is that the distributions depend upon many factors, including the region, the season, the specific merchant, and the specific issuer and acquirer.

¹Data reflects all Visa programs except Interlink, PLUS, and commercial funds transfers in China As reported by member financial institutions globally and therefore may be subject to change.

Value	%	Value	%
00	76.94	00	76.94
01	21.60	01	20.67
02	0.99	02	0.90
03	0.27	03	0.25
04	0.20	04	1.24
Total	100.00	Total	100.00

Figure 2. The distribution on the left is the baseline distribution. The distribution on the right is the observed distribution. In this example, the value 04 is over 6x more likely in the observed distribution, although the two dominant values 00 and 01 still account for over 97% of the distribution.

3 Change Detection Using Data Cubes of Models (CDCM)

In this section, we describe a methodology called Change Detection using Cubes of Models or CDCM that is designed to detect changes in large, complex data sets.

3.1 Change Detection Models

Change detection models are a standard approach for detecting deviations from baselines [3].

We first describe the cumulative sum or CUSUM change detection algorithm [3]. We assume that we know the mean and variance of a distribution representing normal behavior, as well as the mean and variance of another distribution representing behavior that is not normal.

More explicitly, assume we have two Gaussian distributions with mean μ_i and variance σ_i^2 , $i = 0, 1$.

$$f_i(x) = \frac{1}{\sqrt{2\pi\mu_i}} \exp \frac{-(x - \mu_i)^2}{2\sigma_i}$$

The log odds ratio is then given by

$$g(x) = \log \frac{f_1(x)}{f_0(x)}.$$

and can now define a CUSUM algorithm as follows [3]:

$$Z_0 = 0.$$

$$Z_n = \max\{0, Z_{n-1} + g(x_n)\}.$$

An alert is issued where the Z_n exceeds a threshold.

Quite often the statistical distribution of the anomalous distribution is not known. In this case, if the change is reflected in the mean of the observations and the standard deviation is the same pre- and post-change, the generalized likelihood ratio or GLR algorithm can be used [3]:

$$G_k = \frac{1}{2\sigma^2} \max_{1 \leq j \leq k} \frac{1}{k-j+1} \left[\sum_{i=j}^k (x_i - \mu_0) \right]^2, \quad k > 1$$

where μ_0 is the mean of the normal distribution and σ is the standard deviation of the both the normal and abnormal distributions,

which are assumed to be Gaussian. Here k is fixed and determines the size of the window used to compute the score. Again, the detection procedure is to announce a change at the first up-crossing of a threshold by the GLR score.

3.2 Cubes of Models

The basic idea of the CDCM algorithm is that for each cell cell in a multi-dimensional data cube, we estimate a separate change detection model.

For the purposes here, we can define a *data cube* as usual, namely a multi-dimensional representation of data in which the cells contain measures (or facts) and the edges represent data *dimensions* which are used for reporting the data.

We define a *cube of models* as a data cube in which each cell is associated with a baseline model.

3.3 Learning and Scoring Change Detection Models

In our model, we assume that there are a stream of events, which in our case are transactions, and each event can be assigned to one or more cells in cube of models. For example, for each project, each transaction is assigned to the appropriate cell(s), as determined by one of six regions, by one of over 800 Merchant Category Codes (MCCs), by one of 8 terminal types, etc. We also assume that various derived data attributes are computed from the event data to form feature vectors, which are sometimes called profiles in this context. The profiles contain state information and derived data and are used as inputs to the models as usual.

Estimating baseline models. To learn the baseline models, we take a collection of event data and process it as follows.

1. First, we assign each event to one or more cells in the data cube as appropriate.
2. Second, we transform and aggregate the events and compute the required profiles for each cell using the event data.
3. Third, for each cell, we use the resulting profiles to estimate the parameters for the baseline model, and output the baseline model.

Scoring baseline models. To score a stream of event data, we proceed as follows.

1. First, we retrieve the appropriate XML file describing the segmentation.
2. Next, we assign each event to one or more cells in the data cube as appropriate.
3. We then access the profile associated with each cell, and update the profiles using the new event.
4. We then use the profile as the input to the appropriate baseline model and compute a score.
5. Next, we process the resulting score using the appropriate rules to determine whether an alert should be produced.
6. Next, we apply XSLT transformations to the score to produce a report describing the alert.

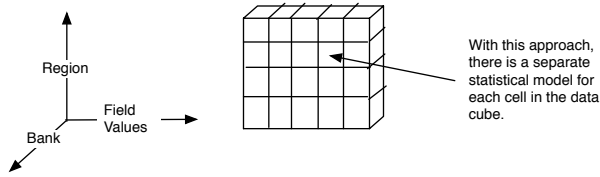


Figure 3. The basic idea with change detection using cubes of models or CDCM is that there is a separate change detection model for each cell in a multi-dimensional data cube. In the work described here we estimated and maintained over 15,000 different baseline statistical models and monitored them monthly.

7. Finally, if an alert is produced, we pass the alert to the required application or analyst.

4 Some Typical Alerts

In this section, we describe some typical alerts that have been generated by the CDCM system we developed. Currently, we compute alerts each month and re-estimate baselines several times a year. The system has been in operation for approximately two years.

It is important to remember when reading the case studies in this section that the issues identified by these alerts represent both a very small fraction of the transactions and a very small fraction of the total purchase dollars.

4.1 Dimensions of Cube

For the various alerts that we describe below, we used the following dimensions to define a data cube:

- The geographical region, specifically the US, Canada, Europe, Latin America, Asia Pacific, Middle East/Africa, and other.
- The field value or combination of values being monitored.
- The time period, for example monthly, weekly, daily, hourly, etc.
- The type of baseline report, for example a report focused on declines or a report describing the mixture of business for a merchant.

Today (January, 2007), for each of 324 field values times 7 regions times 1 time period times 3 report types, we estimate a separate baseline, which gives $324 \times 7 \times 1 \times 3 = 6816$. In addition, for 623 field values times 7 regions times 1 time period times 2 report types, we estimate a separate baseline, which gives an additional $623 \times 7 \times 1 \times 2 = 8726$ separate baseline models. So in total, we are currently estimating 15,542 ($=6816+8726$).

Actually, the description above is a simplified version of what actually takes place. For example, the 6816 baselines mentioned arise from $324 \times 7 = 2272$ different field values, but the 2272 different field values are not spread uniformly across the 7 regions as indicated, although the total is correct.

We are in the process of doubling the number of field values and increasing the number of time periods, so shortly we will be estimating and monitoring over 40,000 separate baselines.

4.2 Incorrect Merchant Category Code

In this example, an airline was coding some of its transactions using a Merchant Category Code (MCC) B instead of the preferred MCC A, which lowered the approval rate for the transactions. This resulted in a baseline alert that in turn resulted in a manual investigation by an analyst. Following this, a conference call was arranged with the individual responsible for the relationship with the bank who was the acquiring bank for the airline. As a result of this call, a fix was installed by the airline. This fix resulted in an annual recoverable costs of over \$300,000.

4.3 Testing of Counterfeit Cards

In this example, the decline rate for a large bank was essentially the same month to month but the baseline model identified a particular category of transactions (specified by a combination of five fields) for which the decline rate sharply peaked in September 2006 compared to an earlier baseline period. One way of thinking about this, is that for this bank, most of the 50,000+ or so baselines were normal for September, but one was not. When investigated, this particular baseline was elevated due to suspected testing of stolen/counterfeit cards to determine whether they were still valid. After several discussions and further investigation, changes were made so that this type of testing was not possible.

4.4 Incorrect Use of Merchant City Name

In this example, a European merchant's transactions were coded incorrectly so that the merchant city name field contained incorrect information. This resulted in a lower acceptance rate for the transactions from this merchant. This lower acceptance rate was picked up by a baseline alert that for each acquirer monitored decline levels for each MCC. After an investigation, the merchant corrected the problem.

5 Discussion and Lessons Learned

Lesson 1. The most important lesson we learned was that thus far it has been more fruitful to examine many individual baseline and change detection models, one for each different segment of the event stream, even if these models are very simple, than to build a single, relatively complex model and apply it to the entire event stream.

Lesson 2. The time and effort required to get the alert format right is substantial. Although it was certainly expected, the business return on the project was dependent to a large degree on the ability to deliver to the analysts information in a format that they could readily use. After quite a bit of experimentation, a report format was developed that reported:

- What is the issue? This part of the report identifies the relevant business unit and the relevant business issue.
- Who has the issue? This part of the report identifies the relevant subsystem of VisaNet, the relevant attribute, and the relevant attribute value.
- What is the business opportunity? This part of the report identifies the daily business value associated with the issue and the statistical significance of the alert (Low, Medium High).
- What is the business impact? This part of the report describes the a business measure as currently measured, the historical

measure of the business measure during the baseline period, and the number of transaction affected.

- The final part of the report contains additional information, such as the alert ID, alert creation date, whether the alert is new, and whether the alert is associated with an issue that has been previously identified and now is being monitored for compliance.

One way to summarize the report is that the items in alerts gradually changed from those items related to the statistical models and how the alerts were generated to items directly related to how the alerts were investigated and how the business impact was estimated. The surprise for us was not that this transition had to be made, but rather the time and effort required to get it right.

Lesson 3. It turned out that some of the most important alerts we found were alerts that had low statistical significance. For each report, we include an estimate of the statistical significance of the alert (low, medium, high and very high) as well as an estimate of the business significance of the alert (in dollars). It turned out that after investigation, the alerts that generated by most dollars saved, were often the alerts with low statistical significance. For this reason, it was usually not a good idea to investigate alerts in the order of most statistically significant to least statistically significant. Rather, the analysts used a more complex prioritization that thus far we have not tried to formalize.

Lesson 4. As a result of analysis of an alert, it was sometimes possible to create specialized baselines and reports that would look for similar problems in the future. We quickly learned that even a few specialized reports like this could easily occupy most of our available time. The lesson we learned was that it was important to devote some of our time to looking for new opportunities (think of this as a survey activity), since some of these turned out to be even more important than what we were currently doing.

6 Related Work

There is a large amount of research on change detection algorithms per se. The monograph by Basseville and Nikiforov [3] is a good summary of this research field. In particular, the change detection algorithms that we use here, including CUSUMs, Generalized Likelihood Ratios, and related algorithms are covered in this reference.

The work described in this paper differs from classical change detection and contingency tables in that it uses a separate change detection model for each cell in a cube of models.

More recently, Ben-David, Gehrke and Kifer [4] introduced a non-parametric change detection algorithm that is designed for streams. The methods used here are parametric. In contrast to their approach which uses a single change detection model, we build a large number of models in order to handle complex, heterogeneous data, one for each cell in a multi-dimensional data cube.

The paper by Fawcett and Provost [6] has a similar goal — detecting unusual activity in large data sets — but uses a much different approach. Their approach is to introduce operating characteristic style measures in order to identify unusual behavior.

Guralnik and Srivastava [8] study event detection in time series data by using a new change detection algorithm they introduce, which involves iteratively deciding whether to split a time series interval to look for further changes.

In contrast to all these methods, our approach is to use relatively simple classical change detection algorithms, such as CUSUM and GLR, but to build thousands of them, one for each cell in a multi-dimensional data cube. As far as we are aware of, our paper is also one of the few papers in the data mining literature that presents a case study of change detection involving a system as large and heterogeneous as VisaNet.

7 Summary and Conclusion

In this paper, we have shared our experiences and some of the lessons learned over the past two years developing and operating a baseline and change detection system for Visa. Because of the complex and highly heterogeneous nature of Visa's transactional data, we did not build a single change detection model, but rather over 15,000 individual change detection models. Indeed we built a separate change detection model for each cell in a multi-dimensional data cube. This is an example of we have been calling Change Detection using Cubes of Models or CDCM.

Overall, the approach seems to work quite well. Indeed, substantial business value is being generated using this methodology, and thus far we have not been able to achieve the same performance using a single baseline or change detection model.

To summarize, we have demonstrated through this case study that change detection using data cubes of models (CDCM) is an effective framework for computing changes on large, complex data sets.

8 References

- [1] Alan Agresti, *An Introduction to Categorical Data Analysis*, John Wiley and Sons, Inc., New York, 1996.
- [2] The Augustus open source data mining system can be downloaded from www.sourceforge.net/projects/augustus.
- [3] M. Basseville and I. V. Nikiforov. *Detection of Abrupt Changes: Theory and Application*. Prentice Hall, 1993.
- [4] Shai Ben-David, Johannes Gehrke, Daniel Kifer, *Detecting Change in Data Streams*, Proceedings of 2004 VLDB Conference, 2004.
- [5] Joseph Bugajski, Robert Grossman, Eric Sumner, Tao Zhang, *A Methodology for Establishing Information Quality Baselines for Complex, Distributed Systems*, 10th International Conference on Information Quality (ICIQ), 2005.
- [6] Tom Fawcett and Foster Provost, *Activity monitoring: noticing interesting changes in behavior*, KDD '99: Proceedings of the fifth ACM SIGKDD international conference on Knowledge discovery and data mining, 53–62, ACM Press, New York, 1999.
- [7] Robert L. Grossman, *PMML Models for Detecting Changes*, Proceedings of the KDD-05 Workshop on Data Mining Standards, Services and Platforms (DM-SSP 05), ACM Press, New York, 2005, pages 6-15.
- [8] Valery Guralnik and Jaideep Srivastava, *Event detection from time series data*, Proceedings of the Fifth ACM SIGKDD international Conference on Knowledge Discovery and Data Mining, ACM Press, New York, NY, 33-42, 1999.