

KDD-Cup 2000 Organizers' Report: Peeling the Onion

Ron Kohavi
Blue Martini Software
2600 Campus Dr.
San Mateo, CA 94403

ronnyk@bluemartini.com

Carla E. Brodley
School of Electrical and Computer Engineering
Purdue University
West Lafayette, IN 47907

brodley@ecn.purdue.edu

Brian Frasca
Blue Martini Software
2600 Campus Dr.
San Mateo, CA 94403

brianf@bluemartini.com

Llew Mason
Blue Martini Software
2600 Campus Dr.
San Mateo, CA 94403

lmason@bluemartini.com

Zijian Zheng
Blue Martini Software
2600 Campus Dr.
San Mateo, CA 94403

zijian@bluemartini.com

ABSTRACT

We describe KDD-Cup 2000, the yearly competition in data mining. For the first time the Cup included insight problems in addition to prediction problems, thus posing new challenges in both the knowledge discovery and the evaluation criteria, and highlighting the need to “peel the onion” and drill deeper into the reasons for the initial patterns found. We chronicle the data generation phase starting from the collection at the site through its conversion to a star schema in a warehouse through data cleansing, data obfuscation for privacy protection, and data aggregation. We describe the information given to the participants, including the questions, site structure, the marketing calendar, and the data schema. Finally, we discuss interesting insights, common mistakes, and lessons learned. Three winners were announced and they describe their own experiences and lessons in the pages following this paper.

Keywords

KDD-Cup, e-commerce, competition, data mining, real-world data, insight, data cleansing, peeling the onion, best practices.

1. INTRODUCTION

The KDD-Cup is a yearly competition in data mining that started in 1997. KDD-Cup 2000, the fourth competition, involved multiple problems, following the suggestions of previous organizers [1]. For the first time, the Cup included insight questions in addition to prediction problems.

The domain for the KDD-Cup was e-commerce, considered a “killer domain” for data mining because it contains all the ingredients necessary for successful data mining [2]. The ingredients include (i) wide records (many attributes), (ii) many records (large volume of data), (iii) controlled data collection (e.g., electronic collection), (iv) ability to evaluate results and

demonstrate return on investment, and (v) a domain where action can easily be taken (e.g., change the site, offer cross-sells). Blue Martini Software approached several clients using its Customer Interaction System to volunteer their data, and a small dot-com company called Gazelle.com, a legwear and legcare retailer, agreed to volunteer their data, properly sanitized.

After studying the data and consulting with Gazelle.com and retail experts at Blue Martini Software, five questions were defined. Two questions were prediction questions while the remaining three were insight questions. Only a portion of the available data was made available to competitors (about the first two months) while a test-set (the third month) was kept for evaluation, in line with standard best practices of having a separate test set.

To make the problem more realistic, we collected background information from Gazelle.com and made their marketing calendar available to competitors. The events (e.g., a TV advertisement) help explain the changes in the number of visitors over time.

Data was made available to in two formats: original data and aggregated data. While the original data was collected at the page request level, the questions were at the session and customer level. Because most tools do not have sufficiently powerful aggregation capabilities, we used the Blue Martini Customer Interaction System to generate the aggregated data, summarizing session-level and customer-level behavior. Further details about the data and aggregations are provided in Section 4.

The evaluation of the insight questions was done in consultation with Blue Martini's retail experts. We created a standardized scoring mechanism described in Section 3. As we evaluated the submissions whose statistics can be found in Section 5, we found many observations that were “shallow,” i.e., they involved patterns that did not lead to deep understanding of the issues. We would like to highlight the need for “peeling the onion” when doing data mining investigations. Results and insights are described in Section 6.

We conclude the paper with lessons learned. Also in this issue are three reports from the winners of the competition.

2. BACKGROUND INFORMATION

It is helpful to know the following background information about the Gazelle.com webstore:

- The home page contained more than 70 images. This made downloads extremely slow for modem-based visitors.
- As with many dot-coms, Gazelle.com's initial goal was to attract customers, even if it meant losing money in the short term. They had many promotions that are relevant for mining, because promotions affected traffic to the site, the type of customers, etc. The important promotions were
 - FREE - Free shipping (\$3.95 value). Active from March 20 to April 30 (shipping was normally free if sale was above \$40).
 - MARCH1 - \$10 off from March 1 to April 1.
 - FRIEND - \$10 off from March 1 to April 30.
 - FREEBAG - A free bag from March 30 to April 30.

Note that both the MARCH1 and FRIEND promotions offered \$10 off. They were used for different purposes, and were run with different promotion codes.

- Gazelle.com ran a TV advertisement during a prime-time episode of the popular comedy show, Ally McBeal, on February 28.
- Gazelle.com changed their registration form significantly on February 26, so some customer attributes were only collected prior to this date and some were collected only after this date.

3. THE QUESTIONS AND EVALUATION CRITERIA

There were five independent questions for KDD-Cup 2000. Two of the questions were standard prediction problems with objective evaluation criteria, while the remaining three were subjective "insight" questions.

Question 1

Given a set of page views, will the visitor view another page on the site or will the visitor leave?

This question was motivated by the idea that knowing whether a visitor is likely to leave can help determine the "best" page to display (e.g., special promotions could be shown to encourage the visitor to stay). The evaluation criterion for this question was simply the number of correct predictions on the test set. The winner was the entry with the highest accuracy.

Question 2

Given an initial set of page views, which product brand ("Hanes", "Donna Karen", "American Essentials", or "Other") will the visitor view in the remainder of the session?

This question was motivated by the problem of improving navigation by automatically placing a hyperlink on the current page pointing to a particular brand page. To make the problem more manageable, we restricted the task to predicting one of three most commonly sold brands, or "Other" (defined as not viewing any of the three brands in the remainder of the session). The

evaluation criterion was a weighted prediction score where points were awarded as follows:

2 points: If they predicted one of the three specific brands and one of the remaining pages in the session included the predicted brand.

1 point: If they predicted "Other" and none of the remaining pages in the session included a visit to one of the three specific brands.

0 points: All other cases.

The winner was the entry with the highest score.

For the remaining three questions, the competitors were required to submit text and graphs that a business user would be able to understand and find useful. Each submission was limited to 1,000 words and ten graphs.

Question 3

Given a set of purchases over a period of time, characterize visitors who spend more than \$12 on an average order at the site.

The motivation for this question was that insight about a website's more valuable customers could be useful for determining marketing directions, product selection, etc.

Question 4

This was the same as Question 1, but the goal was to provide insight, rather than predict accurately.

Question 5

This was the same as Question 2, but the goal was to provide insight, rather than predict accurately.

For Questions 3, 4, and 5, for which no simple objective measure existed, we talked to retail experts at Gazelle.com and Blue Martini Software about the submissions. We then formalized the evaluations by collecting all of the significant insights, weighting them, and creating a combined score based on the insights found, the correctness of the submission, and the presentation of their submission (keeping in mind that *business users* were the target audience). The actual formula used for computing an entrant's score was

$$Score = 3P + 3C + \sum_{i=1}^N w_i I_i$$

where P is the entrant's presentation score (0-10), C is the entrant's correctness score (0-10), and for each insight i , w_i is the weight assigned to the insight and I_i is the entrant's score for the insight (0-2). The number of insights and their weights varied for each question.

The presentation score captured the effectiveness of presentation of the entrant's submission. This included factors like:

- How readable and easy to understand was the submission?
- Were there graphs, tables, and figures that business people could understand?
- Was there an effort to distill the important information, or was too much irrelevant information presented?

The correctness score was based on whether the entrant's claims were correct and whether the claims had sufficient data to support them.

For each question, we defined a complete set of insights based on all of the insights provided by every competitor. These insights were weighted to reflect how interesting they would be to a business user (based on conversations with retail experts from Gazelle.com and Blue Martini Software). Many insights were given low weight (and sometimes even zero weight) because they simply correlated with more fundamental insights. For each entry, every insight was awarded an insight score which was either zero (if they didn't discover the insight), one (if they partially described the insight) or two (if they fully described the insight). Due to the large number of insights (over 30 each for Questions 3 and 4), we do not include a list here. A complete list of insights with detailed explanations and weights can be found on the KDD-Cup 2000 home page [4].

4. THE DATA

In this section, we describe what data was collected in the webstore, how we generated the initial star schema for the data warehouse, what types of data cleansing/obfuscating were performed, and which data transformations were applied. Finally, we summarize the final schema and data formats provided for the KDD-Cup.

4.1 Initial Data Collection

Gazelle.com went live with Blue Martini's Customer Interaction System (CIS) on January 30, 2000 with soft-launch to friends and families. On the webstore, an application server in the Blue Martini architecture generates web pages from Java based templates. Among the other things, the architecture logs customer transactions and clickstreams at the application server layer. Since the application server generates the content (e.g., images, products and articles), it has detailed knowledge of the content being served. This is true even when the content is dynamically generated or encrypted for transmission commonly used for checkout. Weblog data is not needed. Application servers use cookies (or URL encoding in the absence of cookies) to keep track of a user's session, so there is no need for "sessionizing" clickstreams as there is for standard weblogs. Since the application server also keeps track of users using login mechanisms or cookies, it is easy to associate individual page views with a particular visitor.

Among the data collected by the Blue Martini application server, the following three categories are related to this KDD-Cup:

- Customer information, which includes customer ID, registration information, and registration form questionnaire responses.
- Order information at two levels of granularity: 1) Order header, which includes date/time, discount, tax, total amount, payment, shipping, status, and session ID; and 2) Order line, which includes quantity, price, product, date/time, assortment, and status.
- Clickstream information at two levels of granularity: 1) Session, which includes starting and ending date/time, cookie, browser, referrer, visit count, and user agent; and 2) Page view, which includes date/time, sequence number, URL, processing time, product, and assortment.

In general, each customer can have multiple sessions. Each session can have multiple page views and multiple orders. Each order can have multiple order lines. Each order line is a purchase record of one product with a quantity of one or more.

4.2 Star Schema Creation

The data collector in the Blue Martini application server is implemented within an On-Line Transaction Processing (OLTP) system. OLTP systems are designed for efficient handling of a large number of small updates and short queries. This is critical for running an e-commerce business, but is not appropriate for analysis, which usually requires full scans of several very large tables and a star schema¹ design [7][8] which business users can understand. For data mining, we need to build a data warehouse using dimensional modeling techniques. Both the data warehouse design and the data transfer from the OLTP system to the data warehouse system are very complex and time-consuming tasks. Because Blue Martini's architecture contains metadata about tables, columns, and their relationships, it can automatically construct the data warehouse from the OLTP system [6].

When preparing the data for the KDD-Cup, we integrated syndicated data from Acxiom into the schemas, which enriched the customer information for analysis by introducing more than fifty new attributes such as Gender, Occupation, Age, Marital Status, Estimated Income, and Home Market Value.

Two star schemas used for generating the KDD-Cup data are the Clickstream star and the Order Lines star. The Clickstream star consists of one fact table: "Clickstream" and six dimension tables: "Customer Profiles", "Acxiom", "Web Sessions", "Products", "Assortments", and "Contents". The Order Lines star consists of the one fact table: "Order Lines" and six dimension tables: "Customer Profiles", "Acxiom", "Order Headers", "Products", "Assortments", and "Promotions".

4.3 Data Cleansing/Obfuscating

To protect customer privacy, we removed attributes containing information about individuals such as Login Name, Password, Credit Card, Customer Name, and Session IP Address. We also removed attributes containing profit-related information such as Product Unit Cost. For attributes that we believe are important for mining this data (solving the KDD-Cup questions), we scrambled the data. For example, the values of the Email attribute were mapped to keep only the domain suffix such as COM, EDU, ORG, and GOV. In addition, we kept "Gazelle.com" for email addresses with the suffix of gazelle.com. All company names were mapped to "COMPANY" and a number, so that it is possible to tell people are from the same company without knowing which company it is. Session Cookie IDs were encoded, so that each Cookie ID appears as a different number, while it is still possible to determine that several sessions are from the same cookie.

Data cleansing is usually a part of the KDD process. We chose to do some initial data cleansing ourselves for three reasons. Firstly, unlike a real data mining project, the participants of the KDD-Cup did not have direct contact with the domain experts. Secondly, data obfuscating must be done before releasing the data, and thirdly, the questions are challenging enough even after this initial data cleansing. To clean the data, we

- Removed Keynote records. Keynote hit the Gazelle.com home page 3 times a minute, 24 hours a day, 7 days a

¹ A star schema is a join of database tables with one central fact table joined to several other tables (called dimensions).

week, generating about 125,000 sessions per month. These records can skew mining results.

- Removed test users. We used criteria such as with “test” in customer names or purchased using a credit card that were used by more than 15 different users. Note that the test users have very different purchasing and browsing behaviors.
- Removed returned and uncompleted orders. The number of these orders is small, but they may cause confusion.

4.4 Data Transformations

We provided two types of data for the KDD-Cup questions, namely unaggregated and aggregated.

The data transformation for the unaggregated data is very simple. Questions 1, 2, 4, and 5 share the same unaggregated dataset. It is a flat table created by joining the Clickstream star. In this table, each record is a page view. Session attributes are repeated multiple times if the session has multiple page views. Similarly, customer information is also repeated in the table. To define the targets for these four questions, we added three Boolean attributes in the table as follows. Three example sessions are given in Table 1, showing how the sessions were clipped.

- “Question 1 Test Set” indicating whether you will see this page view if the session is in the test set for Questions 1 and 4. This is defined based on a clipping point in half of the randomly selected sessions. For a selected clipping session, we randomly generated a clipping point between one and the session length minus one. No clipping was performed for sessions of length one.
- “Question 2 Test Set” indicating whether you will see this page view if the session is in the test set for Questions 2 and 5. This is defined based on a clipping point in all the sessions. The clipping point is generated in the same way as for Question 1.
- “Session Continues” as the target of Questions 1 and 4.

Session ID	Request Sequence	Question 1 Test Set	Question 2 Test Set	Session Continues
29	1	T	T	F
29	2	T	T	F
29	3	T	F	F
56	1	T	T	T
56	2	T	T	T
56	3	F	F	T
68	1	T	T	F

Table 1: How sessions got clipped.

The unaggregated dataset for Question 3 is also a flat table created by joining the Order Lines star. Each order line is a record in the table. Attribute values for order headers and customers may repeat multiple times. A Boolean attribute “Spend Over \$12 Per Order On Average” is added to the table as the target. This attribute is defined at the customer level.

These two unaggregated datasets contain the raw data, providing enough information for those people with data transformation ability to do the data mining. Note that the first dataset does not

contain the order information while the second dataset does not contain the clickstream information. Participants could join them together if they thought doing so could help them to solve the questions.

Considering that many researchers, especially those working on data mining algorithms, do not have software readily available to transform (including aggregate) the raw data, we provided an aggregated version of the data. The aggregated data consists of three datasets: one for Questions 1 and 4, one for Questions 2 and 5, and the other for Question 3. These datasets are derived by aggregating the two unaggregated datasets to the level of granularity appropriate for mining. That is, the session level for Questions 1, 2, 4, and 5 and the customer level for Question 3. At the same time, we added new attributes based on examination of existing attributes. For example, we extracted the session browser family names, the browser names, and the top three browser family names. In the two aggregated datasets for Questions 1, 2, 4, and 5, each session is a single record. During the generation of these two datasets, all page views marked “not in the corresponding test sets” in the unaggregated datasets were removed before the aggregation operation. In the aggregated dataset for Question 3, each customer is a single record.

The aggregation operations generated 151 and 153 new attributes for Questions 1 & 4 and Questions 2 & 5, respectively. Examples include the number of views of individual top products which were selected based on the statistics of the datasets, the number of views of assortments, the number of views of different templates, and information about the last page, which includes information appearing on it and its date/time information. For questions 2 and 5, we defined three numeric attributes indicating the number of views of the respective brands (Hanes, Donna Karan, American Essentials) in the remainder of the session. In addition, we also defined a Boolean attribute that was set to true if none of the brands were viewed in the remainder of the session and false otherwise.

When generating the aggregated dataset for Question 3, we joined clickstream data to the order lines data since we believed that clickstream can help to answer Question 3 and it is hard to join them after aggregation. The aggregation for this dataset was carried out at two levels: first to the session level and then to the customer level, generating 434 new attributes in total such as “Average Session Request Count”, “First Session First Referrer Top 5”, and “Percent of Products Purchased on Sunday”.

4.5 Final Data Schema and Formats

The datasets were released in flat files using C5 format (www.rulequest.com), a widely used data format for data mining. There was no training/test split for Question 3 data, as it was a pure insight question. Questions 1 and 2 had training and test datasets. The training datasets contain the target information while the test datasets do not. To avoid leaks (with respect to targets), we did the training/test splits using time. The data we got from Gazelle.com was collected from January 30, 2000 to April 30, 2000 (3 months). We used the data before April 1, 2000 (2 months) for training for all of the questions. Since Questions 1 and 2 share information, their test sets could not overlap. We used the data after April 14, 2000 (half a month) as the test set for Question 1, and the data from April 1, 2000 to April 14, 2000 (half a month) as the test set for Question 2.

Table 2 summarizes the number of attributes and the number of records in the datasets. Questions 4 and 5 do not appear in the table because Question 4 used the same data as Question 1, while Question 5 used the same data as Question 2. It is worth mentioning that for Question 2, we had four target attributes in the training set, and only one dummy target attribute in the test set.

Question	Training set		Test set	
	Attributes	Records	Attributes	Records
1: Unaggregated	217	777,480	215	164,364
2: Unaggregated	217	777,480	215	142,204
3: Unaggregated	232	3,465	-	-
1: Aggregated	296	234,954	296	50,558
2: Aggregated	299	234,954	296	62,913
3: Aggregated	518	1,781	-	-

Table 2: Dataset statistics.

5. SUBMISSION STATISTICS

We received 170 non-disclosure agreements requesting access to the data. Of these, there were 31 participants who submitted an entry for one or more of the questions. The number of entries we received for each question is shown in Figure 1. Since the competition ended, we have received more than 190 additional click-through agreements for access to the data for research or educational purposes.

Each participant was also required to submit a questionnaire answering questions about their efforts. This included questions on the resources they utilized (e.g., the number of people involved and the number of hours spent in each phase of analysis), the software and hardware that they used, and the data mining techniques and methodologies used in both processing and analyzing the data. The statistics presented in this section are based on the answers we received in these questionnaires.

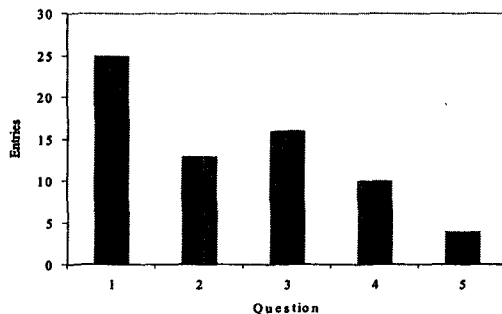


Figure 1: Number of entries for each question.

In total, the 31 participants spent 6,129 person-hours preparing and analyzing the data. This corresponds to about 200 person-hours per participant. One participant spent more than 900 person-hours on their submission. The number of people involved varied from one to thirteen, although most entries came from teams of two or three people. The breakdown of how the hours were spent on average is shown in Figure 2.

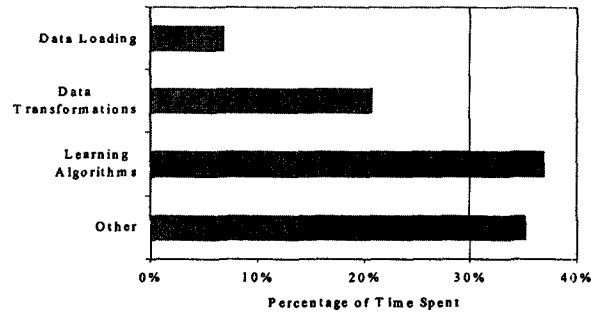


Figure 2: Average time spent on each phase of analysis.

Notice that, in contrast to most studies [5], less than 30% of the time was spent in data loading and transformations. Most likely, this was due to two factors. Firstly, the data was collected within Blue Martini's integrated e-commerce system designed for data mining, and thus was in a form more amenable to analysis [6]. Secondly, as described in Section 4, we spent significant time aiding the contestants by transforming the data and constructing new features for use in analysis.

The breakdown of data mining software origin used by participants is shown in Figure 3. One interesting trend to note is the increase in the use of commercial software for the KDD-Cup: the proportion of entries using commercial or proprietary software has grown from 44% (KDD-Cup 1997) to 52% (KDD-Cup 1998) to 77% (KDD-Cup 2000).

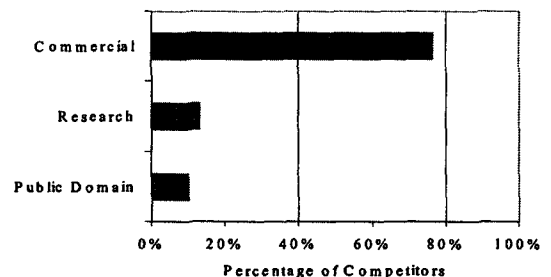


Figure 3: Type of software used by the competitors.

The operating system used by competitors was an even mix of Microsoft Windows (54%) and Unix (46%). Of those competitors using Unix, various flavors of commercial Unix accounted for 65%, while Linux accounted for the remaining 35%. Despite the balance between Microsoft Windows and Unix operating systems, the hardware used was primarily desktop PCs (73%), rather than Unix workstations (27%).

For data access, 32% of competitors used a database, while 68% used flat files. The breakdown of data processing tools used by competitors is shown in Figure 4. From this figure it can be seen that most competitors made use of the data processing tools built into their analysis software rather than developing proprietary data processing tools for the KDD-Cup.

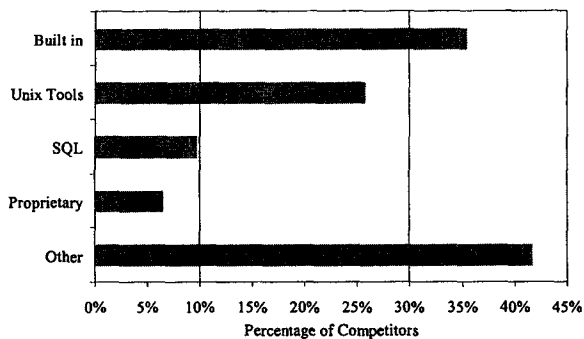


Figure 4: Data processing tools used.

As mentioned in Section 4, we provided both aggregated and unaggregated data. The majority of competitors used the aggregated data (59%) rather than the unaggregated (41%). This suggests that many data mining tools provide only limited support for data aggregation.

Figure 5 shows the top algorithmic techniques used by the competitors. The figure shows both the percentage of competitors who tried that algorithm and the percentage of competitors who submitted a solution to at least one question using that algorithm. As can be seen, decision trees were by far the most popular choice, with more than 50% of the competitors submitting a solution to at least one question using decision trees.

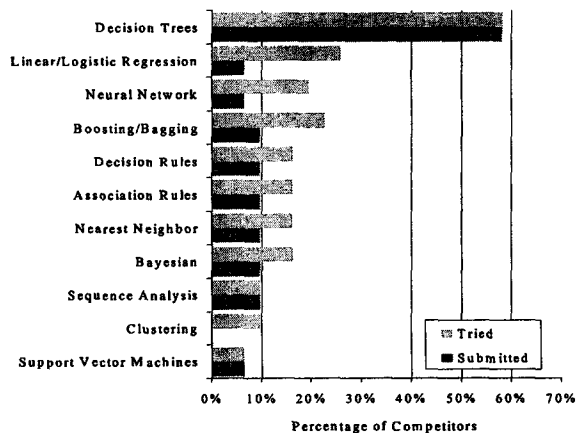


Figure 5: Algorithms tried versus submitted.

6. RESULTS AND INSIGHTS

In this section we present the results for each of the five questions.

Participants discovered many important actionable insights, including which referrers resulted in heavy spending, which pages cause abandonment, and what segments of the population are heavy spenders. Many seemingly interesting insights were obvious once one discovered the underlying cause, which was usually related to time or session length. For example, many participants noted a correlation between being a heavy spender and a visitor's answer as to whether they would like email from Gazelle.com. When this response is plotted against time, it is easy to see that it varies dramatically -- this is because Gazelle changed the default for this field twice. Predicting who would leave the site was made particularly challenging because many sessions

were of length one -- in this data web crawlers that viewed a single page in each session accounted for 16% of sessions. Despite this, surprisingly few participants identified which visitors were actually web crawlers rather than real people. In examining the results when shorter sessions were removed, we noted that was possible to predict accurately when the prediction confidence was high.

For Question 1 (Given a set of page views will the visitor view another page on the site or will the visitor leave), the accuracy values ranged from 77.06% to 59.56% with a mean of 73.14%. The difference between the top two performers was only 0.10%, which translates into 50 sessions. In fact, the difference in accuracy of the top five participants was statistically insignificant (a 95% confidence interval corresponds to $\pm 0.37\%$). Despite this result, if we restrict the evaluation to predicting sessions with five or more page views the results are far more significant (the difference between first and second place was 1.5% and a 95% confidence interval corresponds to $\pm 0.79\%$). Figure 6 shows that the gains charts for the top two participants track the optimal gain for 10% of these longer sessions, which account for 43% of the target. The optimal gain is shown by the leftmost curve on the graph.

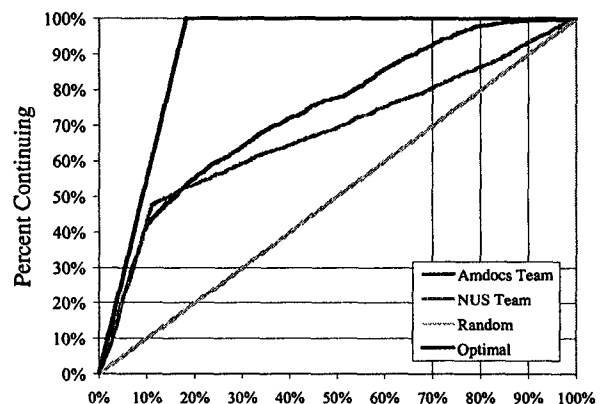


Figure 6: Cumulative gains chart for sessions with five or more page views.

Question 4 was the insight question corresponding to Question 1. Some of the key insights found were that web crawlers and gazelle testers leave and that the length of stay depends on the referrer site (users from Mycoupons had longer sessions, whereas users from ShopNow tended to leave quickly). Participants noted that a returning user's probability of continuing was double that of a first time visitor. Viewing some specific products caused users to leave the site. This is an example of an actionable insight, in that the web site might consider removing those products. Another actionable insight is that 32% of customers left after entering the replenishment section of the site. Many "discoveries" were explained by noticing that the probability of leaving decreases with the number of pages viewed in the session. For example, the insight that "viewing many different products in a session implies low abandonment" is explained by this fact.

For Question 2 (Given a set of page views which product brand will the visitor view in the remainder of the session), the scores ranged from 60956 to 60697 with a mean of 60814.8. Like Question 1, we found the difference between the top participants

to be statistically insignificant. However, like Question 1, we observed very good lift curves when we restricted our evaluation to sessions with five or more page views. One of the best predictors was the referrer URL: Fashionmall and Winnie-Cooper are good referrers for Hanes and Donna Karan, whereas Mycoupons, Tripod, and Deal-finder are good referrers for American Essential. When we look more closely at this result we see that the American Essentials brand primarily contains socks, a low priced item which often falls under the \$10 coupon price. Very few participants realized that the Donna Karan brand was only available starting February 26.

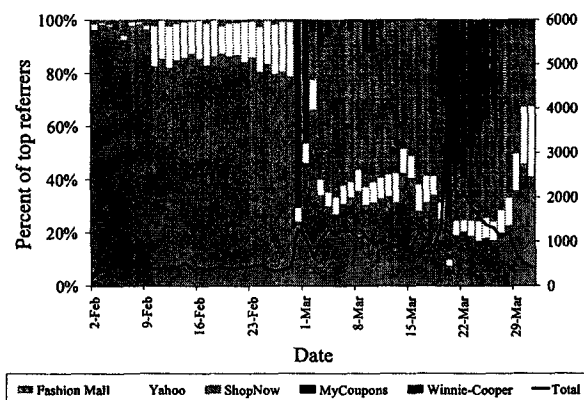


Figure 7: Top referrers by date.

For Question 3 (Characterize visitors who spend more than \$12 on an average order at the site) many interesting insights were simply related to time. For example, noting that significant activity began on February 29th, when the TV ad, Friends promotion and hard launch occurred. Another example is that the referring site traffic changed dramatically over time (see Figure 7). Some of the deeper insights that arose from this observation were related to the conversion rate. While the overall conversion rate for the site was only 0.8%, Mycoupons had an 8.2% conversion rate, but generated low spenders. On the other hand, the conversion rate from Fashionmall and ShopNow was only 0.07% even though they brought 35,000 visitors. Some of the other factors correlating with heavy purchasers were:

- They were not an AOL customer (the Gazelle.com site displayed badly within the AOL browser window).
- They came to the site after seeing a print advertisement.
- They had either a very high or very low income.
- They were living in the Northeastern U.S.

7. LESSONS LEARNED

The KDD-Cup is a great tool for highlighting both to the data mining research community and to the users of data mining the issues faced by the participants and by the organizers. We now describe the main lessons learned.

The most important lesson is that humans are an important part of the KDD process, even if the only interesting measurement is accuracy or score (Questions 1 and 2). The fully automated tools can never match the human insight that allowed the winners to create multi-stage models (see the KDD-Cup 2000: Winner's Reports in this issue), identify crawlers, local testers, and construct additional features. The importance of human understanding was also apparent in the choice of algorithms tried

versus submitted: Decision Trees were used the most often and submitted the most often while Neural Networks, Logistic Regression, and Clustering had the worst try-to-submit ratios. Many participants who thought they found an interesting result did not spend the time to "peel the onion" and find the true underlying causes. For the insight questions, the iterative process is even more important because many of the initial correlations are obvious and not interesting to the business users (e.g., those who purchase an item that costs over the heavy-spender threshold of \$12 are indeed heavy spenders). Many insights that seemed interesting had extremely low support. For example, several participants claimed that all purchasers who came from Shopnow were heavy spenders. While the statement was true, the support was six people! With the human involvement required, it takes time to derive useful insight---hundreds of hours.

The changes to the site created interesting effects and biases. Those that ignored the documentation about special marketing events did not do well. Time is a crucial attribute and changes to the site and products needs to be taken into account. In one case, a competitor claimed that the problem was "too real" and that we should simplify it. Our questions were hard, but they represent real-world problems in a real-world setting. The results showed significant lift, especially on longer sessions, and many insights were extremely interesting and actionable. For many one-click sessions, it was impossible to predict well, but for when the confidence was high, especially on longer sessions, predictions were very good.

The data was collected through the Blue Martini application server, and avoids the use of standard weblogs and allows correlating purchases to clickstreams. The data collector also saves information about the products shown in addition to URLs, making information more stable across site changes. Such data was rich and easier to work with, and the addition of Acxiom attributes certainly helped in deriving insights. Even with all these advantages over weblogs, identifying crawlers and test users remains a hard problem

For future organizers of the KDD-Cup, we offer some suggestions. Before volunteering to organize the KDD-Cup, make sure you understand the amount of effort involved. We estimated that we spent a total of 800 hours on getting the data, cleansing it, obfuscating it, transforming it, setting the web pages, working on the legal agreements, and evaluating the results. Plan on spending significant time in thinking about data obfuscation and identifying "leaks" in the data (giveaway attributes that predict the target because they're downstream in the collection process). For example, our system stored the session length, an attribute that we had to recompute after clipping the data, or else it would giveaway the target. We were very careful about removing leaks this time, having seen problems in previous years, but we still had to re-release the data twice in the initial phase due to mistakes in randomization and cookie obfuscation. We spent significant time writing the introductory material, giving the background knowledge, explaining the columns, yet we still had to develop a FAQ, which had 67 questions at the end of the competition. We gave the participants two question periods, one right after we released the data, and one before submission. We believe this was useful to get people started and also allowed us to plan our time better. The evaluation took a very long time, especially creating the weighted list of insights and validating the insights. We asked the participants to write a report for business

users, but after reading the reports we suspect that many of the authors have never talked to a business user. On the bright side, we learned many things about the data that we did not know and we saw some excellent methods to present results and make them more accessible to a wider audience.

8. ACKNOWLEDGMENTS

We thank Gazelle.com for providing the data. We thank the Axiom Corporation for providing the syndicated data overlay. Catharine Harding and Vahe Katros, our retail experts at Blue Martini Software, were helpful in reviewing submissions and explaining some of the patterns. Sean MacArthur from Purdue University helped write the scoring code.

REFERENCES

- [1] Ismail Parsa, KDD-Cup-97: A Knowledge Discovery and Data Mining Tools Competition talk. Newport beach, CA, USA. <http://www-aig.jpl.nasa.gov/public/kdd97/>.
- [2] Ron Kohavi and Foster Provost, Applications of Data Mining to E-commerce (editorial), Special issue of the International Journal on Data Mining and Knowledge Discovery, Jan 2001. <http://xxx.lanl.gov/abs/cs.LG/0010006>.
- [3] ACM SIGKDD Explorations homepage, the information for authors link: <http://www.acm.org/sigkdd/explorations/>.
- [4] KDD-Cup 2000 homepage. Carla Brodley and Ron Kohavi. <http://www.ecn.purdue.edu/KDDCUP/>.
- [5] Gregory Piatetsky-Shapiro, Ron Brachman, Tom Khabaza, Willi Kloezen, and Evangelos Simoudis, An Overview of Issues in Developing Industrial Data Mining and Knowledge Discovery Applications, *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining*, 1996.
- [6] Suhail Ansari, Ron Kohavi, Llew Mason, and Zijian Zheng. Integrating E-Commerce and Data Mining: Architecture and Challenges. *WEBKDD'2000 workshop: Web Mining for E-Commerce - Challenges and Opportunities*, 2000.
- [7] Ralph Kimball, Laura Reeves, Margy Ross, Warren Thornthwaite, *The Data Warehouse Lifecycle Toolkit: Expert Methods for Designing, Developing, and Deploying Data Warehouses*, John Wiley & Sons, 1998.
- [8] Ralph Kimball and Richard Merz, *The Data Webhouse Toolkit: Building the Web-Enabled Data Warehouse*, John Wiley & Sons, 2000.

About the authors:

Ron Kohavi is the Director of Data Mining at Blue Martini Software. Prior to joining Blue Martini, he managed the MineSet project, Silicon Graphics' award-winning product for data mining and visualization. Kohavi received his Ph.D. in Machine Learning from Stanford University, where he led the MLC++ project, the Machine Learning library in C++. He received his BA from the Technion, Israel.

Carla E. Brodley is an Associate Professor in the School of Electrical and Computer Engineering at Purdue University. She received her bachelors degree from McGill University in 1985 and her PhD in computer science from the University of Massachusetts in 1994. Her research interests include machine learning, computer vision, and content-based image retrieval. She has applied techniques from these areas to problems from a variety of fields including remote sensing, medical images and computer security.

Brian Frasca is a Principal Engineer at Blue Martini Software. He received his M.S. from Stanford specializing in databases. He joined Blue Martini in April 1999 where he has designed the decision support data warehouse and has developed data mining transformations.

Llew Mason is the Manager of Data Mining Analytics at Blue Martini Software. He joined Blue Martini Software in September 1999 after completing his Ph.D. in Systems Engineering from the Australian National University, Canberra, Australia. His research interests include machine learning, computational learning theory, large margins analysis and methods for combining classifiers.

Zijian Zheng received his Ph.D. in computer science from the University of Sydney in 1996. He joined Blue Martini Software as a Senior Software Engineer in May 1999. He is an active researcher and developer in the areas of data mining and machine learning.