

BLUE MARTINI

DIMACS/IBM Workshop on Data Mining in the Internet Age

Mining E-commerce Data: Challenges and Stories from the Trenches

Ronny Kohavi
Director, Data Mining
Blue Martini Software
ronnyk@bluemartini.com
<http://robotics.Stanford.EDU/~ronnyk/>

2 May 2000

© Copyright 1998-2000, Blue Martini Software, San Mateo California, USA

BLUE MARTINI

Overview

- o E-commerce: killer domain for data mining
- o Blue Martini Software value proposition
- o Architectural decisions
- o Stories and Challenges
- o Summary

2

© Copyright 1998-2000, Blue Martini Software, San Mateo California, USA

BLUE MARTINI

E-Commerce: the Killer Domain for Data Mining

- o E-commerce provides all the right ingredients for successful data mining, including:
 - o Large amounts of data (many records)
 - o Rich data with many attributes (wide)
 - o Clean data
 - o Actionable domain
 - o Measurable return-on-investment



© Copyright 1998-2000, Blue Martini Software, San Mateo California, USA

BLUE MARTINI

E-Commerce: the Killer Domain for Data Mining (II)

Why Data Mining and E-Commerce are a match made in heaven (or in the data warehouse):

- o Clickstreams at sites provide data that dwarfs large warehouses built in previous years
- o Designed correctly, sites can assign many attributes to content on web pages, customers, products, purchases.
- o Data is collected electronically at the webstore and it is clean (no painful legacy transformations)
- o Insight derived from data mining analysis is easily turned into action, closing the loop with the transactional systems
- o ROI is easy to measure at a webstore

4

© Copyright 1998-2000, Blue Martini Software, San Mateo California, USA

BLUE MARTINI

Data Mining becomes Important

- o Prior to 2000, horizontal data mining companies were bought for about \$10M
 - o Compression Sciences was bought by Gentia (\$3M)
 - o HyperParallel was bought by Yahoo (\$2.3M)
 - o Clementine was bought by SPSS (\$7M)
 - o Thinking Machines was bought by Oracle (~ \$25M)
- o Around 2000, a phase shift occurred in valuations, which were \$100M-\$500M
 - o NetPerceptions bought KD1 for 2.24 million shares worth \$116M (now worth about \$43 million)
 - o Epiphany bought RightPoint (previously DataMind) for 3.6 million shares worth \$400M (now worth about \$240M)
 - o Vignette bought DataSage for 3.16 million shares worth \$577 (now worth \$456M after 3 to 1 split)
 - o NeoVista bought by Accrue for 2.4 million shares worth \$140M (now worth \$70M)

5

© Copyright 1998-2000, Blue Martini Software, San Mateo California, USA

BLUE MARTINI

Blue Martini Software

Blue Martini's solution is built on the vision of a complete E-commerce solution with integrated data mining, called Micro Marketing



300+ employees
35+ customers including:

- Levi's
- Harley Davidson
- Gymboree
- The Men's Warehouse
- gloss.com
- gazelle.com
- craft.com

© Copyright 1998-2000, Blue Martini Software, San Mateo California, USA

Value Proposition

BLUE MARTINI

- o A company's **brand** is a strategic asset
 - o Avoid diluting it with a mediocre web store
 - o Leverage the internet to build your brand and reinforce the same message from multiple touchpoints (web, call centers, wireless, bricks and mortar)
- o Support personalized recommendations/cross-sells
 - o Increase conversion rates (browse to buy)
 - o Increase basket size (through cross sells, better navigation)
 - o Increase customer retention and loyalty



© Copyright 1998-2000, Blue Martini Software, San Mateo California, USA

The Webstore is an Experimental Laboratory

BLUE MARTINI

- o Opening webstores will not dramatically impact revenues for established retailers, but lessons learned will affect other channels
- o The webstore provides an experimental laboratory and a trend-discovery system
 - o Which cross-sells work?
 - o Which ads are effective?
 - o What are people looking for (failed searches for pokédex)

E-Commerce

Amazon \$20.2 B

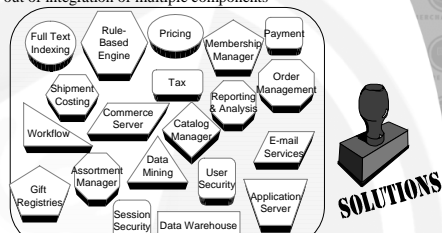
Wal-Mart
1999 revenues:
\$162.8 B

© Copyright 1998-2000, Blue Martini Software, San Mateo California, USA

Perils of Components

BLUE MARTINI

Blue Martini charges close to \$1M for the software because we take the pain out of integration of multiple components

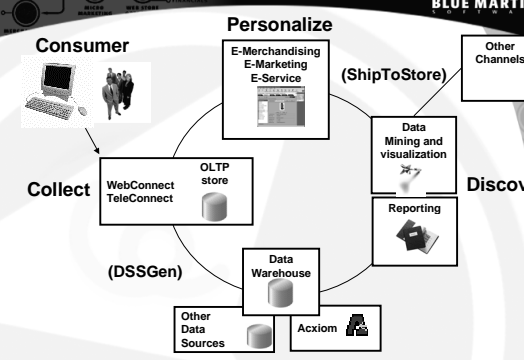


Average E-Commerce site costs \$5.9M to assemble and \$4.3M annually to maintain.
International Data Corporation, January 1999

© Copyright 1998-2000, Blue Martini Software, San Mateo California, USA

Closing the Loop

BLUE MARTINI



© Copyright 1998-2000, Blue Martini Software, San Mateo California, USA

Architectural Decisions

BLUE MARTINI

- o Web pages are dynamic (jsp/jhtml) making API calls everywhere
 - o Everything is in the database, the web pages access it.
 - o Every object can have attributes (e.g., customers, products, assortments, orders, content pages) which are accessible in the webstore and later in mining
- o Webstore's app server writes clickstreams
 - o Knows about sessions - avoid sessionizing problems
 - o Knows what was displayed - products, assortments
 - o Can save additional information not in weblogs
 - o Central DB repository - avoid need to join weblogs
 - o Can investigate orders and clickstreams together

© Copyright 1998-2000, Blue Martini Software, San Mateo California, USA

Architectural Decisions (II)

BLUE MARTINI

- o Provide automatic ETL (Extract/Transform/Load) to build Data Warehouse, cutting the 70-90% data preparation step in KDD dramatically
- o Provide access to syndicated data (Acxiom) About 60 customer attributes
- o Personalization rules engine runs at the webstore, providing cross-sells, recommended images, assortments, etc.
- o ShipToStore closes the loop, allowing model scores to be shipped back to store for personalization

© Copyright 1998-2000, Blue Martini Software, San Mateo California, USA

Story I - Mysterious Birth Years

BLUE MARTINI

- The KDD-98 data contained interesting anomalies for date of birth:
 - Spikes on years ending in zero
 - No individuals were born prior to 1910
 - There were twice as many individuals who were born on even years as on odd years

Why?

Graph adopted from KDD'99 Competition by Jim Georges and Anne Milley, SAS institute SIGKDD Explorations, Jan 2000, Volume 1, Issue 2, page 79

© Copyright 1998-2000, Blue Martini Software, San Mateo California, USA

13

Challenge I - Dates

BLUE MARTINI

- Dates/times are **very** important and appear very frequently in e-commerce, yet most data mining algorithms do not support them
- Common themes:
 - Provide well-used measurements in industry, such as Recency and Frequency (of RFM).
 - Provide strong support for date operations (days between dates, day-of-week, etc) and let the users define date operations

Is there something better we can do?

© Copyright 1998-2000, Blue Martini Software, San Mateo California, USA

Story II - Gender Mystery

BLUE MARTINI

- A site has gender on the registration form
- Acxiom, a syndicated data provider, also provides gender
- There was a very large discrepancy between the percentage of males according to registration and by Acxiom

Why?

© Copyright 1998-2000, Blue Martini Software, San Mateo California, USA

15

Challenge II - Unknown vs. N/A

BLUE MARTINI

- A *pregnant* attribute on a customer can have four values: yes, no (response), unknown (we didn't ask), and not applicable (e.g., for males or young children).
- E-commerce products have attributes that depend on their family:
 - Pants have length, inseam, fabric material
 - Microwaves have voltage, cubic inches, weight
- For data mining, we flatten all attributes and most values are N/A (e.g., voltage for pants)
- How should we handle N/A versus unknown?

© Copyright 1998-2000, Blue Martini Software, San Mateo California, USA

16

Challenge II' - Product Hierarchies

BLUE MARTINI

- Products are typically arranged in a hierarchy. Most algorithms expect same-size records
- Common themes:
 - Flatten product attributes (lots of N/As).
 - Allow users to choose parts of hierarchy for pivots based on product id (SKU). Add Boolean columns from hierarchy

© Copyright 1998-2000, Blue Martini Software, San Mateo California, USA

17

Story III - Low Conversion Rates

BLUE MARTINI

- Conversion rate is the ratio of buyers to browsers.
- High conversion rates are obviously desired
- Reports we generate can show highest and lowest conversion rates
- Why do some products have really low conversion rates?

© Copyright 1998-2000, Blue Martini Software, San Mateo California, USA

18

Challenge III - Changing Sites

- One of the hardest problems is dealing with changing sites - all our I.I.D assumptions do not hold!
- It is very easy to change personalization rules (cross-sells/up-sells), images shown, messages, etc
- Sites changes their registration forms
- New products are introduced, old ones are deleted (especially with perishable products such as wines) or change attributes
- How do we deal with slowly changing dimensions? (or quickly changing?)

© Copyright 1998-2000, Blue Martini Software, San Mateo California, USA 19

Story IV - Analyst numbers

4Q 1999 Online Holiday Shopping Revenues (Billions). Why don't the numbers agree?

Analyst	Revenue (Billions)
BizRate.com	\$4.4
Gartner Group	\$4.5
Forrester Research	\$5.0
US Commerce Dept	\$5.3
Jupiter Communications	\$6.0
IDC	\$7.1
Yankee Group	\$8.0
Dataquest	\$8.5
Boston Consulting Group	\$10.5
Goldman Sachs	\$12.0
Ernst & Young	\$13.0

© Copyright 1998-2000, Blue Martini Software, San Mateo California, USA 20

Challenge IV - Scalability

- Yahoo had 465 million page views per day in December of 1999 (*)
- That's about 2-4GB of clickstream an hour, depending on the amount of clickstream information stored
- What can we do with such volumes?
- Are there useful aggregations of such data that can be done on the fly?

(*) CFO April 2000

© Copyright 1998-2000, Blue Martini Software, San Mateo California, USA 21

More Challenges

- How can we prune the number of associations generated?
- What is the loss metric for displaying a cross-sell?
 - If you display the product most-likely to be purchased, would you display bread on all pages until bread was purchased?
 - If you display product with highest lift, would you still display it if the probability was 0.01% (up from 0.00001%)?
- Are there algorithms that could take a star-schema and mine it without flattening it (e.g., Query flocks)?
- Bots/Crawlers tend to skew statistics dramatically.
- How can marketing campaigns be taken into account?

© Copyright 1998-2000, Blue Martini Software, San Mateo California, USA 22

Business Challenges

- Most companies do not have the expertise to build data mining transformations. Can we automate them? The more vertical, the easier this gets.
- How can we generate comprehensible models? Actionable models?

© Copyright 1998-2000, Blue Martini Software, San Mateo California, USA 23

Summary

- E-commerce is the killer-domain for data mining
- Automatic generation of data warehouse and closing the loop back to store is key to making data mining usable and actionable
- Clickstreams need to be collected at the app-server level where meta information exists
- Challenges: Date handling, unknown vs. N/A, hierarchy support, constantly changing sites, scalability, and more

Some images used herein where obtained from DMSI's MasterClips/Master Photo(C) Collection, 1895 Francisco Blvd East, San Rafael 94901-5506, USA

© Copyright 1998-2000, Blue Martini Software, San Mateo California, USA 24