

The High ROI of Data Mining for Innovative Organizations

UVA CMIT
Chantilly, VA
July 9, 2010

John F. Elder IV, Ph.D.
elder@datamininglab.com

Elder Research, Inc.
300 West Main Street, Suite 301
Charlottesville, Virginia 22903
434-973-7673
www.datamininglab.com

Case Study Lessons Outline

- The Hype Cycle vs. Real Results
- 3 Major ways Data Mining helps:
 - Eliminate the bad:
 - Examples from IRS, Hewlett-Packard, Capital One
 - Discover the good:
 - Ex: Pfizer, Westwind Foundation
 - Streamline / Automate:
 - Ex: Lumidigm, Peregrine, SSA, Anheuser-Busch
- Lessons Learned:
features of successful projects

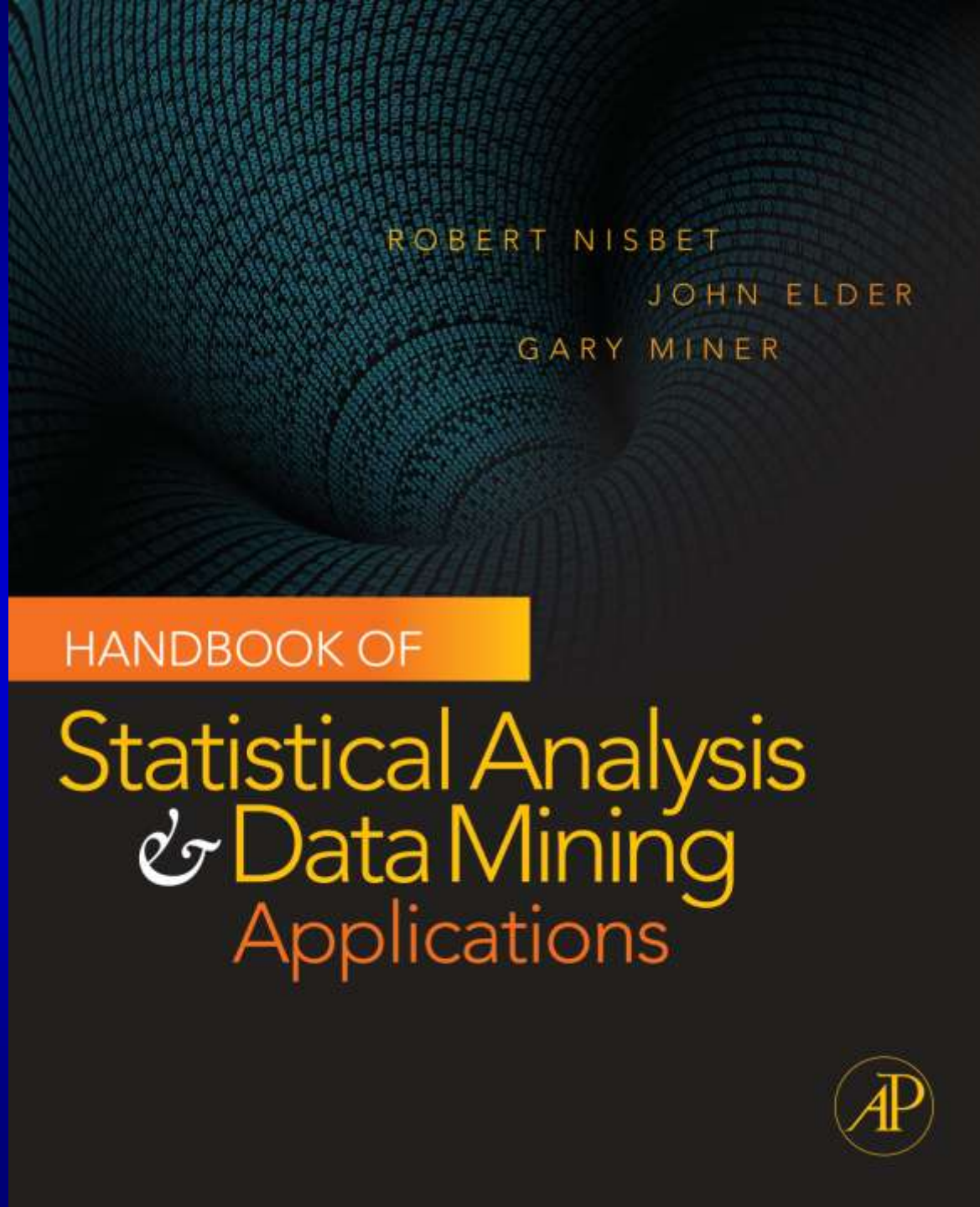
Book written *for* practitioners,
by practitioners

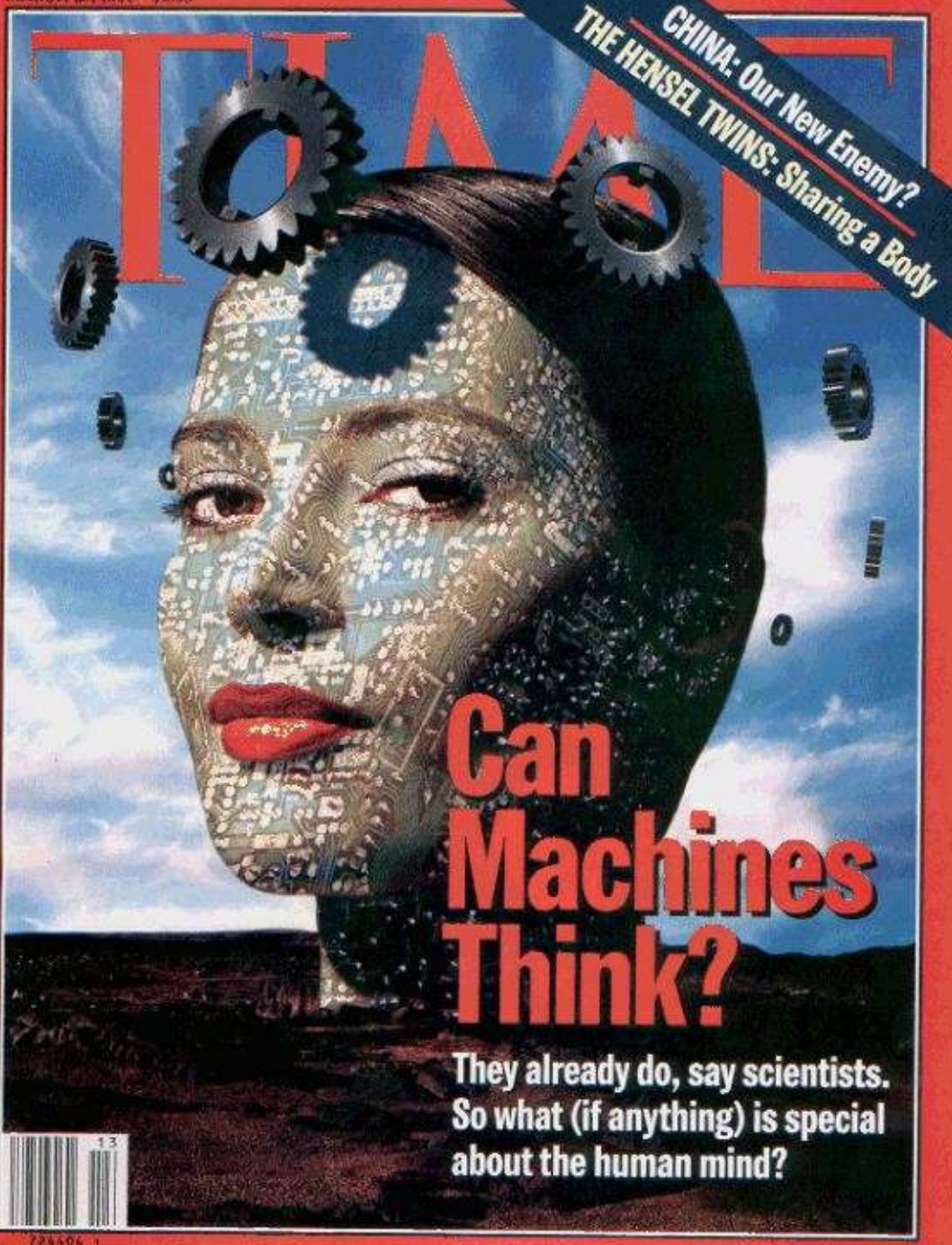
Published by Elsevier's
Academic Press in May 2009

Details and reviews at Amazon:
www.tinyurl.com/bookERI

Won 2009 PROSE Award for
Mathematics

QuickTime™ and a
decompressor
are needed to see this picture.





“Of course machines can think. After all, humans are just machines made of meat.”

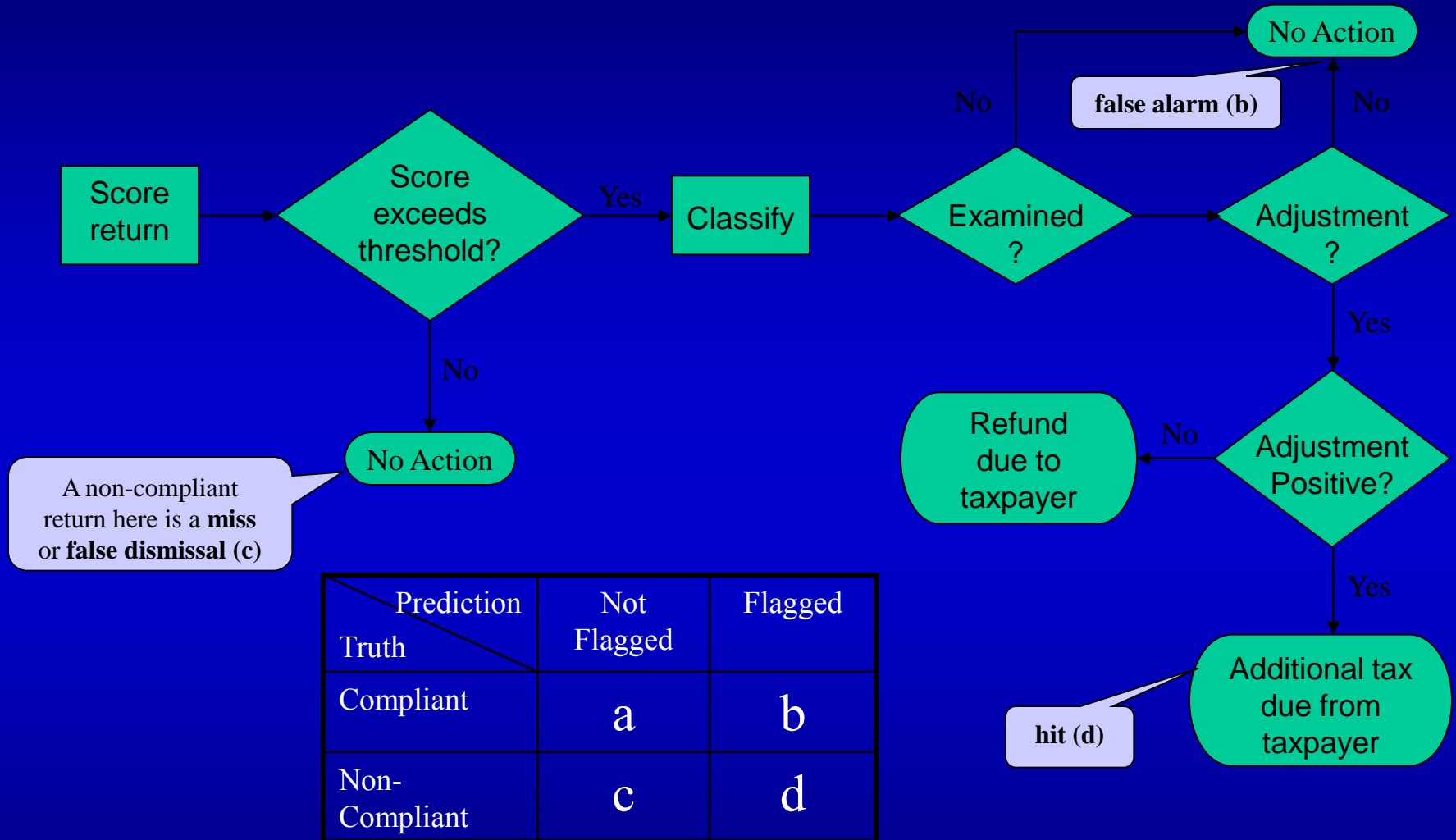
- MIT CS professor

Human and computer strengths are more complementary than alike.

3 Major Ways Data Mining Helps

- Eliminate the Bad
 - Score credit risk, Detect fraud, Find outliers
- Highlight the Good
 - Time market trades, Discover new drugs, Uncover hidden value
- Streamline or Semi-Automate Decisions
 - Intuit products of interest, Assess status, Verify identity, Be aware of “what if” scenarios, Speed service

Case 1) IRS Fraud Detection

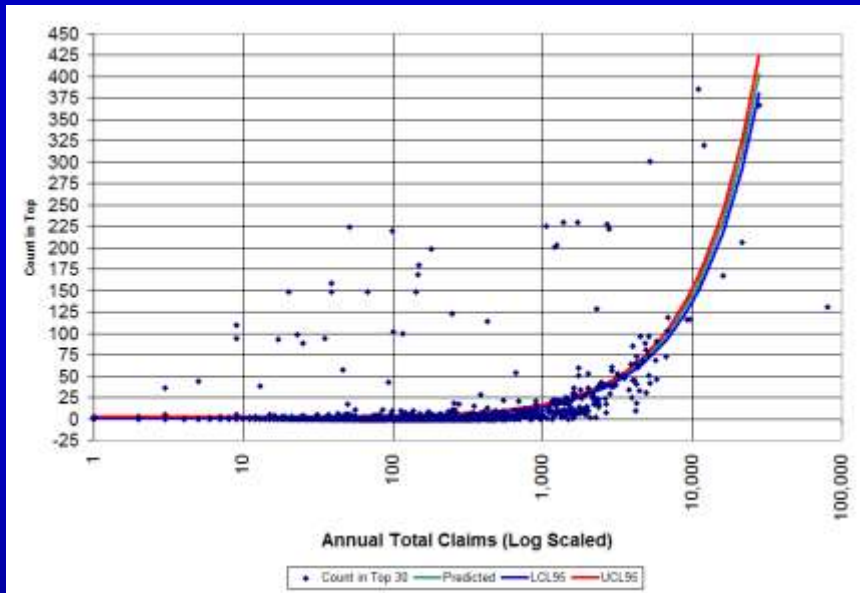


$$\text{detection rate} = d/(c+d)$$

$$\text{hit:scan} = 1 : (b+d)/d$$

$$\text{workload} = b+d$$

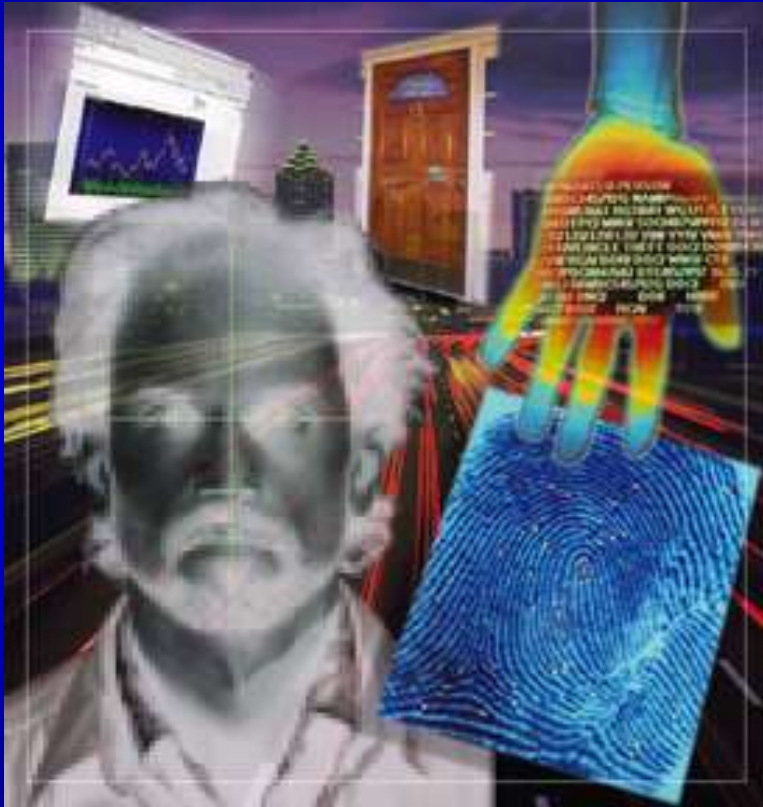
Case 2) Hewlett-Packard: Service Fraud Detection



- Tips indicated fraud exists
- Goal: Learn from known cases to find unknown
- Automate current process, build model on known, score all data, investigate top
- Recovered **\$20M** in 9 months
- Awards + promotions + growth
Became profit center

Case 3) **Lumidigm**: Bio-Metrics

(id not by what you have or know, but what you are)



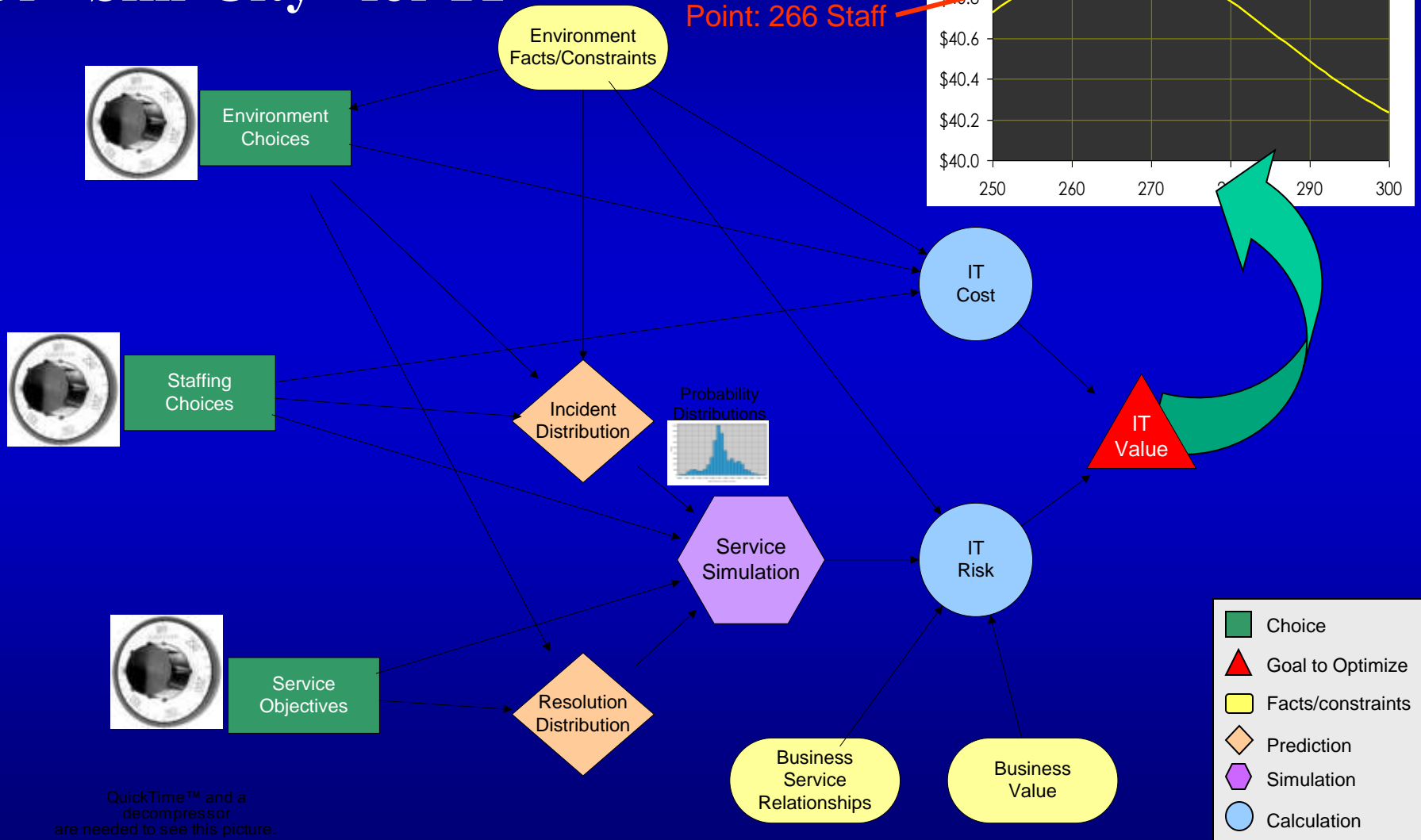
- Original Goal:
Diagnose disease
- Data: infra-red skin reflections
at multiple frequencies
- Problem: hard to remove
person specifics
- Idea: Turn problem into gain:
new bio-metric
- Q: *Who* are you?
- Use to validate ID for safety or
convenience

Solution in search of the right Application



QuickTime™ and a
decompressor
are needed to see this picture.

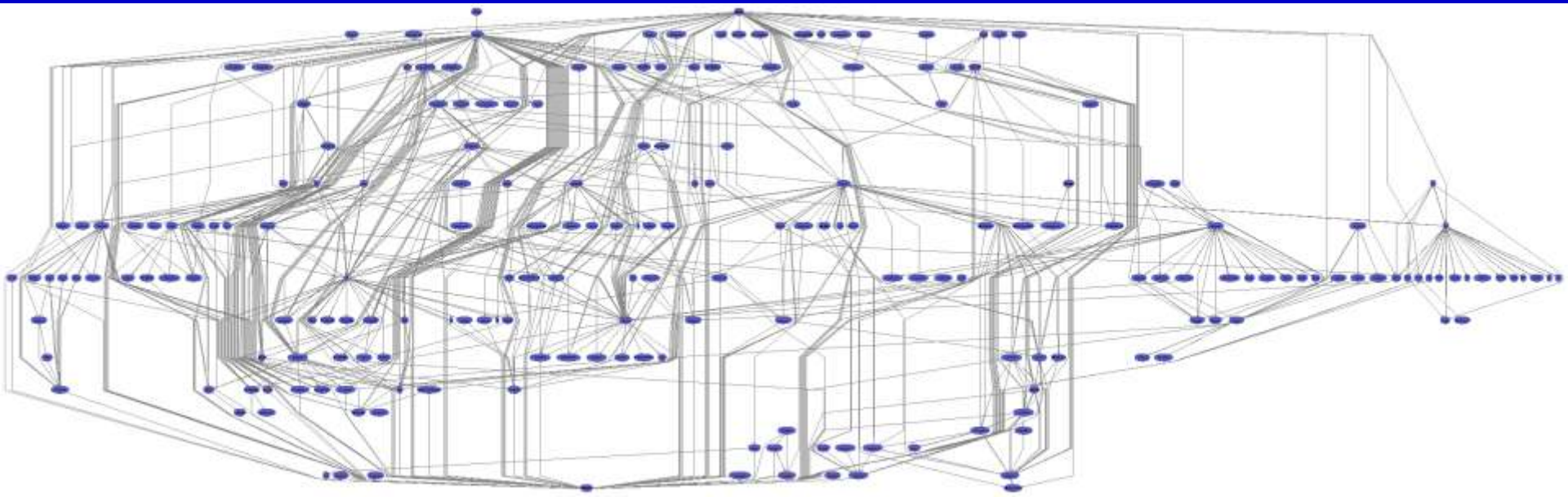
Case 4) Peregrine Systems: Business Service Modeling Or “Sim-City” for IT



QuickTime™ and a
decompressor
are needed to see this picture.

Goal: Simulation & Optimization

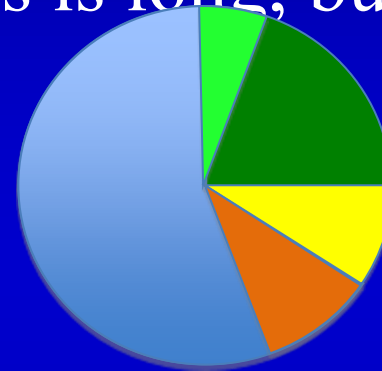
- Build tool to manage IT Service (help tickets)
- New product critical to company (& worked!)
- Analytic breakthrough: keep uncertainty to end
- Success: purchased by HP; increased HP sales



Case 5) Social Security Administration Disability Approval

- **Pain:** Approval process is long, bureaucratic

Up to
2 Years !



With text mining,
eventually approved
immediately!

1/2 of appeals overturn
original decision

- **Goal:** Fast-track “easy” cases
- **Challenge:** Free-text on disability application
- **Result:** 20% of Approvals possible immediately and with greater consistency

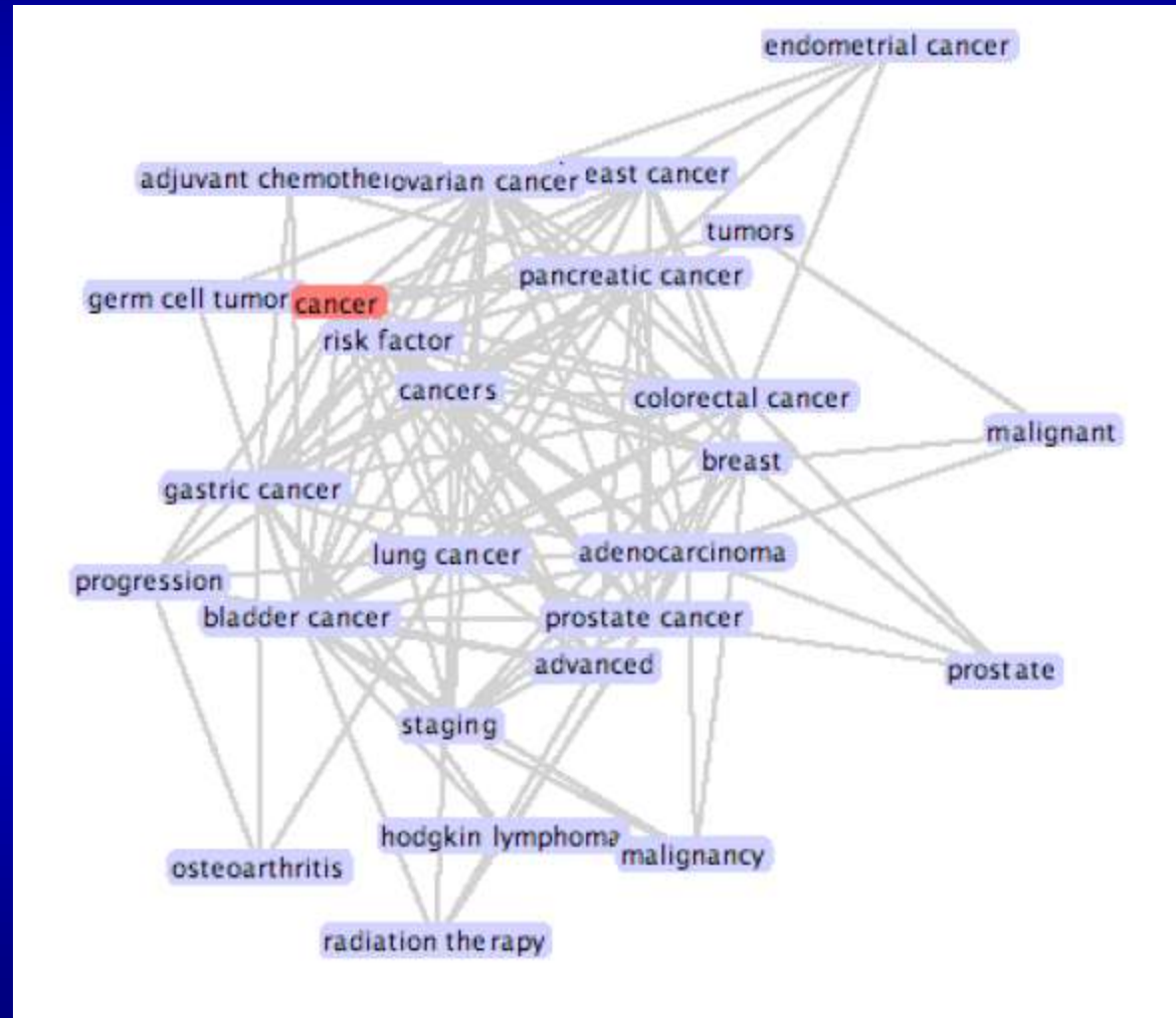
Text is Messy and Complex

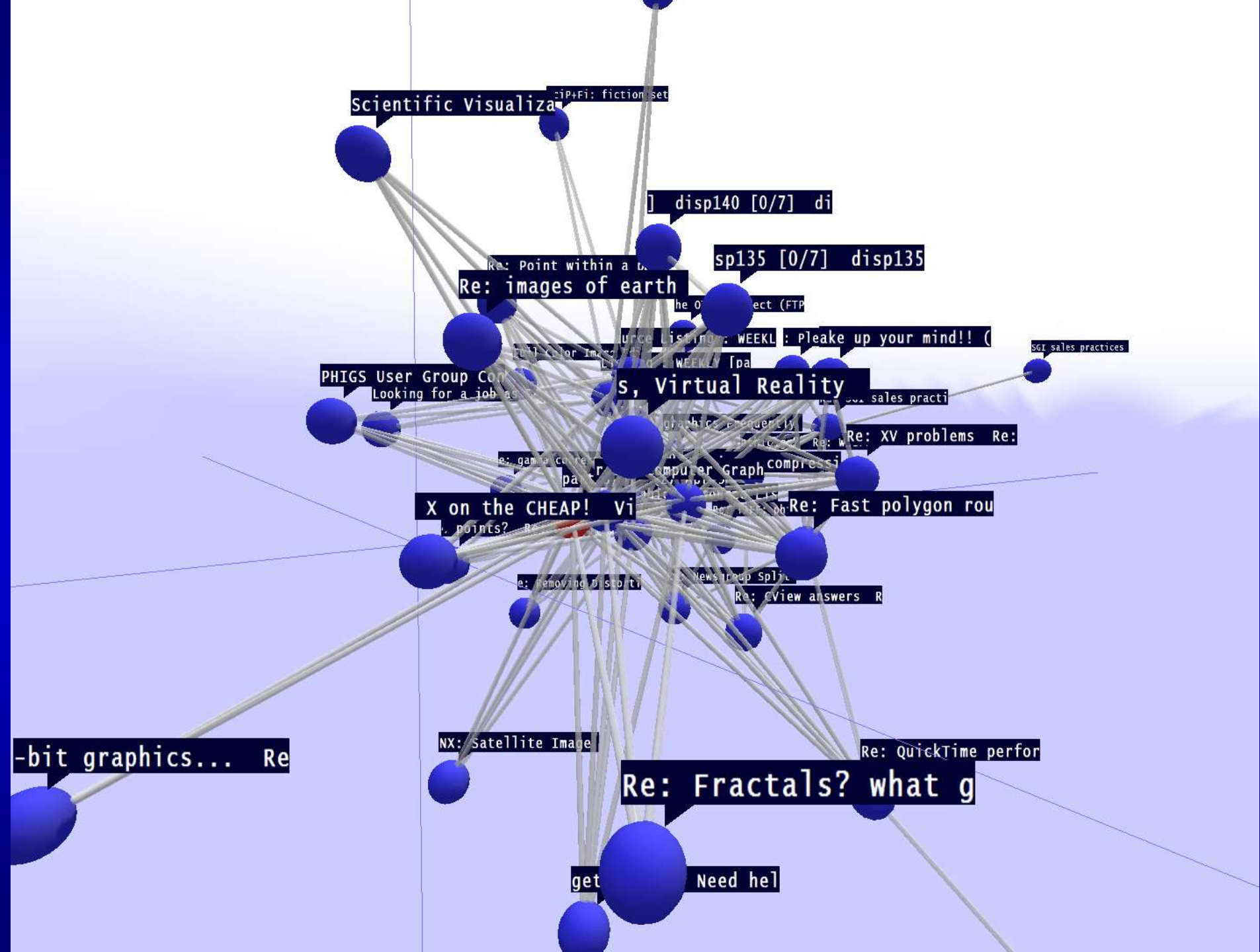
- Multi-Word Phrases (Concepts or Lemmas)
- Stemming (ex: Learn = Learning, Learned, Learns...)
- Synonyms (ex: ALS = Lou Gherig's Disease)
- Misspellings

Ex: 51 phrases found for “Learning Disability”

learning disablituy, learning deisability, learning disablity, learining disabilities, learning disabiality, learningdisabilty, learning disabiilty, learning disabilitty, learning disabilty, learning disablety, learning disabilitiy, learnoing disability, learning disabilities, learning disabilitiy, learning dsblty, learning disibility, learnings disabilty, learningdisability, larning disabilities, learning disabilities, learning disabilitties, learning disibilities, learning diasability, learning dasability, learnning disability, learning disabilities, lerning disability, learning disabiltes, learneing disability, learninig disability, learning disaiblities, leraning disability, learning disaiblity, learnings disability, learning disabilitys, learning disabillity, learnings disabilities, learning diasability, learning disabiliites, learning dsiability, learning disabliity, learning disibilty, learning disibilities, learning disbality, learning disbility, learning disabilit, learningdisabilities, learningi disability, lerniung disabilities, learning disabliities, learning disaability, learning disabilities

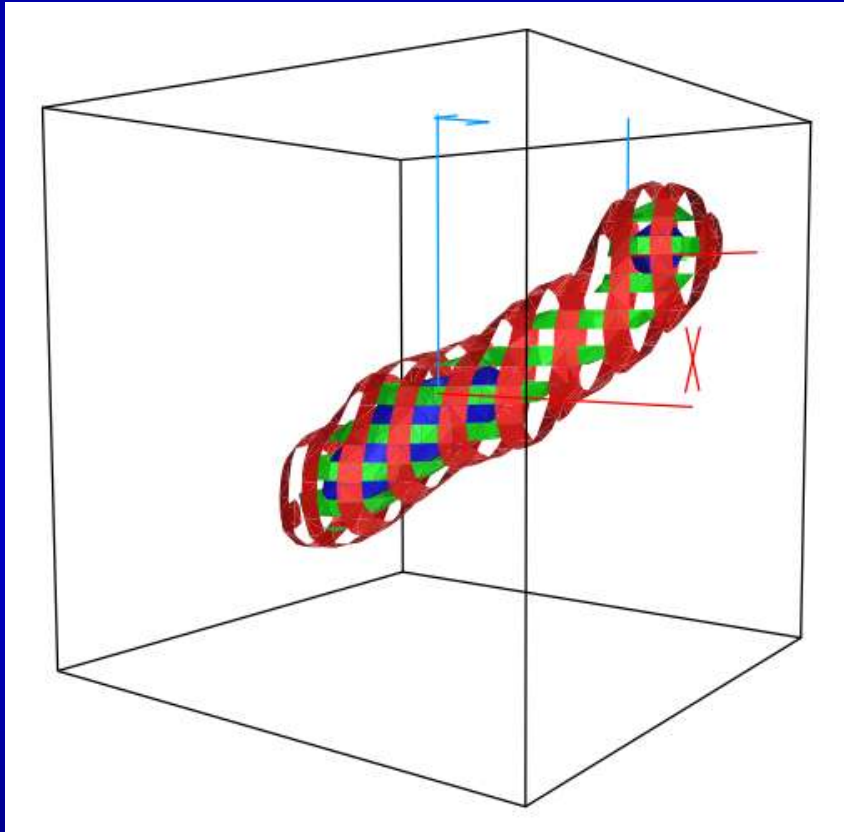
ERI new technology: Discover web of word relations



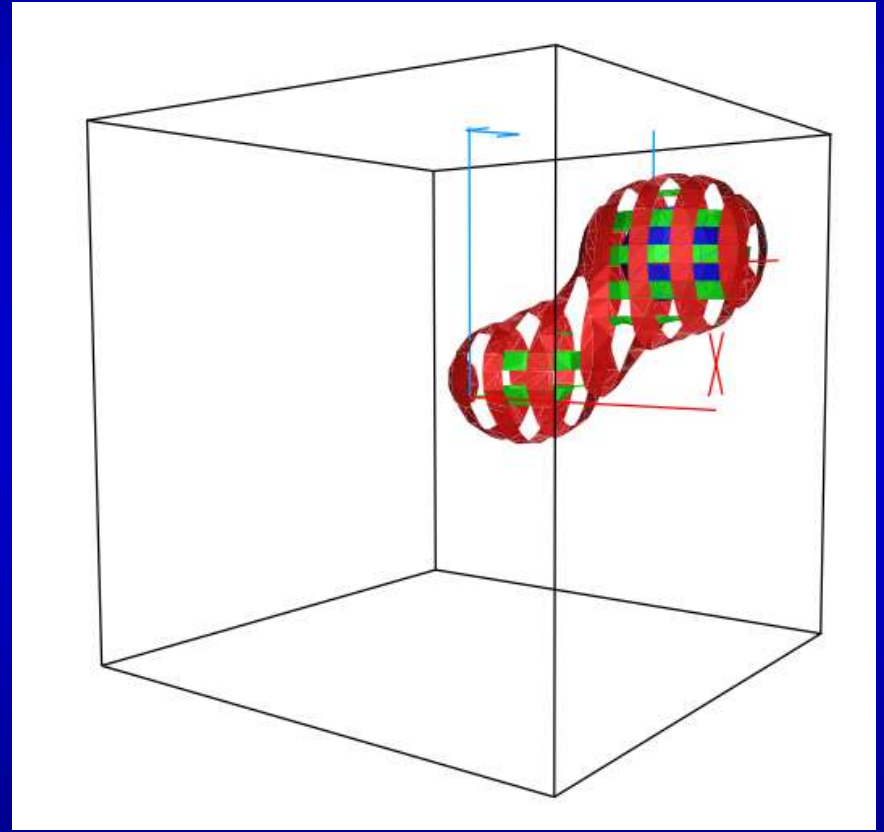


Case 6) Pharmacia & Upjohn (Pfizer) New Drug Effective?

Placebo



Drug



Density surfaces enclose ascending quartiles of data

Bottom Line

- IRS fraud detection
- HP consumer fraud detection
- Lumidigm biometric identification
- Peregrine optimization
- SSA text mining
- Pharmacia & Upjohn (Pfizer) drug efficacy

Lessons Learned: Necessary Ingredients for Analytic Project Success

- Gain Expected: either:
 - Leverageable - an incremental improvement will matter, OR
 - “Low-hanging fruit” - nobody’s yet (dared) attack the problem
- Interdisciplinary Team: experts needed in
business area, statistics, algorithms, and databases
- Data Vigilance: capture and maintain the
accumulating information stream
- Time: learning occurs over multiple cycles
- Business Champion essential!

... then Data Mining can add extraordinary value



John F. Elder IV

Chief Scientist, Elder Research, Inc.

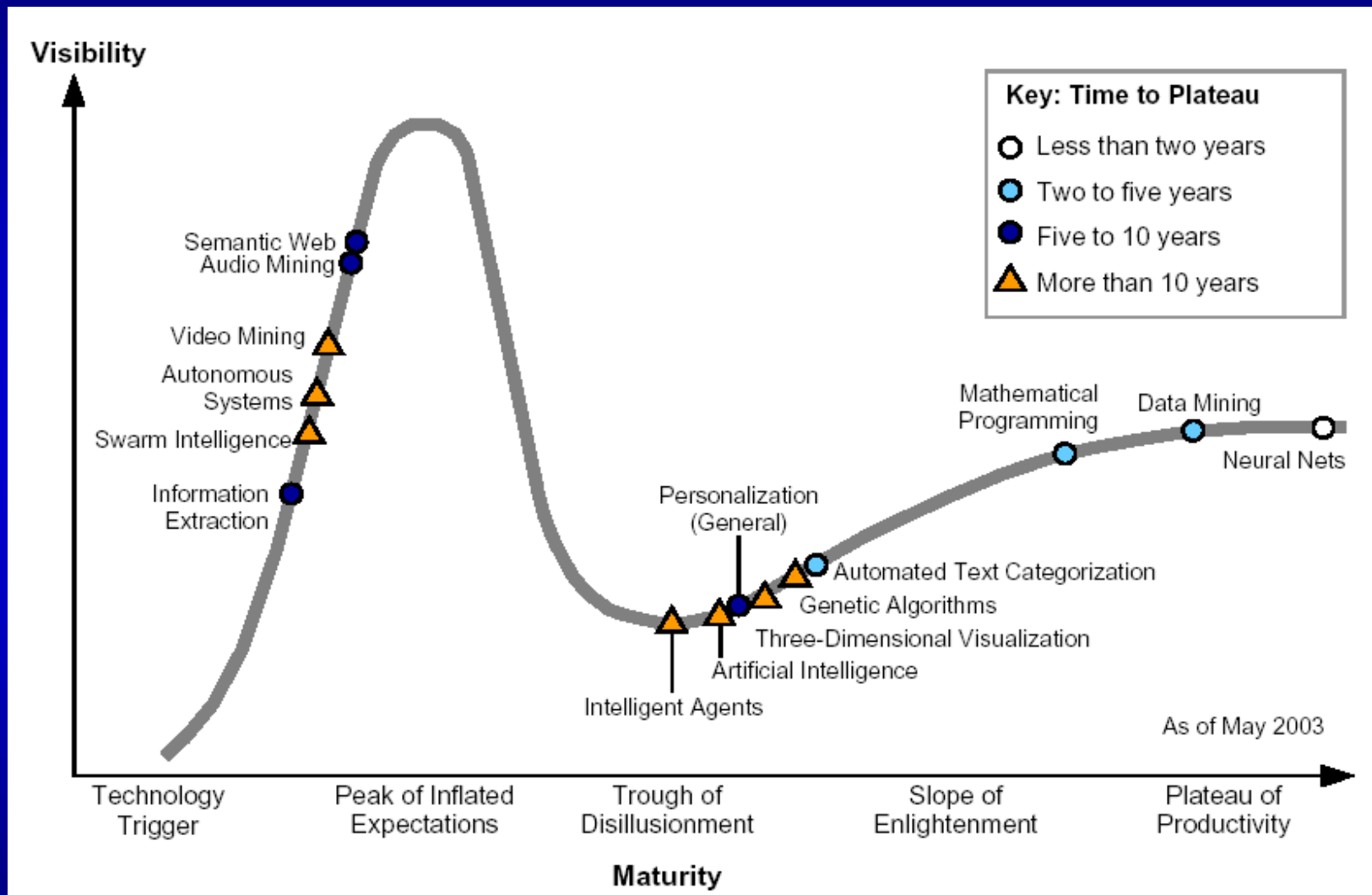
DR. JOHN ELDER HEADS A DATA MINING CONSULTING TEAM WITH OFFICES IN CHARLOTTESVILLE, VIRGINIA, WASHINGTON DC, MOUNTAIN VIEW, CALIFORNIA AND MANHASSET NEW YORK (WWW.DATAMININGLAB.COM). FOUNDED IN 1995, ELDER RESEARCH, INC. FOCUSES ON FEDERAL, COMMERCIAL, INVESTMENT, AND SECURITY APPLICATIONS OF ADVANCED ANALYTICS, INCLUDING TEXT MINING, STOCK SELECTION, IMAGE RECOGNITION, BIOMETRICS, PROCESS OPTIMIZATION, CROSS-SELLING, DRUG EFFICACY, CREDIT SCORING, RISK MANAGEMENT, AND FRAUD DETECTION.

JOHN OBTAINED A BS AND MEE IN ELECTRICAL ENGINEERING FROM RICE UNIVERSITY, AND A PHD IN SYSTEMS ENGINEERING FROM THE UNIVERSITY OF VIRGINIA, WHERE HE'S AN ADJUNCT PROFESSOR TEACHING OPTIMIZATION OR DATA MINING. PRIOR TO 15 YEARS AT ERI, HE SPENT 5 YEARS IN AEROSPACE DEFENSE CONSULTING, 4 HEADING RESEARCH AT AN INVESTMENT MANAGEMENT FIRM, AND 2 IN RICE'S *COMPUTATIONAL & APPLIED MATHEMATICS* DEPARTMENT.

DR. ELDER HAS AUTHORED INNOVATIVE DATA MINING TOOLS, IS A FREQUENT KEYNOTE SPEAKER, AND WAS CO-CHAIR OF THE 2009 *KNOWLEDGE DISCOVERY AND DATA MINING* CONFERENCE, IN PARIS. JOHN'S COURSES ON ANALYSIS TECHNIQUES - TAUGHT AT DOZENS OF UNIVERSITIES, COMPANIES, AND GOVERNMENT LABS - ARE NOTED FOR THEIR CLARITY AND EFFECTIVENESS. DR. ELDER WAS HONORED TO SERVE FOR 5 YEARS ON A PANEL APPOINTED BY THE PRESIDENT TO GUIDE TECHNOLOGY FOR NATIONAL SECURITY. HIS BOOK ON DATA MINING, WITH BOB NISBET AND GARY MINER, WON THE 2009 PROSE AWARD FOR MATHEMATICS. A BOOK ON ENSEMBLES WITH GIOVANNI SENI WAS PUBLISHED IN 2010.

JOHN IS A FOLLOWER OF CHRIST AND THE PROUD FATHER OF 5.

Data Mining and the Hype Cycle

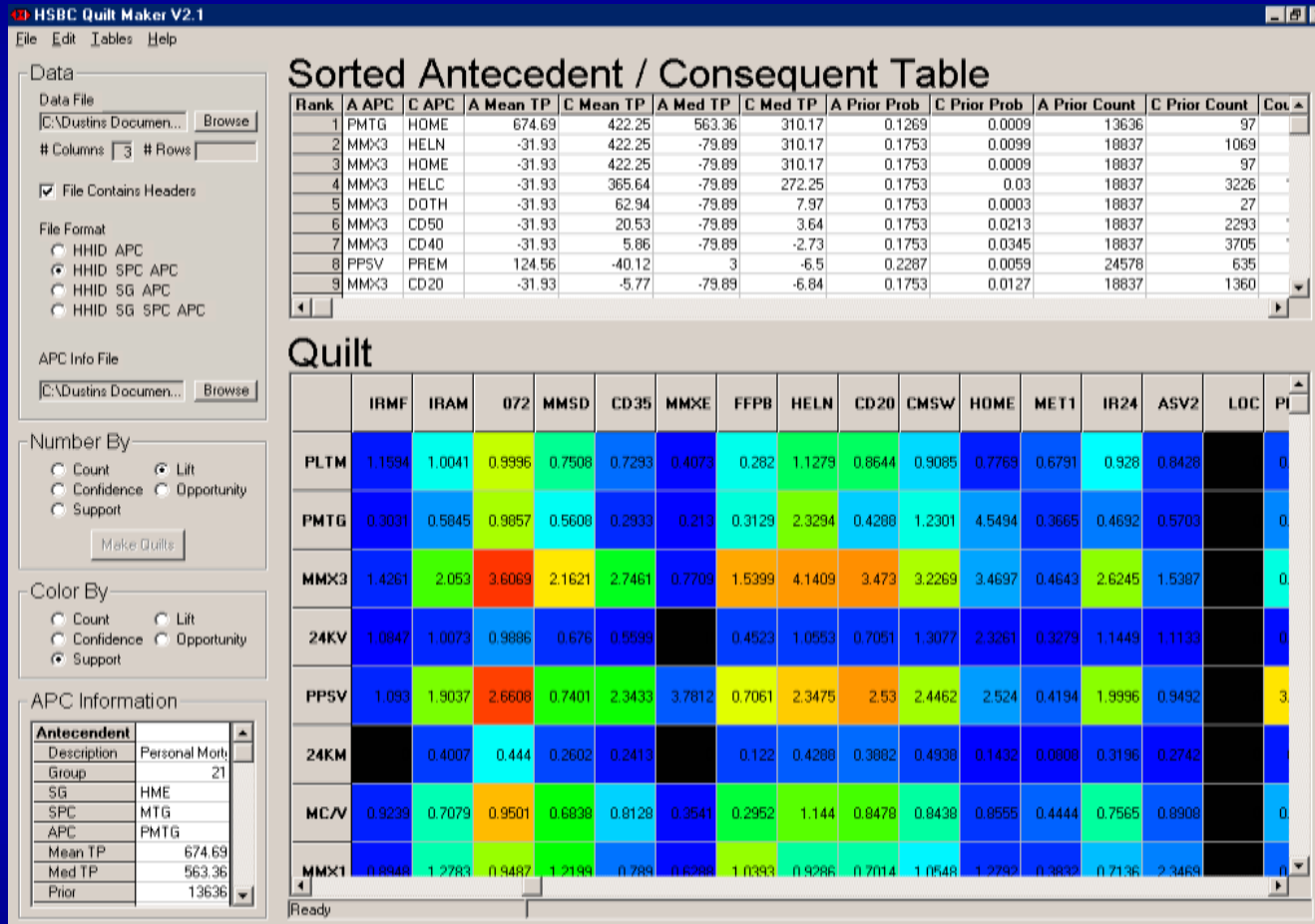


Case 7) HSBC Cross-sell / Up-sell



- Q: What product will interest a customer next?
- Reason: Better target marketing campaigns
- New: Turn contact (a cost) into gain
- Data: transaction history

Case 7) Key Technology: Custom Association Discovery Tool: *QuiltMaker*

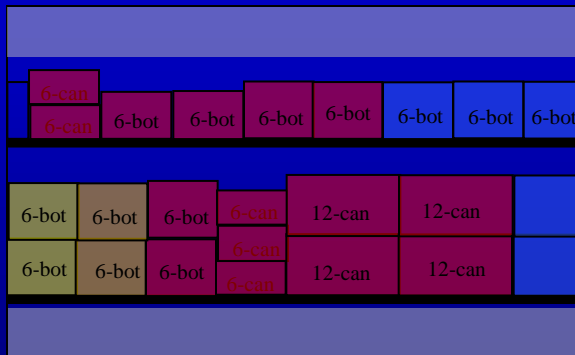


Visualize and Quantify next-most-interesting profitable product

Case 8) Anheuser-Busch Image Recognition



- Q: How do A-B products appear in the store?
- Reason: Discover and implement configurations that work
- Currently: Takes 4 hours to record configurations
- Goal: automate
- Data: images



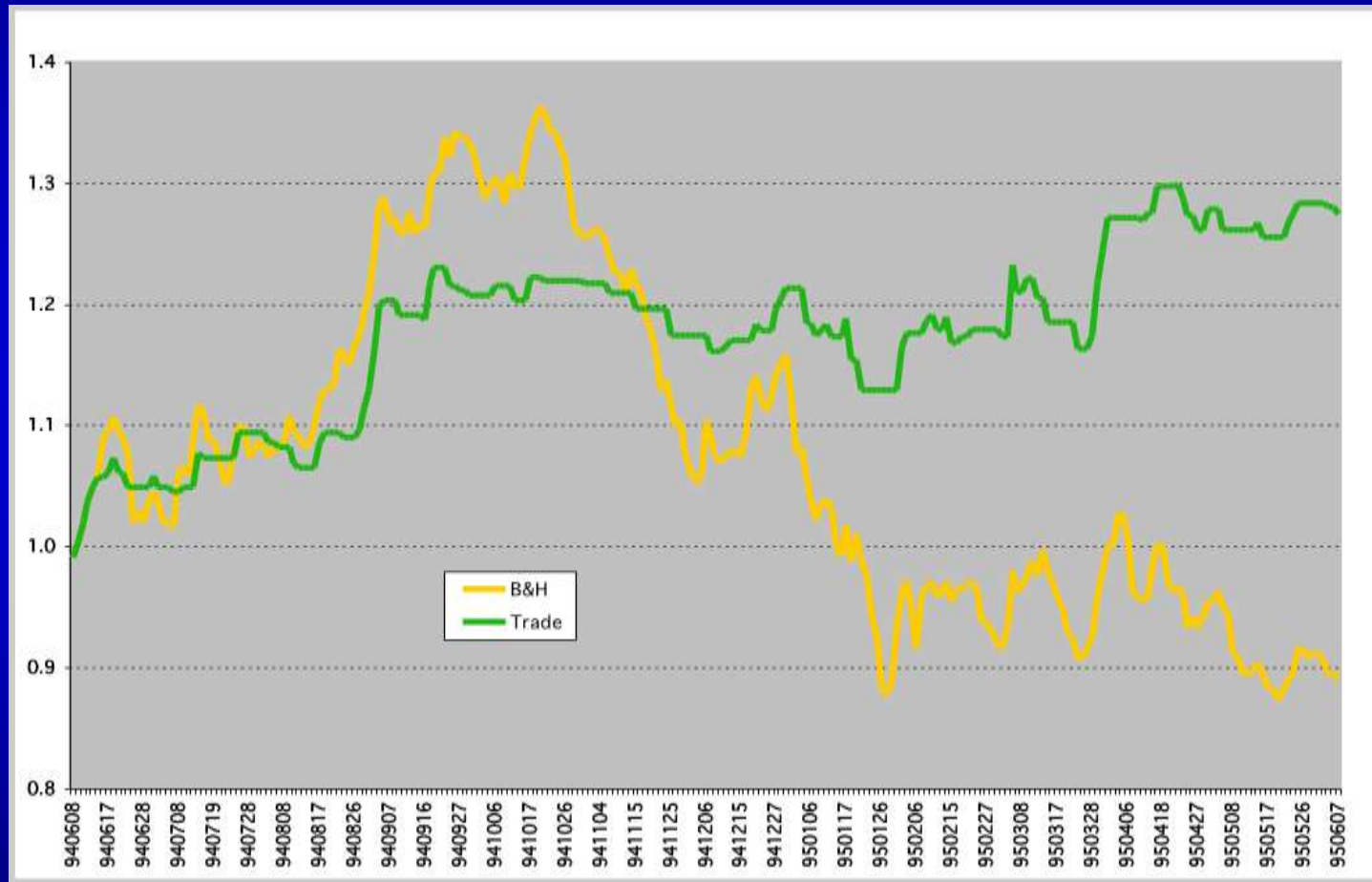


- Automation 90% accurate
-> 10-fold speed-up
- Id'd “stock outs” and “competitor creep”
- Key technological breakthrough:
iterative identification

“Planogram”:
symbolic summary of
product
configuration. Used
for planning and
analysis



Case 9) WestWind Foundation: Hedge Fund Strategy to Time Market



Resampling Technology used to determine whether gains were “real”:

READ file “fund_1yr” date position return
MULTIPLY position return trade
SUM trade original
PRINT original

REPEAT 1000
 SHUFFLE position pos
 MULTIPLY pos return trade
 SUM trade total
 SCORE total Z
END
HISTOGRAM Z

COUNT $Z > \text{original better}$
DIVIDE better 1000 prop_bet
PRINT prop_bet

