

KMAP – A Visualizer for Kohonen Self-Organizing Map Data Mining Results

George S. Almasi

Staff Member Emeritus, IBM T. J. Watson Research Center
Vice President, Michael Rothman & Associates
almasi@gsalmasi.com, (914) 232-2378

Richard D. Lawrence

Staff Member, IBM T. J. Watson Research Center

Michael J. Rothman

President, Michael Rothman & Associates

DIMACS Workshop on Visualization and Data Mining, October 24, 2002

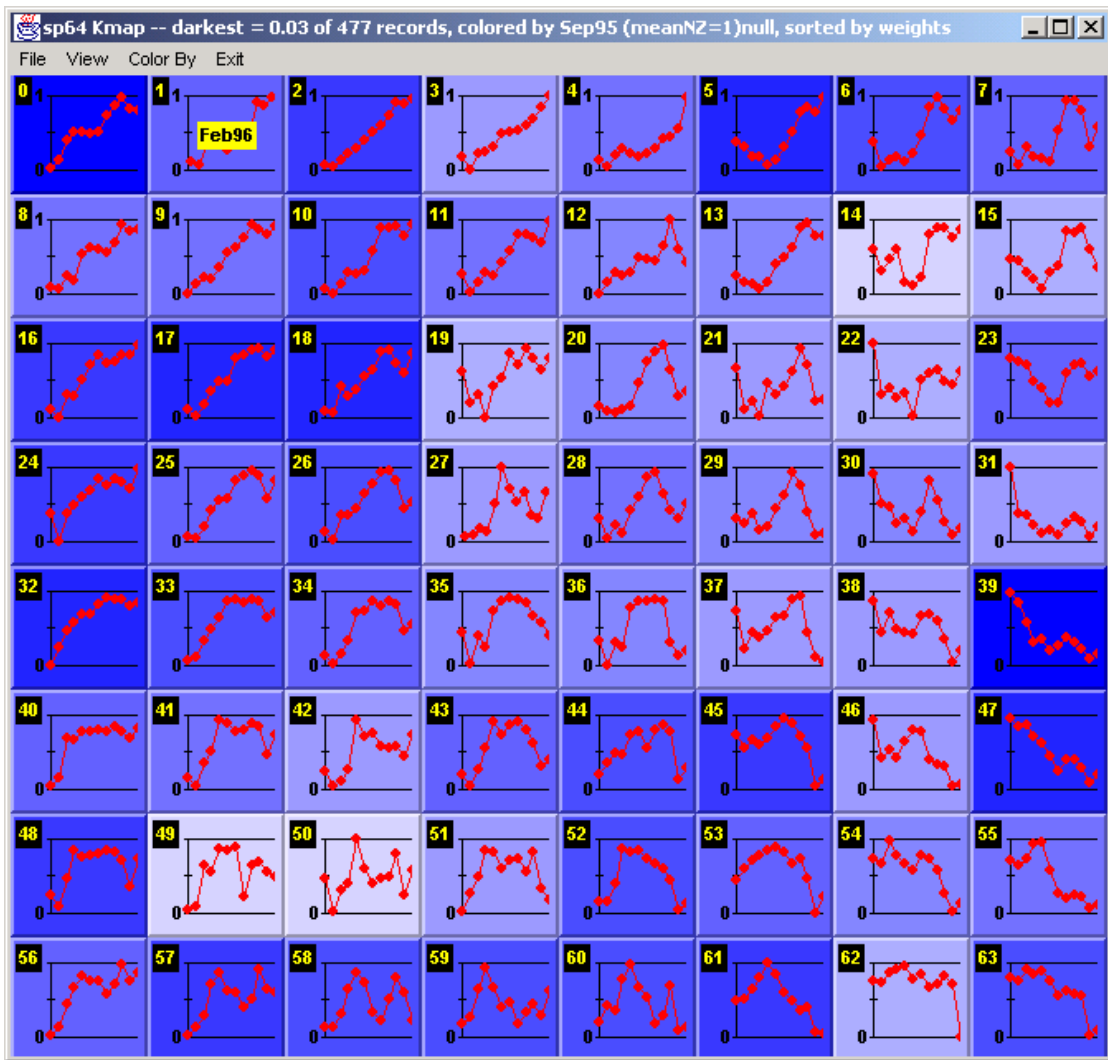
Kohonen Self-Organizing Map is a data mining technique for breaking a data set into clusters

- Similar records are clumped together
- Similar clumps are placed near each other in a 2-D map
- More precisely, similar records from a dataset are separated into clusters using a Euclidian metric, and then the multi-dimensional clusters are put in a 2-D map in such a way as to preserve the underlying topology

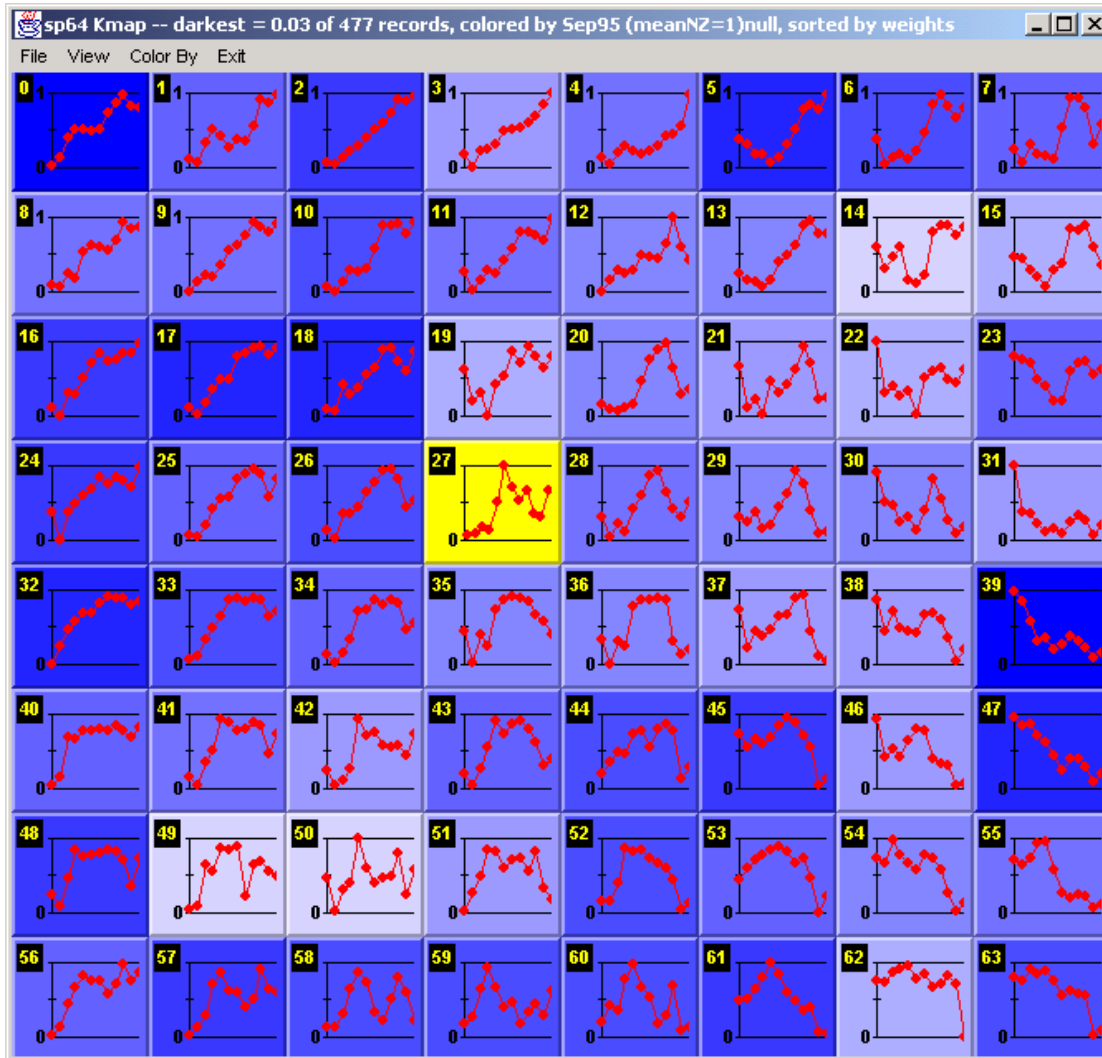
The KMAP visualizer shows Kohonen Map topology and allows interactive exploration of the clusters

- Written in Java for portability
- Uses Sitraka JCChart Java Beans
- Can run as Web applet
- Several demos and view modes will be shown:
 1. Stock Market (S&P 500 “birds of a feather”)
 - The straightforward view, but with some surprises
 2. Supermarket Shoppers (for personalized recommendations)
 - The “top ten” view for data with many attributes
 3. Credit Bureau (bankruptcy prediction)
 1. Custom “patterns within patterns” view for “extracted scores”, a method for dealing with still more attributes and disparate data bases.

A years' comparison of S&P 500 month-end closing prices, broken into 64 clusters...



S&P 500 drilldown: where's IBM?



Record Finder

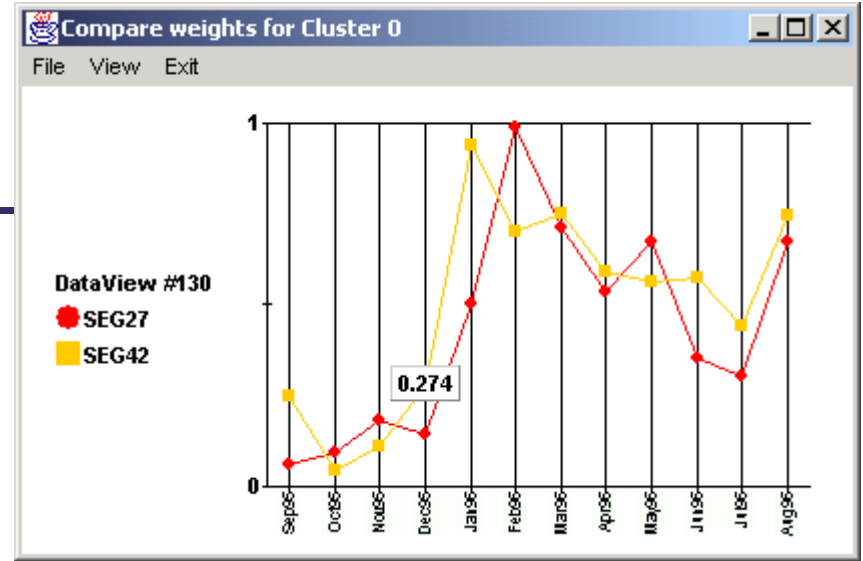
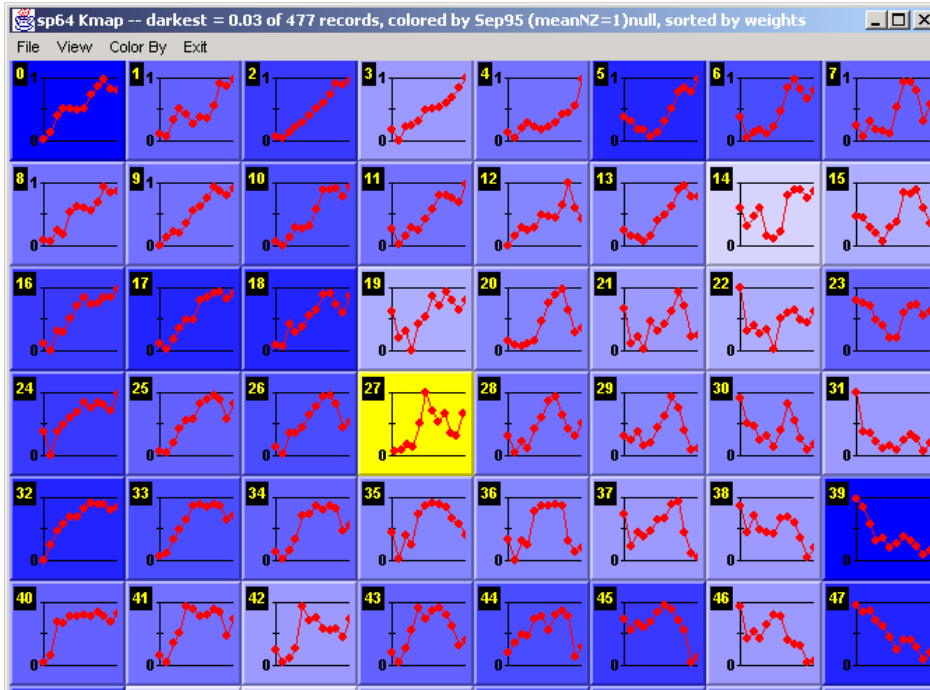
OK

Select the record you wish to locate:

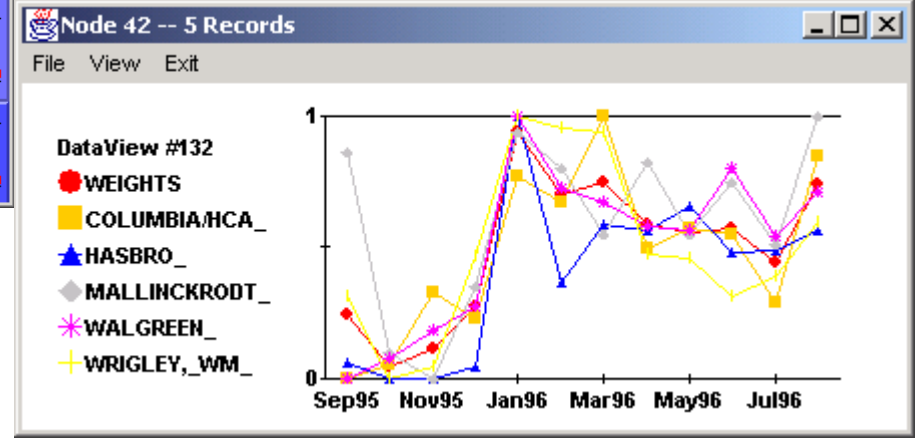
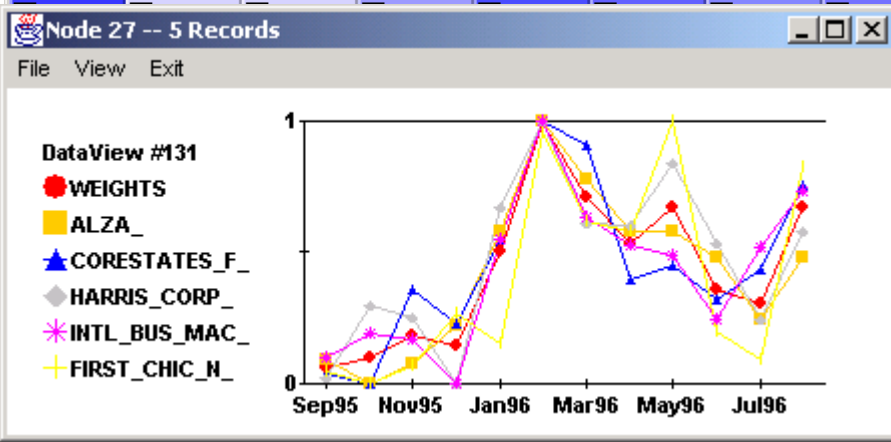
Find

- ILLINOIS_TOO_
- INCO_LTD_
- INGERSOLL-RA_
- INLAND_STEEL_
- INTEL_
- INTERGRAPH_
- INTL_BUS_MAC**
- INTL_FLAV_+_

Can one stock predict another's behavior?



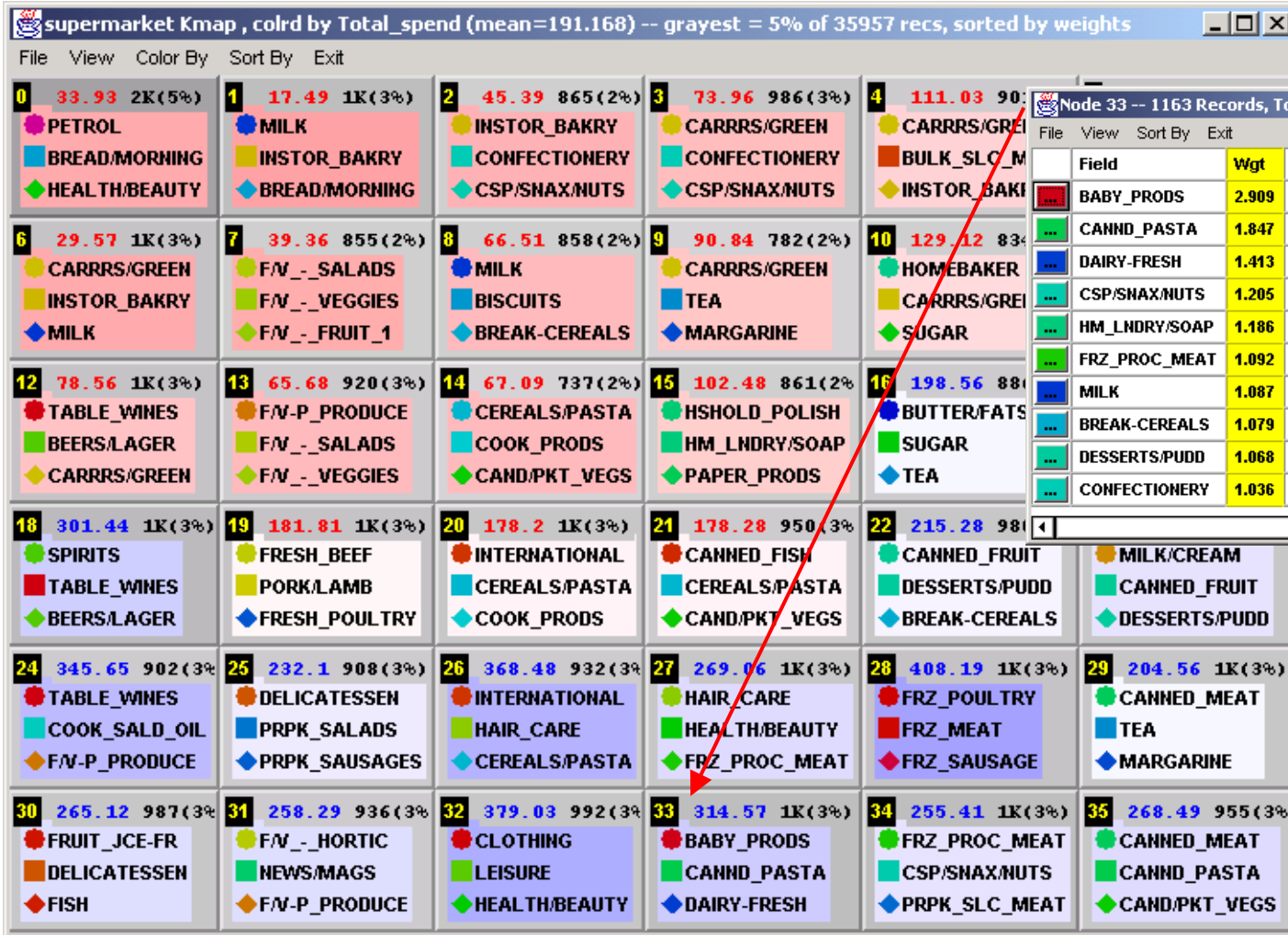
Does Hasbro predict IBM?



The KMAP visualizer shows Kohonen Map topology and allows interactive exploration of the clusters

- Written in Java for portability
- Uses Sitraka JCChart Java Beans
- Can run as Web applet
- Several demos and view modes will be shown:
 1. Stock Market (S&P 500 “birds of a feather”)
 - The straightforward view, but with some surprises
 - ➔ 2. Supermarket Shoppers (for personalized recommendations)
 - The “top ten” view for data with many attributes
 3. Credit Bureau (bankruptcy prediction)
 1. Custom “patterns within patterns” view for “extracted scores”, a method for dealing with still more attributes and disparate data bases.

Supermarket shoppers fall into distinct clusters



Drilldown:

Node 33 -- 1163 Records, Total_spend=315, Age=32.055

File View Sort By Exit

Field	Wgt	mean	R	meanNZ	R	partcp	R
BABY_PRODS	2.909	34.16	12.90	34.22	2.40	1.00	5.3
CANND_PASTA	1.847	2.25	4.94	2.80	1.68	0.80	2.9
DAIRY-FRESH	1.413	10.76	2.55	11.59	1.73	0.93	1.4
CSP/SNAX/NUTS	1.205	7.93	2.05	8.70	1.39	0.91	1.4
HM_LNDRY/SOAP	1.186	10.99	1.97	12.18	1.34	0.90	1.4
FRZ_PROC_MEAT	1.092	8.64	1.93	10.96	1.22	0.79	1.5
MILK	1.087	9.16	1.82	9.75	1.55	0.94	1.1
BREAK-CEREALS	1.079	7.11	1.83	8.15	1.29	0.87	1.4
DESSERTS/PUDD	1.068	1.94	2.24	3.10	1.29	0.63	1.7
CONFECTIONERY	1.036	7.03	1.85	8.13	1.28	0.86	1.4

P. S.: This cluster's favorite candy bar is quite different from the overall favorite!

“Wallet share” offers a second view...


BABY_PRODS Kmap; darkest shade = 42% market share of 95,238

0	1	2	3	4	5
02% Val	01% Val	01% Val	01% Val	01% Val	02% Val
m 1 0.33	m 0 0.17	m 1 0.38	m 1 0.32	m 1 0.23	m 2 0.67
mNZ 12 0.85	mNZ 9 0.63	mNZ 10 0.67	mNZ 7 0.46	mNZ 6 0.45	mNZ 10 0.67
ptp 0.07 0.3	ptp 0.05 0.2	ptp 0.11 0.5	ptp 0.13 0.7	ptp 0.1 0.51	ptp 0.18 0.9
6	7	8	9	10	11
01% Val	00% Val	01% Val	00% Val	01% Val	02% Val
m 1 0.34	m 0 0.13	m 1 0.43	m 0 0.17	m 1 0.37	m 1 0.52
mNZ 12 0.83	mNZ 10 0.72	mNZ 10 0.73	mNZ 6 0.45	mNZ 8 0.57	mNZ 8 0.56
ptp 0.08 0.4	ptp 0.03 0.1	ptp 0.11 0.5	ptp 0.07 0.3	ptp 0.12 0.6	ptp 0.17 0.9
12	13	14	15	16	17
01% Val	01% Val	00% Val	02% Val	01% Val	01% Val
m 1 0.25	m 1 0.29	m 1 0.23	m 2 0.73	m 1 0.49	m 1 0.31
mNZ 9 0.63	mNZ 11 0.78	mNZ 7 0.52	mNZ 12 0.84	mNZ 7 0.5	mNZ 6 0.41
ptp 0.07 0.4	ptp 0.07 0.3	ptp 0.08 0.4	ptp 0.16 0.8	ptp 0.18 0.5	ptp 0.14 0.7
18	19	20	21	22	23
01% Val	01% Val	02% Val	01% Val	01% Val	01% Val
m 1 0.43	m 1 0.41	m 2 0.58	m 1 0.49	m 1 0.31	m 1 0.31
mNZ 8 0.56	mNZ 8 0.59	mNZ 11 0.76	mNZ 9 0.62	mNZ 6 0.39	mNZ 6 0.4
ptp 0.14 0.7	ptp 0.13 0.6	ptp 0.14 0.7	ptp 0.15 0.7	ptp 0.15 0.7	ptp 0.15 0.8
24	25	26	27	28	29
03% Val	02% Val	04% Val	01% Val	04% Val	01% Val
m 4 1.35	m 2 0.68	m 4 1.63	m 1 0.48	m 4 1.49	m 1 0.21
mNZ 16 1.11	mNZ 11 0.78	mNZ 13 0.93	mNZ 5 0.34	mNZ 12 0.85	mNZ 5 0.35
ptp 0.23 1.2	ptp 0.16 0.8	ptp 0.33 1.7	ptp 0.27 1.4	ptp 0.33 1.7	ptp 0.11 0.6
30	31	32	33	34	35
03% Val	02% Val	09% Val	42% Val	01% Val	03% Val
m 3 1.14	m 2 0.68	m 8 3.11	m 34 12.9	m 1 0.26	m 3 1.28
mNZ 13 0.91	mNZ 9 0.61	mNZ 19 1.32	mNZ 34 2.4	mNZ 4 0.28	mNZ 11 0.77
ptp 0.23 1.2	ptp 0.21 1.1	ptp 0.44 2.3	ptp 1 5.36	ptp 0.17 0.5	ptp 0.31 1.6

SPIRITS Kmap; darkest shade = 47% market share of 113,756

0	1	2	3	4	5
02% Val	01% Val	00% Val	00% Val	01% Val	02% Val
m 2 0.53	m 1 0.21	m 0 0.11	m 1 0.17	m 1 0.26	m 2 0.55
mNZ 33 1.25	mNZ 19 0.72	mNZ 14 0.53	mNZ 13 0.49	mNZ 19 0.71	mNZ 15 0.58
ptp 0.05 0.4	ptp 0.03 0.2	ptp 0.02 0.2	ptp 0.04 0.3	ptp 0.04 0.3	ptp 0.11 0.9
6	7	8	9	10	11
01% Val	00% Val	01% Val	01% Val	01% Val	04% Val
m 1 0.44	m 0 0.11	m 1 0.25	m 1 0.23	m 1 0.25	m 4 1.11
mNZ 30 1.14	mNZ 16 0.62	mNZ 20 0.76	mNZ 13 0.51	mNZ 14 0.52	mNZ 23 0.89
ptp 0.05 0.3	ptp 0.02 0.1	ptp 0.04 0.3	ptp 0.05 0.4	ptp 0.06 0.4	ptp 0.15 1.2
12	13	14	15	16	17
03% Val	00% Val	00% Val	01% Val	02% Val	02% Val
m 4 1.18	m 1 0.16	m 0 0.09	m 1 0.28	m 3 0.83	m 3 0.92
mNZ 27 1.02	mNZ 19 0.72	mNZ 10 0.39	mNZ 18 0.68	mNZ 21 0.79	mNZ 19 0.7
ptp 0.14 1.1	ptp 0.03 0.2	ptp 0.03 0.2	ptp 0.05 0.4	ptp 0.13 1.1	ptp 0.16 1.1
18	19	20	21	22	23
47% Val	01% Val	01% Val	01% Val	01% Val	02% Val
m 44 14	m 1 0.37	m 1 0.41	m 1 0.28	m 1 0.27	m 2 0.7
mNZ 44 1.68	mNZ 20 0.76	mNZ 13 0.47	mNZ 14 0.52	mNZ 15 0.58	mNZ 21 0.8
ptp 1 8.34	ptp 0.06 0.4	ptp 0.1 0.8	ptp 0.06 0.5	ptp 0.06 0.4	ptp 0.1 0.8
24	25	26	27	28	29
05% Val	01% Val	02% Val	01% Val	02% Val	01% Val
m 7 2.09	m 2 0.53	m 3 0.85	m 2 0.49	m 2 0.76	m 1 0.4
mNZ 25 0.96	mNZ 15 0.58	mNZ 15 0.58	mNZ 16 0.59	mNZ 20 0.75	mNZ 18 0.68
ptp 0.26 2.1	ptp 0.11 0.9	ptp 0.18 1.4	ptp 0.1 0.8	ptp 0.12 1.1	ptp 0.07 0.9
30	31	32	33	34	35
01% Val	02% Val	03% Val	01% Val	02% Val	02% Val
m 2 0.52	m 2 0.69	m 3 1.04	m 1 0.38	m 2 0.48	m 2 0.7
mNZ 22 0.81	mNZ 17 0.66	mNZ 20 0.76	mNZ 13 0.5	mNZ 16 0.62	mNZ 19 0.71
ptp 0.08 0.6	ptp 0.13 1.1	ptp 0.16 1.3	ptp 0.09 0.7	ptp 0.09 0.7	ptp 0.12 1

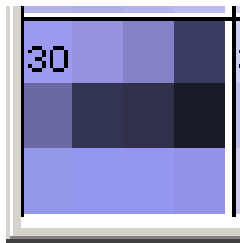
The KMAP visualizer shows Kohonen Map topology and allows interactive exploration of the clusters

- Written in Java for portability
- Uses Sitraka JCChart Java Beans
- Can run as Web applet
- Several demos and view modes will be shown:
 1. Stock Market (S&P 500 “birds of a feather”)
 - The straightforward view, but with some surprises
 2. Supermarket Shoppers (for personalized recommendations)
 - The “top ten” view for data with many attributes
 -  3. Credit Bureau (bankruptcy prediction)
 1. Custom “patterns within patterns” view for “extracted scores”, a method for dealing with still more attributes and disparate data bases.

Example: Bankruptcy Prediction

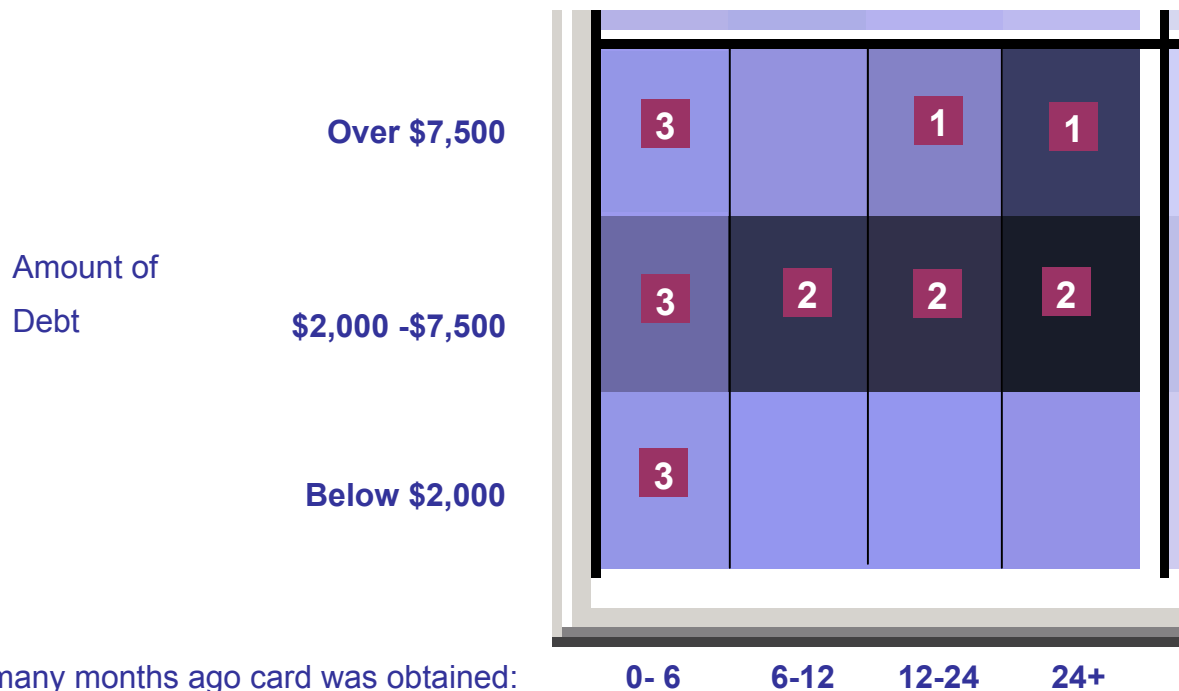
- **We applied Score Extraction to bankruptcy prediction:**

1. Complex Credit Bureau data (200+ items per record) was boiled down by dividing each household's total debt into 12 categories, based on the amount of each debt and how long ago it was incurred.
2. A data mining program divided these households into 36 groups, based on these new categories.
3. Checking subsequent bankruptcy history of each group revealed one group with unusually high bankruptcy rate
4. Placing each group's debt attributes in a small 3x4 checkerboard pattern shows the characteristic "signature" below for the high bankruptcy group:



This is explained in the next chart...

The Analyst can see a Bankruptcy Signature



What an analyst sees here:

1. they used to be able to get credit lines over \$7500, but not in the last year
2. a “sweet spot” in credit lines between \$2000 & \$7500... easy to get and they quickly add up to a lot of money
3. have been unable to get new cards in the last 6 months as their “risk” score elevates

Score (debt in that category): ■ None ■ Max

efax650_36 Kmap -- grayest = 29% of 95157 recs, colored by BANKRUPTCIES (mean=0.003), sorted by weights

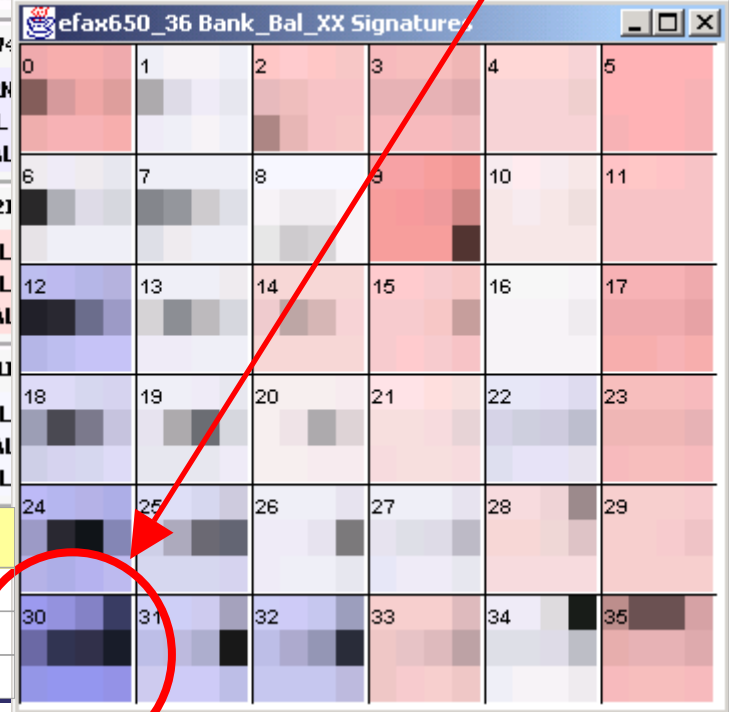
File View Color By Sort By BB Signatures Exit

0 0 1K(2%) BANK_BAL02 MORT_BAL	1 0.005 3K(3%) BANK_BAL02 BANK_BAL05	2 0.001 3K(3%) BANK_BAL01 BANK_BAL04	3 0.001 921(0%) STUDENT_BAL INSTFIN_BAL	4 0.002 2K(2%) INSTFIN_BAL STUDENT_BAL	5 0 28K(29%) MORT_BAL OTHDEPT_BAL
--------------------------------------	--	--	---	--	---

We've developed Special Tools...
...for spotting patterns fast in the extracted scores

*These 530 households
(0.6% of the total)
Have high
Probability
Of bankruptcy*

18 0.013 994(1%) BANK_BAL05 BANK_BAL08 BANK_BAL02	19 0.006 2K(2%) BANK_BAL08 BANK_BAL05 BANK_BAL11	20 0.003 2K(3%) BANK_BAL08 BANK_BAL11 MORT_BAL	21 0.002 3K(4%) POPDEPT_BAL BANK_BAL11 OTHDEPT_BAL	22 0.009 7K(7%) HOMEFURN_BAL ELEC_BAL HOME_BAL
24 0.026 734(0%) BANK_BAL08 BANK_BAL05 BANK_BAL11	25 0.013 1K(1%) BANK_BAL11 BANK_BAL08 BANK_BAL05	26 0.006 3K(3%) BANK_BAL11 MORT_BAL POPDEPT_BAL	27 0.006 1K(1%) POPDEPT_BAL BANK_BAL11 OTHDEPT_BAL	28 0.002 21K(21%) BANK_BAL BANK_BAL MORT_BAL
30 0.038 530(0%) BANK_BAL11 BANK_BAL08 BANK_BAL05	31 0.019 957(1%) BANK_BAL11 BANK_BAL12 BANK_BAL08	32 0.02 608(0%) BANK_BAL11 POPDEPT_BAL OTHDEPT_BAL	33 0.002 1K(1%) OTHDEPT_BAL BANK_BAL11 POPDEPT_BAL	34 0.005 11K(11%) BANK_BAL MORT_BAL BANK_BAL



	under 6 mo.	6 to 12 mo.	12 to 24	over 24
over \$7500	Bal_03	Bal_06	Bal_09	Bal_12
\$2000 to \$7500	Bal_02	Bal_05	Bal_08	Bal_11
under \$2000	Bal_01	Bal_04	Bal_07	Bal_10

Summary: KMAP is a versatile tool for visualizing data mining results

For more information:

1. Demos & documents at our web site: www.preference-engine.com/
2. Credit bureau data study: www.twocrows.com/largedb.pdf
3. Personalized recommendations: R. D. Lawrence, G. S. Almasi et al, “Personalization of Supermarket Product Recommendations”, *Data Mining and Knowledge Discovery* 5(1/2): 11-32 (2001) ...
4. ... and: Almasi et al, U. S. Patent 6,260,036 (2001)
5. High-speed Kohonen clustering program that made Kmap necessary: R. D. Lawrence, G. S. Almasi, H. E. Rushmeier, “A Scalable Parallel Algorithm for Self-Organizing Maps with Applications to Sparse Data Mining Problems”, *Data Mining and Knowledge Discovery* 3(2):171-195 (1999).
6. George S. Almasi, almasi@gsalmasi.com, 914-232-2378