```
1   TOE
2      TOE = Timm's theory of everything
3      It aims to  simplify knowledge-level modeling with a little
4      data mining.
5
6   KNOWLEDGE-LEVEL PROBLEM SOLVING METHODS
7
8      note
9         the following text references certain terms that aren't
10        explained till below. So just relax and go with the flow.
11
12     anomaly detector (hmmm... that's odd)
13        - walk through data in "eras" of, say, 100 instances
14        - report if median "likelihood(1)" of era i < era[i-1]/2
15
16     verification (do I trust what is going on now?)
17        - alert if any app runs on an "era" with anomalies
18
19     classification (give me the executive summary)
20        - "likelihood(n)"
21
22     mode identification (what is happening now?)
23        - classification using  labels of previous eras
24        - if classification is anomalous, declare a new label
25
26     prediction (what will happen now?)
27        - classification of this era, then return in the current
28          era are the expected values
29
30     planning (how to get there?)
31        - find a "contrast set" between a current and goal era.
32
33     control (how to sail upwards)
34        - find a "contrast set" between a current era and all
35           eras with a higher weight.
36
37     monitor (are we currently smiling?)
38        - classification over the utility labels
39
40     explanation
41        - contrast set between two eras
42
43     diagnosis (how did we go bad?)
44        - explanation, from an eras with a lower to
45          a higher utility
46
47     repair (how can we go good?)
48        - diagnosis, but flip the weights.
49        - also "contrast set" between bad and good,
50        - favoring attributes that have the highest frequency
51          difference and are cheapest to control
52
53     insert your own here
54
```

```
55   FUNCTIONS
56
57      supervised
58         count
59            - build a frequency table for all
60              attribute/range/class  values f[Attr,Range,Class].;
61            - e.g. f[sex,male,pregnant]  = 0
62            - Note that f[class,label,class] is the
63              frequency of class label "Range", which we'll
64              denote f[class] (and "F" is
65              the sum of all "f").
66
67         likelihood(1)
68            - every instance is labeled "seen"
69            - compute likelihood that you have seen this before.
70            - prod(f[a,r,"seen"])/f("seen")*(f("seen")/f) = 1)
71
72         likelihood(N)
73            - every instance is labeled L
74            - compute likelihood that new instance has label L
75            - report label with highest likelihood
76
77         contrast
78            - given two populations
79            - find ranges more frequent in one than the other
80            - for top ranked ranges, try with rule generation
81
82      unsupervised
83         discretization
84            - N bin, equal Fred
85
86         bore (best or rest)
87            - discretization on a numeric utility score
88            - label top score "best" and the others "rest"
89
90         distance
91            - reports distance between two rows
92
93         median
94            - propose a node halfway between two others (for discrete
95              attributes, move half to the other value, at random)
96
97         GAC
98            - builds a tree of nearest pairs
99            - if too slow, use sub/micro sampling as a pre-cursor
100
101     sampling
102        randomizer
103           - randomly re-order rows of the data
104
105        eras
106           - spits our data, X instances at a time
107
108        utility
109           - add a label to each row based on a scoring function
110           - note: simplest one is to just apply the class symbol
111
112        sub-sampling
113           - report all rows of the minority class
114           - use same number of every other class (at random)
115
116        over-sampling
117           - report all rows of the majority class
118           - use same number of every other class (repeat at random)
119
120        micro-sampling
121           - pick N instances (at random) of all classes
122
```

```
123   EXPERIMENT
124
125     hypothesis
126        - once the above is working, the building a whole
127          range of knowledge-level tasks is a trivial process
128
129     tools
130        - we'll need a generator of data to test this all out
131
132     generator
133        sampler(L,P)
134           - ascend levels L in the GAC
135           - find the average distance of things at level L
136             returns random instances within D*L .
137
138        alienator
139           - take classified data
140           - generates eras of the same class frequency
141             as the original data set
142           - at interval I, injects a different frequency
143             classes at probability P
```