

Tim Menzies, WVU, USA
Forrest Shull, Fraunhofer, USA
(with John Hoskings, UoA, NZ)
Jan 27-2011



Empirical Software Engineering, Version 2.0

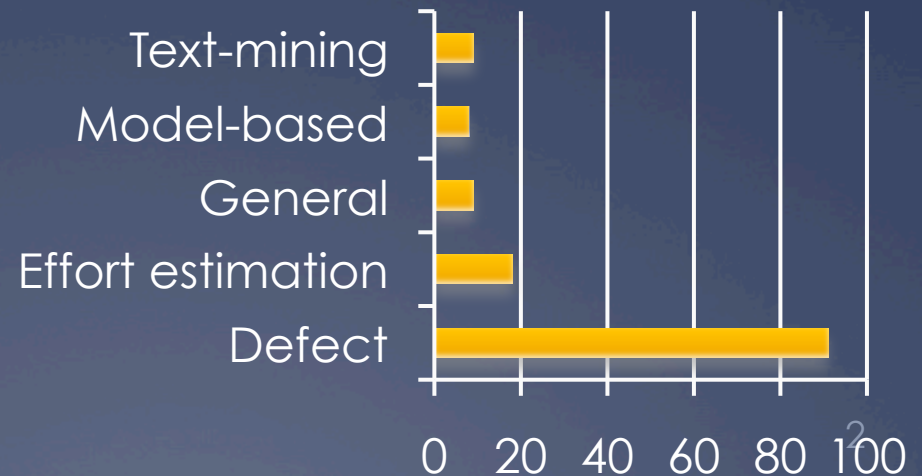
About us



- * Curators of large repositories of SE data
 - * Searched for conclusions
- * Shull: NSF-funded CeBase 2001- 2005
 - * No longer on-line
- * Menzies: PROMISE 2006-2011
 - * If you publish, offer data used in that pub
 - * <http://promisedata.org/data>



- * Our question:
 - * What's next?



Summary

- * We need to do more “data mining”
 - * Not just on different projects
 - * But again and again on the same project
- * And by “data Mining” we really mean
 - * Automated agents that implement
 - * prediction
 - * monitoring
 - * diagnosis,
 - * Planning
 - * Adaptive business intelligence

Adaptive Business Intelligence

- * learning, and re-learning,
- * How to....
 - * Detect death march project
 - * Repair death march projects
 - * Find best sell/buy point for software artifacts
 - * Invest more (or less) in staff training/dev programs
 - * Prioritize software inspections
 - * Estimate development cost
 - * Change development costs
 - * etc

This talk

- * A plea for industrial partners to join in
- * A roadmap for my next decade of research
 - * Many long term questions
 - * A handful of new results

Data Mining & Software Engineering

So many applications of data mining to SE

* Process data

- * Input: Developer skills, platform stability
- * Output: effort estimation

* Product data

- * Input: static code descriptions
- * Output: defect predictors

* Usage data

- * Input: what is everyone using?
- * Output: recommendations on where to browse next

* Social data

- * Input: e.g. which tester do you most respect?
- * Output: predictions of what bugs gets fixed first

* Trace data

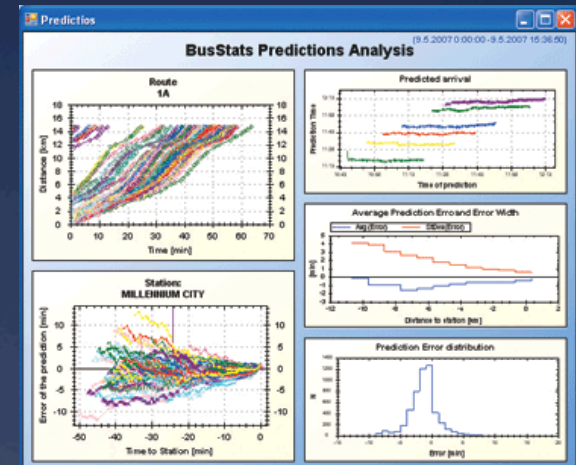
- * Input: what calls what?
- * Output: call sequences that lead to a core dump

* Any textual form

- * Input: text of any artifact
- * Output: e.g. fault localization

The State of the Art

- * If data collected, then usually forgotten
- * Dashboards : visualizations for feature extraction; intelligence left the user
- * MapReduce, Hadoop et. al : systems support for massive, parallel execution.
 - * <http://hadoop.apache.org>
 - * Implements the bus, but no bus drivers
- * Many SE data mining publications
 - * e.g. Bird, Nagappan, Zimmermann and last slide
 - * But, no agents that recognize when old models are no longer relevant,
 - * Or to repair old models using new data



Of course, DM gets it wrong, sometimes

- * Heh, nobody's perfect
- * E.g. look at all the mistakes people make:
 - * Wikipedia: list of cognitive biases
 - * 38 decision making biases
 - * 30 biases in probability
 - * 18 social biases
 - * 10 memory biases
- * At least with DM, can repeat the analysis, audit the conclusion.
- * Create agent communities, each with novel insights and limitations
 - * Data miners working with humans
 - * See more together than separately
 - * Partnership



Does this change
empirical SE
research?

Ben Shneiderman, Mar'08

- * The growth of the World Wide Web ... continues to reorder whole disciplines and industries. ...
- * It is time for researchers in science to take network collaboration to the next phase and reap the potential intellectual and societal payoffs.
 - * -B. Shneiderman.
 - * Science 2.0.
 - * Science, 319(7):1349–1350, March 2008

A proposal



- * Add roller skates to software engineering
- * Always use DM (data mining) on SE data

What's the difference?

SE research v1.0

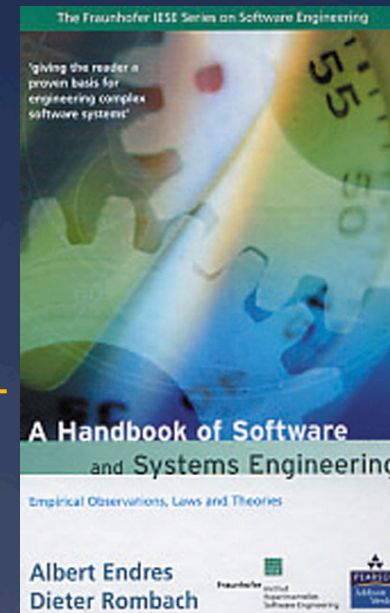
- * Case studies
 - * Watch, don't touch
- * Experiments
 - * Vary a few conditions in a project
- * Simple analysis
 - * A little ANOVA, regression, maybe a t-test

SE research v2.0

- * Data generators
 - * Case studies
 - * Experiments
- * Data analysis
 - * 10,000 of possible data miners
- * Crowd-sourcing
 - * 10,000 of possible analysts
 - *

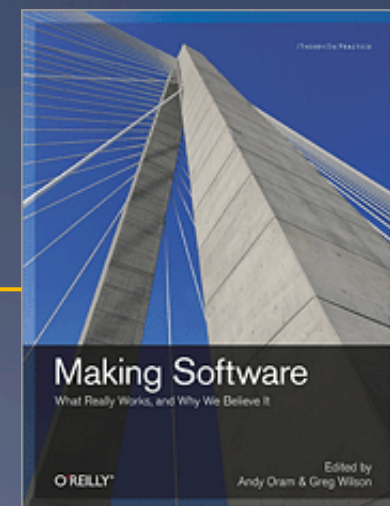
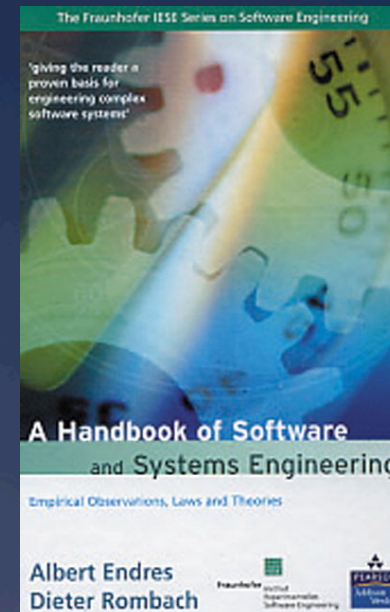
Value-added (to case-study-based research)

- * Case studies: powerful for defining problems, highlighting open issues
- * Has documented 100s of candidate methods for improving SE _____
- * e.g. Kitchenham et. Al IEEE TSE, 2007,
 - * Cross versus Within-Company Cost Estimation
 - * Spawned a sub-culture of researchers
 - * checking if what works here also works there.



Case-Study-based Research: Has Limits

- * Too slow
 - * Years to produce conclusions
 - * Meanwhile, technology base changes
- * Too many candidate methods
 - * No guidance on what methods to apply to particular projects
- * Little generality
 - * Zimmermann et. al, FSE 2009
 - * 662 times : learn here, test there
 - * Worked in 4% of pairs
 - * Many similar no-generality results
 - * Chpt1, Menzies & Shull



Case-studies + DM = Better Research

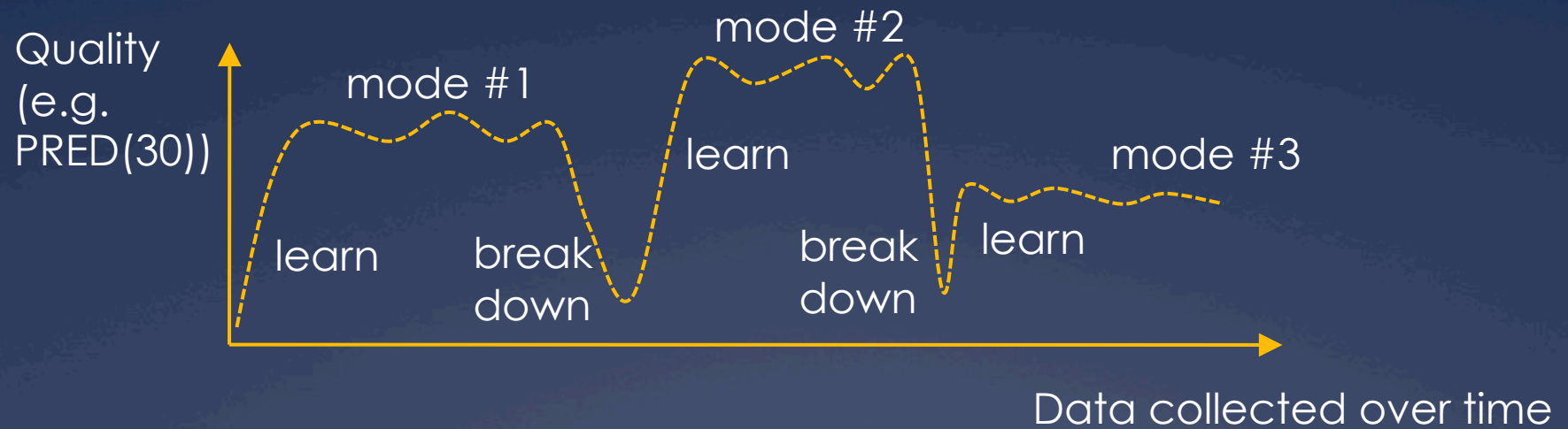
- * Propose a partnership between
 - * case study research
 - * And data mining
- * Data mining is stupid
 - * Syntactic, no business knowledge
- * Case studies are too slow
 - * And to check for generality? Even slower
- * Case study research (on one project) to raise questions
 - * Data mining (on many projects) to check the answers

Adaptive Agents

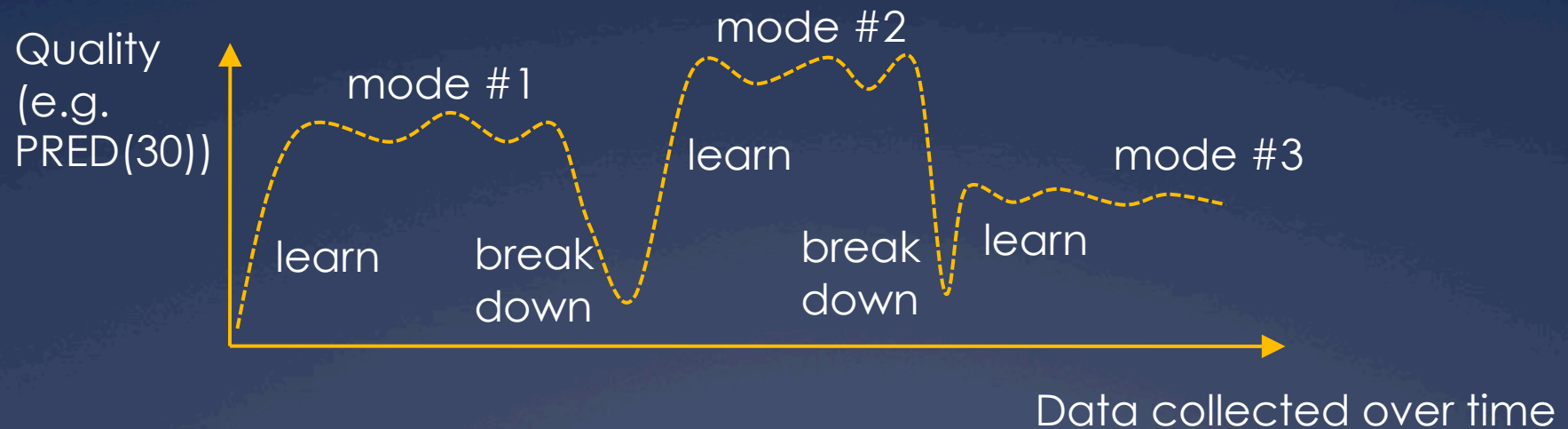
Need for adaptive agents

- * No general rules in SE
 - * Zimmermann FSE, 2009
- * But general methods to find the local rules
- * Issues:
 - * How quickly can we learn the local models?
 - * How to check when local models start failing?
 - * How to repair local models?
- * An adaptive agent watching a stream of data, learning and relearning as appropriate

Agents for adaptive business intelligence

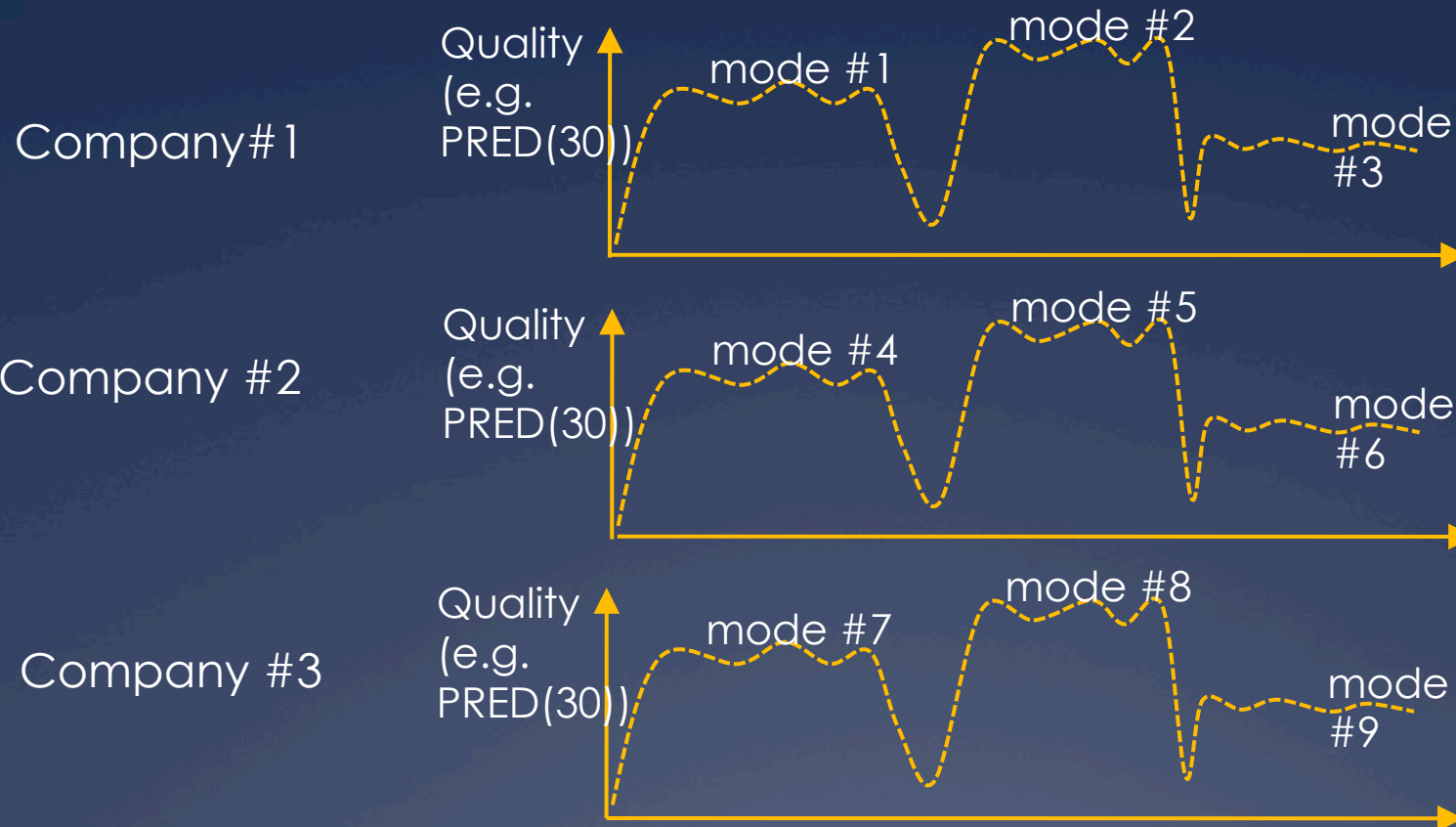


Agents for adaptive business intelligence



- * What is different here?
 - * Not “apply data mining to build a predictor”
 - * But add monitor and repair tools to recognize and handle the breakdown of old predictors
 - * Trust = data mining + monitor + repair

If crowd sourcing



With DM, we could recognize that e.g. 1=4=7

- i.e. when some “new” situation has happened before
- So we'd know what experience base to exploit

Research Questions. How to handle....

- * Anonymization
 - * Make data public, without revealing private data
- * Volume of data
 - * Especially if working from “raw” project artifacts
 - * Especially if crowd sourcing
- * Explanation : of complex patterns
- * Noise: from bad data collection
- * Mode recognition
 - * Is when new stuff is new, or a repeat of old stuff
- * Trust : you did not collect the data
 - * Must surround the learners with assessment agents
 - * Anomaly detectors
 - * Repair

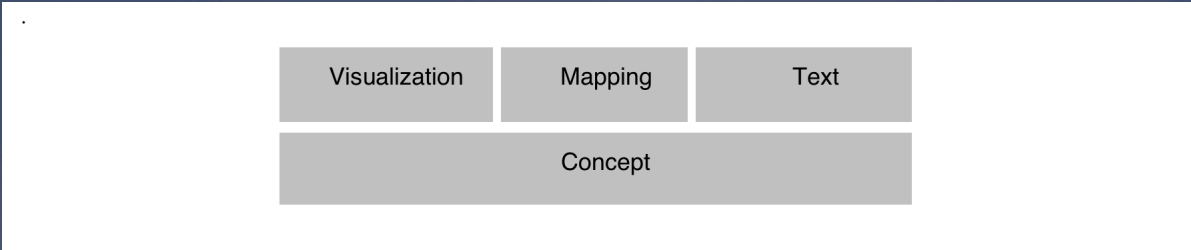
**Most of the technology
required for this approach
can be implemented via
data mining**

**So it would scale
to large data sets**

Organizing the artifacts

Visual Wiki ("Viki") concept

The screenshot shows two web pages side-by-side. The left page is Thinkbase, displaying a complex network graph centered on the movie 'Avatar'. The graph includes nodes for 'James Cameron', 'Gross revenue', 'Notable filming locations', 'James Cameron's Avatar Universe', and 'Avatar: Music from the Motion Picture'. The right page is Freebase, showing a structured data entry for 'Avatar' with fields for 'Initial release date', 'Directed by', 'Rating', 'Runtime', 'Estimated budget', 'Produced by', and 'Screenplay by'. Below this is a section for 'Directed by' featuring a photo and bio of James Cameron.



Enterprise Artifacts example

The image displays three screenshots of the ThinkFree web application interface, illustrating its capabilities in managing enterprise artifacts.

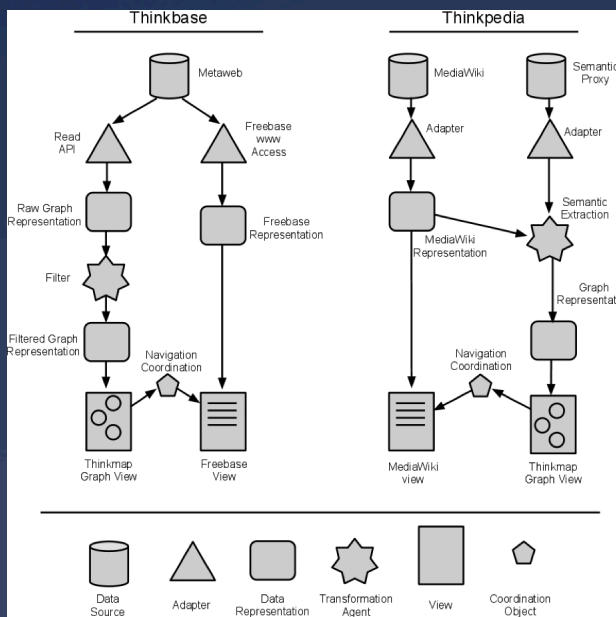
Left Screenshot: Shows a complex network diagram of business functions. The diagram is a hub-and-spoke model with a central node labeled 'Student Management Systems' and various other nodes representing different business functions and services.

Middle Screenshot: Shows a page titled "Student Management Systems" with an overview, business functions, and links to related information. The page includes a search bar, a list of results, and a detailed view of a "child Mega Process".

Right Screenshot: Shows a search interface for "Student" with a list of results. The results include "Student Administration", "Student Services", "Student Recruitment", "Student Cashiering", and "Student Invoicing". A detailed view of a "child Mega Process" is shown, including a table of related items and a "Back to Browse" button.

- * Add documents; organize; search; navigate
- * Edit properties, documents, add links, extract links

VikiBuilder – generating Wikis



The screenshot shows the VikiBuilder v1 web application. The left sidebar contains a list of Visual Wiki Models (VW Lostpedia, VW Thinkbase II, VW Thinkpedia II, VW Thinkpedia II, VikiBuilder II) and Tools (New, Run, [Run Example], Preview). The main area displays a graph for VW Thinkpedia II, showing a flow from VW Wikipedia and VW SemanticProxy through various adapters and representations (Unstructured Representation, Semantic Extractions, Graph Representation) to navigation coordination and TM Thinkpedia graph view. The right sidebar shows the Freebase search interface and the content of VW Thinkpedia II, including facts from the community and a table of data source and adapter information.

from data source	to adapter
VW Wikipedia	Wikipedia Adapter
VW SemanticProxy	SemanticProxy Adapter

from adapter	to data representation
Wikipedia Adapter	Unstructured Representation
SemanticProxy Adapter	Semantic Extractions

from data representation	to transformation agent
Unstructured Representation	Semantic Extractions
Semantic Extractions	Graph Representation

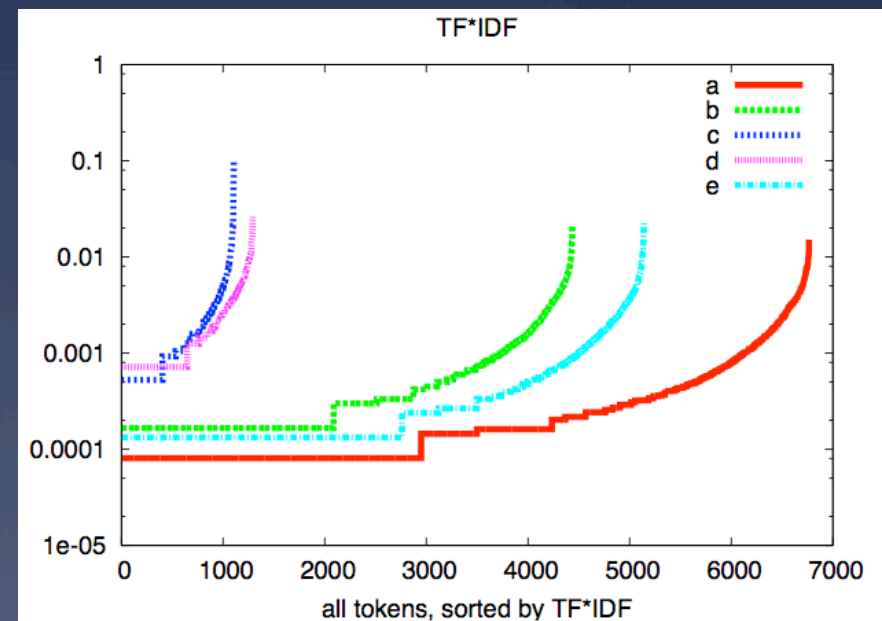
from transformation agent	to data representation
Semantic Extractions	Graph Representation

Text mining

- * Key issue: dimensionality reduction
- * In some domains, can be done in linear time

$$w_{i,j} = tf_{i,j} \times \log\left(\frac{N}{df_i}\right)$$

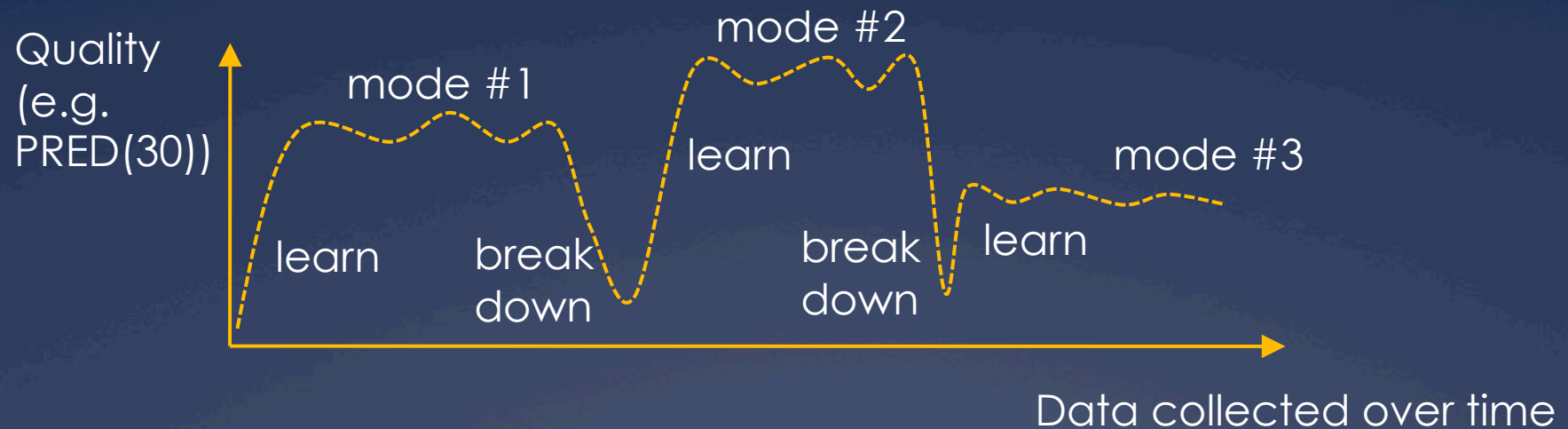
tf_{ij} = number of occurrences of i in j
 df_i = number of documents containing i
 N = total number of documents



- * Use standard data miners, applied to top 100 terms in each corpus

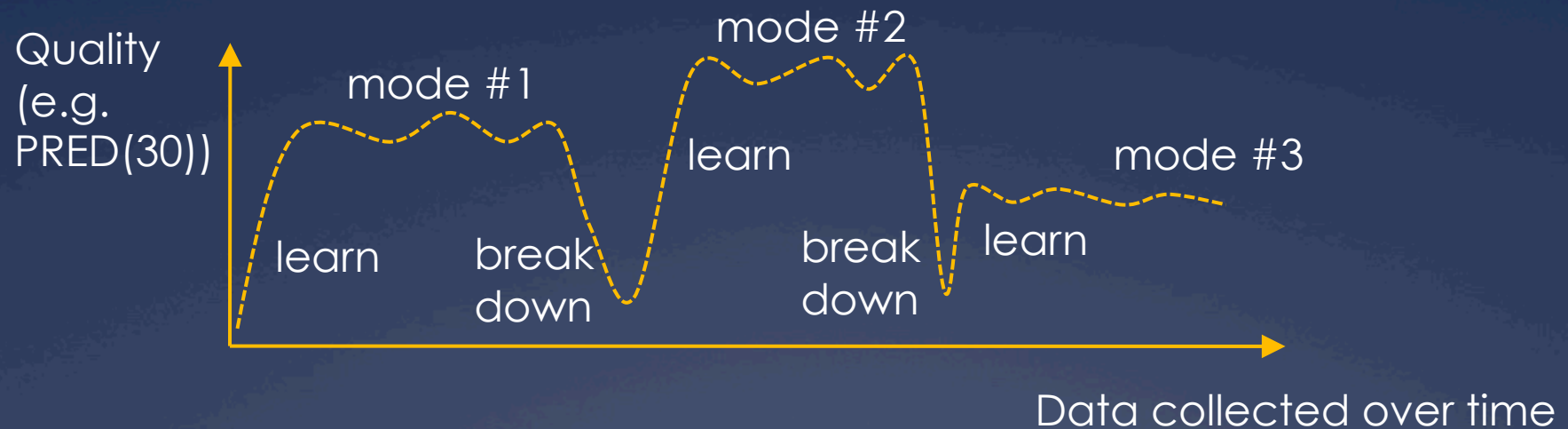
Details

Agents for adaptive business intelligence



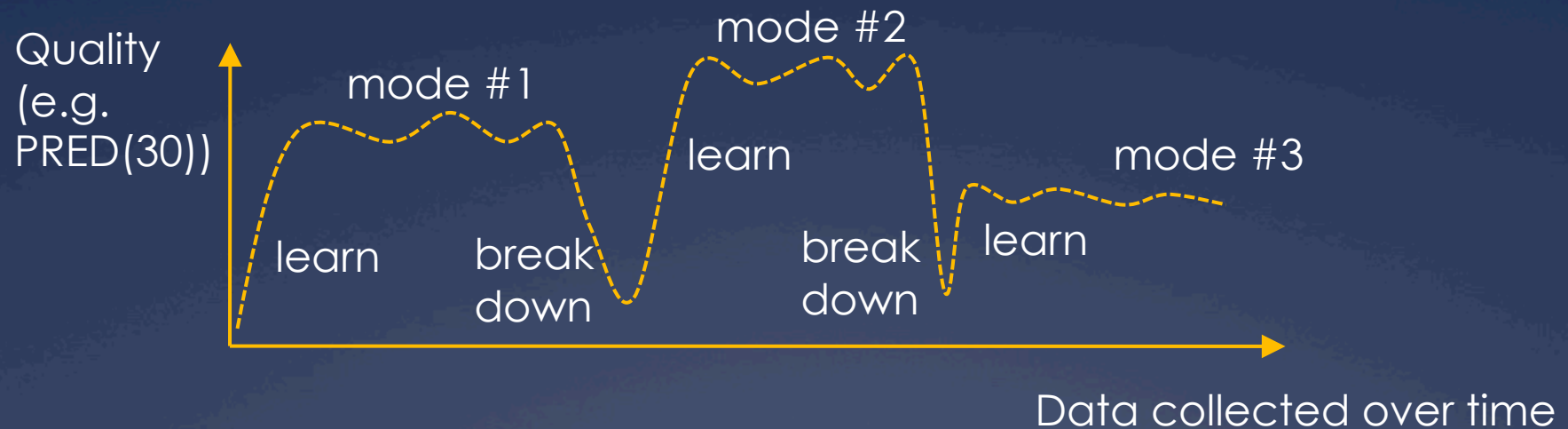
- * Q1: How to learn faster?
 - * Technology: active learning: reflect on examples to date to ask most informative next question
- * Q2: How to recognize breakdown?
 - * Technology: bayesian anomaly detection
 - * Focusing on frequency counts of contrast sets

Agents for adaptive business intelligence



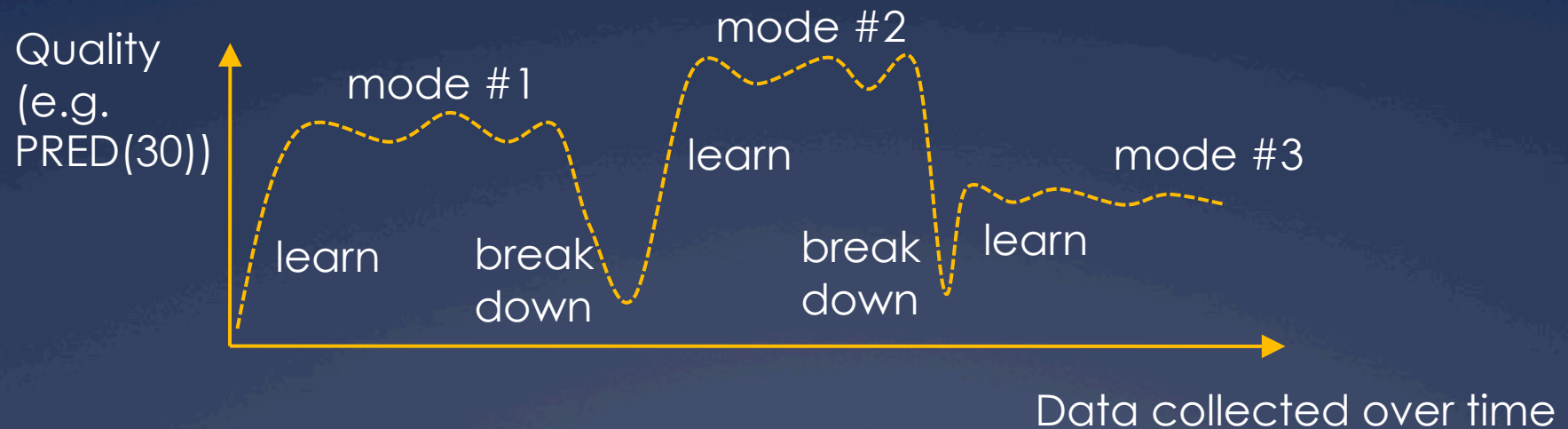
- * Q3: How to classify a mode?
 - * Recognize if you've arrived at a mode seen before
 - * Technology: Bayes classifier
- * Q4: How to make predictions?
 - * Using the norms of a mode, report expected behavior
 - * Technology: table look-up of data inside Bayes classifier

Agents for adaptive business intelligence



- * Q5: What went wrong? (diagnosis)
 - * Delta between current and prior, better, mode
- * Q6: What to do? (planning)
 - * Delta between current and other, better, mode
- * Technology: contrast set learning

Agents for adaptive business intelligence



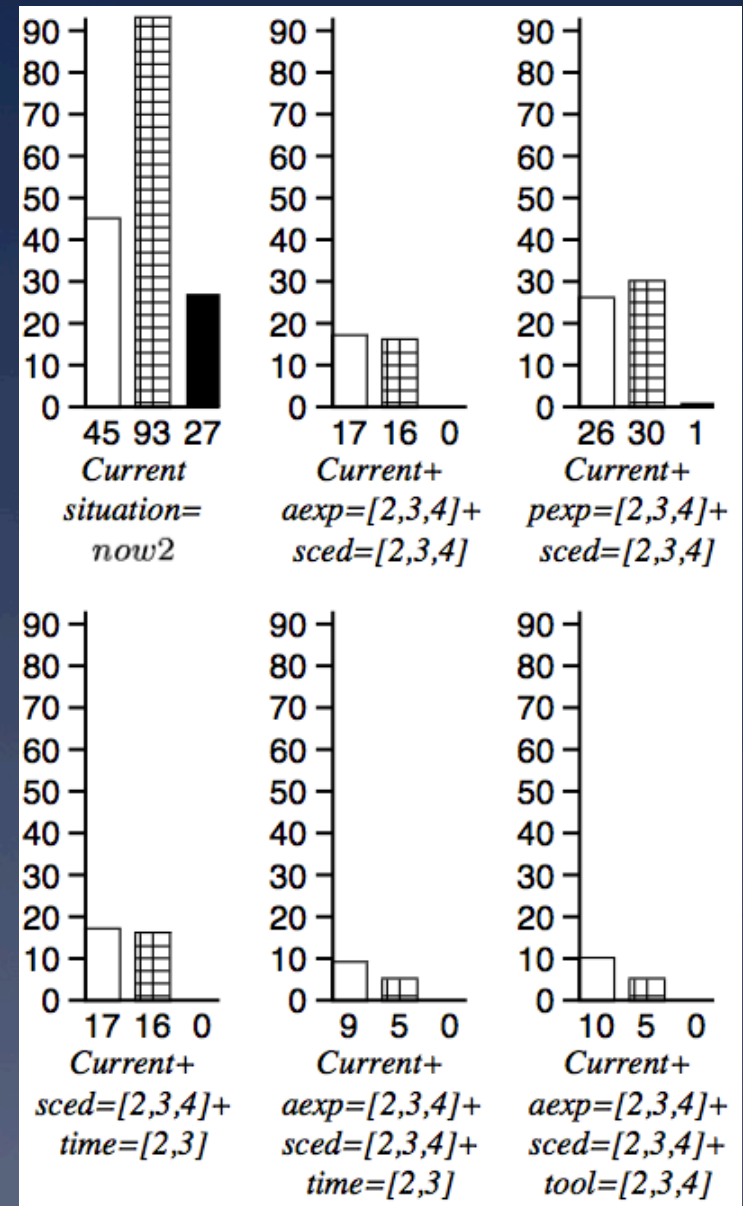
- * Q7: How to understand a mode (explanation)
- * Presentation of essential features of a mode
- * Technology: contrast set learning

Bits and pieces

Prototypes

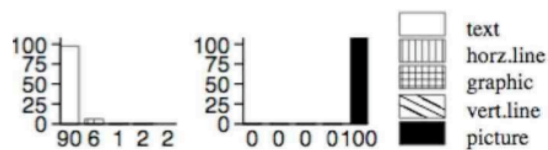
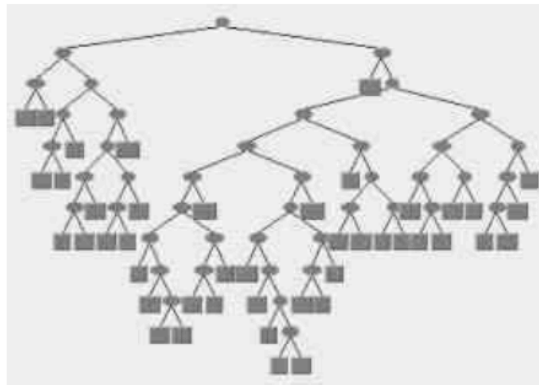
Contrast set learning

- * Minimal contrast set learning = diagnosis & planning
- * A decision tree with weighted leaves
 - * Treatment = decisions that prune branches
 - * Culls bad weights
 - * Keeps good weights
- * E.g. Simulator + C4.5 + 10-way
 - * 10 * 1000 node trees
 - * TAR1: tiny rules: decision on 4 ranges
- * Why so small?
 - * Higher decisions prune more branches
 - * touch fewer nodes



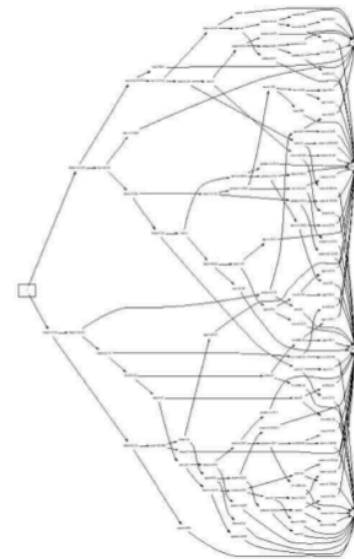
Contrast Set Learning → Succinct Explanations

find graphics on a page from 11 features



$34 \leq \text{height} < 86 \wedge$
 $3.9 \leq \text{mean_tr} < 9.5$

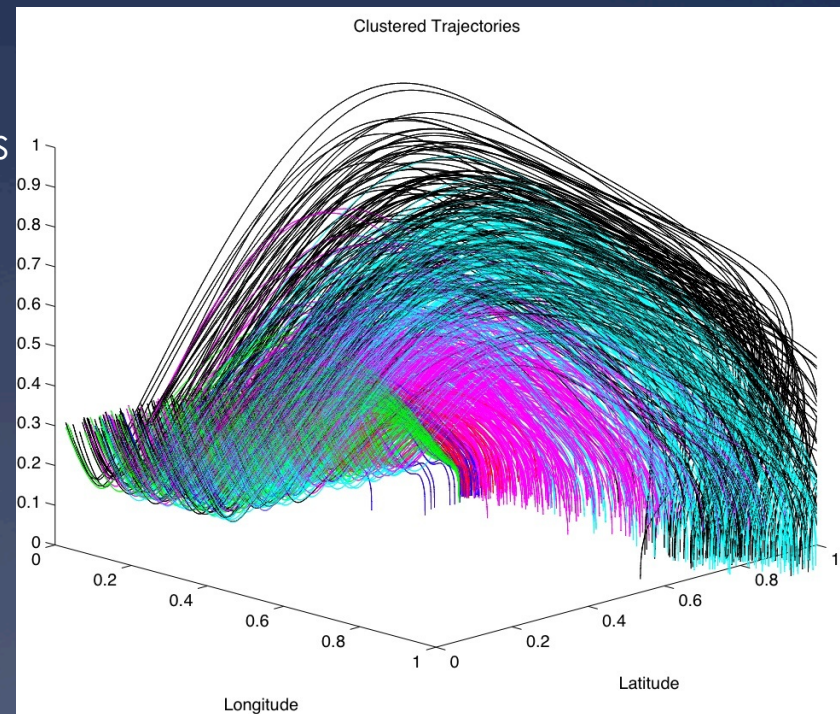
find good housing in Boston



$6.7 \leq \text{RM} < 9.8 \wedge$
 $12.6 \leq \text{PTRATION} < 15.9$

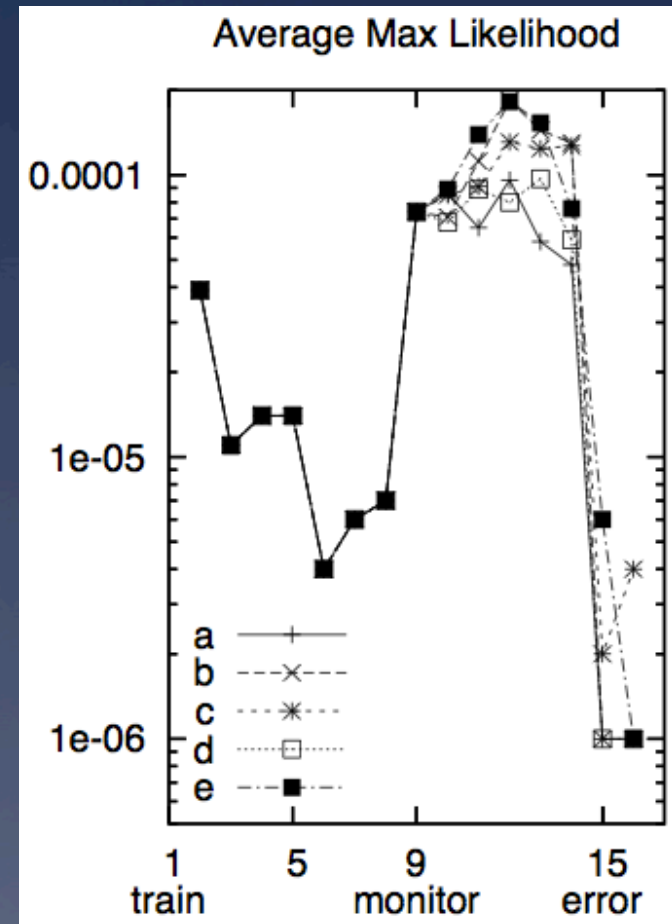
Contrast Set Learning (10 years later)

- * No longer a post-processor to a decision tree learner
 - * TAR3: Branch pruning operators applied directly to discretized data
- * Summer'09
 - * Shoot 'em up at NASA AMES
 - * State-of-the-art numerical optimizer
 - * TAR3
 - * Ran 40 times faster
 - * Generated better solutions
- * Powerful succinct explanation tool



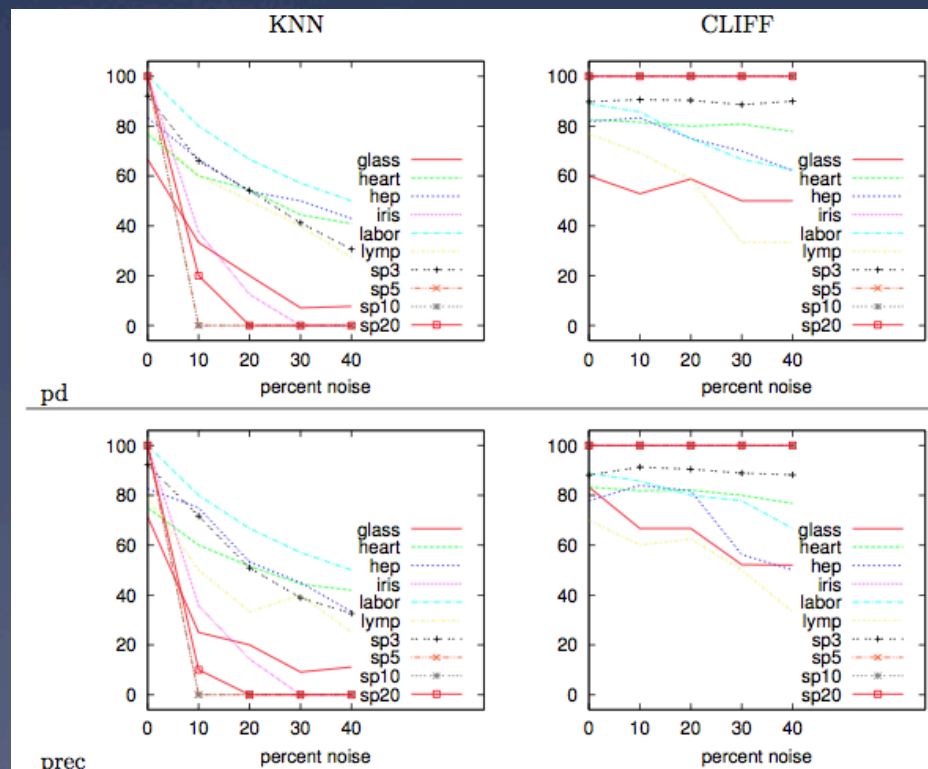
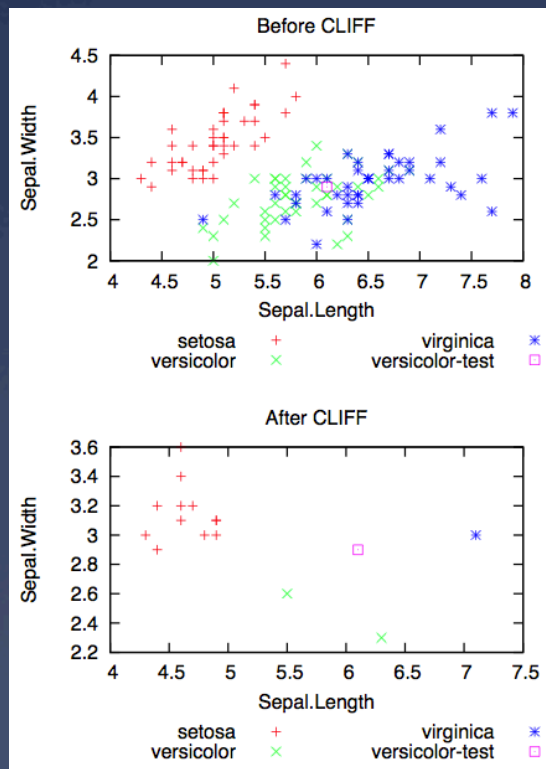
Contrast Set Learning → Anomaly Detection

- * Recognize when old ideas are now out-dated
- * SAWTOOTH:
 - * read data in “eras” of 100 instances
 - * Classify all examples as “seen it”
- * SAWTOOTH1:
 - * Report average likelihood of examples belong to “seen it”
 - * Alert if that likelihood drops
- * SAWTOOTH2:
 - * Back-end to TAR3
 - * Track frequency of contrast sets
 - * Some uniformity between contrast sets and anomaly detection



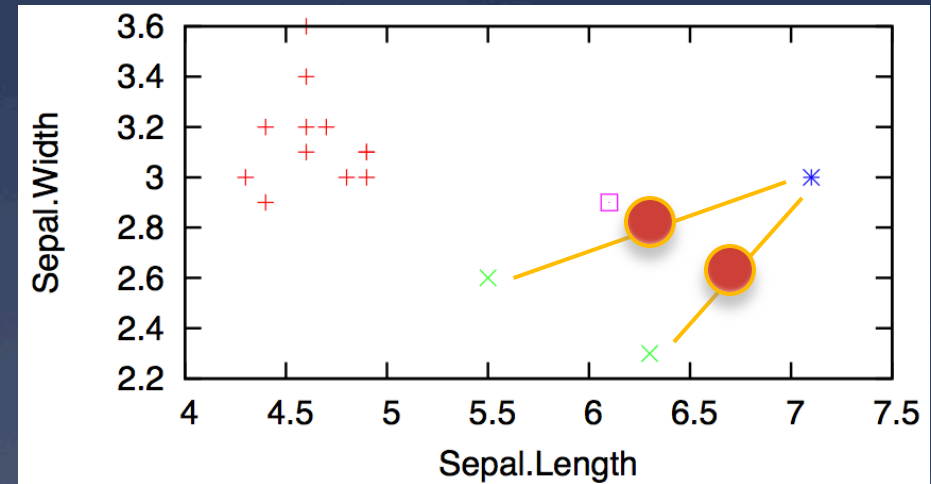
Contrast sets → noise management

- * CLIFF: post-processor to TAR3
 - * Linear time instance selector
- * Finds the attribute ranges that change classification
- * Delete all instances that lack the “power ranges”



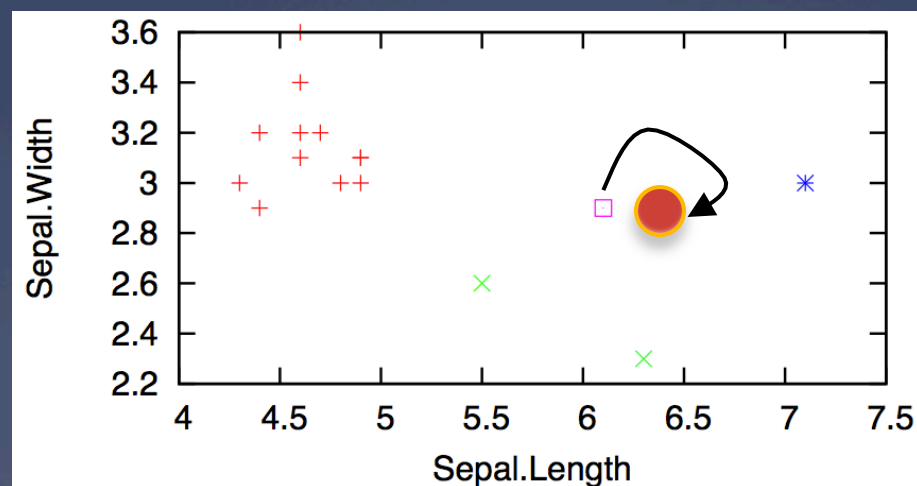
Contrast Sets → CLIFF → Active Learning

- * Many examples, too few labels
- * Reduce the time required for business users to offer judgment on business cases
- * Explore the reduced space generated by CLIFF.
 - * Randomly sample the instances half-way between different classes
- * Fast (in the reduced space)



Contrast sets → CLIFF → Statistical databases

- * Anonymize the data: Preserving its distributions
- * For KNN, that means keep the boundaries between classes
 - * Which we get from CLIFF
- * Also, CLIFF empties out the instance space
 - * Leaving us space to synthesize new instances



And so...

We seek industrial partners

1. That will place textual versions of their products in a wiki
2. That will offer joins of those products to quality measures
3. That will suffer us interviewing their managers, from time to time, to learn the control points.

(Note: 1,2 can be behind your firewalls.)

In return, we offer

- * Agents for
 - * automatic, adaptive, business intelligence
 - * that tunes itself to your local domain

Questions?
Comments?