

```

1
2
3
4
5
6
7
8
9
10 SAM : Semi-supervised And Multi-objective explanations
11 (c) 2022 Tim Menzies <tim@ieee.org> BSD-2 license
12
13
14
15
16 -- In this code:
17 -- Line strive to be 80 chars (or less)
18 -- Two spaces before function arguments denote optionals.
19 -- Four spaces before function arguments denote local variables.
20 -- Private functions start with '_'
21 -- Arguments of private functions do anything at all
22 -- Local variables inside functions do anything at all
23 -- Arguments of public functions use type hints
24 -- Variable 'x' is anything
25 -- Prefix 'is' is a boolean
26 -- Prefix 'fun' is a function
27 -- Prefix 'f' is a filename
28 -- Prefix 'n' is a string
29 -- Prefix 's' is a string
30 -- Prefix 'c' is a column index
31 -- 'col' denotes 'num' or 'sym'
32 -- 'x' is anything (table or number of boolean or string)
33 -- 'v' is a simple value (number or boolean or string)
34 -- Suffix 's' is a list of things
35 -- Tables are 't' or, using the above, a table of numbers would be 'ns'
36 -- Type names are lower case versions of constructors; e.g 'col' isa 'Cols'.
37
38 -- All demo functions 'eg.fun1' can be called via 'lua eg.lua -e fun1'.
39 local eg = {}
40
41 local l=require"lib"
42 local require"sam"
43 local o,oo,per,push,rnd = l.o,l.oo,l.per,l.push,l.rnd
44 local add,adds,dist,div = _add,_add,_dist,_div
45 local mid, records, the = _mid,_records,_the
46 local Num,Sym = _Num,_Sym
47
48 -- Settings come from big string top of "sam.lua"
49 -- (maybe updated from command line)
50 function eg.the(i) oo(the); return true end
51
52
53 -- The middle and diversity of a set of symbols is called "mode"
54 -- and "entropy" (and the latter is zero when all the symbols
55 -- are the same).
56 function eg.ent( sym,ent)
57   sym= adds(Sym(), {"a","a","a","a","b","b","b","c","c"})
58   ent= div(sym)
59   print(ent,mid(sym))
60   return 1.37 <= ent and ent <= 1.38 end
61
62 -- The middle and diversity of a set of numbers is called "median"
63 -- and "standard deviation" (and the latter is zero when all the nums
64 -- are the same).
65 function eg.num( num)
66   num=Num()
67   for i=1,100 do add(num,i) end
68   local med,ent = mid(num), rnd(div(num),2)
69   print(mid(num),rnd(div(num),2))
70   return 50<= med and med<= 52 and 30.5 <ent and ent <32 end
71
72 -- Nums store only a sample of the numbers added to it (and that storage
73 -- is done such that the kept numbers span the range of inputs).
74 function eg.bignum( num)
75   num=Num()
76   the.nums = 32
77   for i=1,1000 do add(num,i) end
78   oo(_nums(num))
79   return 32==#num._has end
80
81 -- We can read data from disk-based csv files, where row1 lists a
82 -- set of columns names. These names are used to work out what are Nums, or
83 -- ro Syms, or goals to minimize/maximize, or (indeed) what columns to ignore.
84 function eg.records()
85   oo(records("./data/au93.csv").cols.y); return true end
86
87 -- Any two rows have a distance 0..1 that satisfies equality, symmetry
88 -- and the triangle inequality.
89 function eg.dist( data,t)
90   data=records("./data/au93.csv")
91   t={}
92   for i=1,100 do
93     local A,B,C = l.any(data.rows), l.any(data.rows), l.any(data.rows)
94     local a,b,c = dist(data,B,C), dist(data,A,C), dist(data,A,B)
95     assert(a<1 and b<1 and c<1)
96     assert(a>0 and b>0 and c>0)
97     assert( dist(data,A,A) == 0) -- equality
98     assert( dist(data,A,B) == dist(data,B,A)) -- symmetry
99     assert(a+b>=c) -- triangle inequality
100    for _,x in pairs(a) do push(t,rnd(x,2)) end end
101    table.sort(t)
102    oo(t)
103    return true end
104
105
106 the = l.cli(the)
107 os.exit( l.runs(the.eg, eg, the))

```

```

108
109
110
111
112 -- For a list of coding conventions in this file, see
113 -- [eg.lua](https://github.com/timm/lua/blob/main/src/sam/eg.lua).
114 local l=require"lib"
115 local the,settings={}
116 SAM : Semi-supervised And Multi-objective explanations
117 (c) 2022 Tim Menzies <tim@ieee.org> BSD-2 license
118
119 USAGE: lua eg.lua [OPTIONS]
120
121 OPTIONS:
122 -e --eg      start-up example      = nothing
123 -h --help    show help              = false
124 -n --nums    how many numbers to keep = 256
125 -p --p       distance coefficient   = 2
126 -s --seed    random number seed     = 10019]]
127 -- Commonly used lib functions.
128 local o,oo,per,push = l.o,l.oo,l.per,l.push
129
130
131 ----- Classes
132 local Data,Cols,Sym,Num,Row
133 -- Holder of 'rows' and their summaries (in 'cols').
134 function Data() return {cols=nil, -- summaries of data
135   rows={}, -- kept data
136 } end
137
138 -- Holds of summaries of columns.
139 -- Columns are created once, then may appear in multiple slots.
140 function Cols() return {
141   names={}, -- all column names
142   all={}, -- all the columns (including the skipped ones)
143   klass=nil, -- symbolic klass column (if it exists)
144   x={}, -- independent columns (that are not skipped)
145   y={} -- depedent columns (that are not skipped)
146 } end
147
148 -- Summarizes a stream of symbols.
149 function Sym(c,s)
150   return {n=0,at=c or 0, -- items seen
151     at=c or 0, -- column position
152     names=or "", -- column name
153     _has={} -- kept data
154 } end
155
156 -- Summarizes a stream of numbers.
157 function Num(c,s)
158   return {n=0,at=c or 0, names=or "", _has={}, -- as per Sym
159     isNum=true,
160     lo= math.huge, -- lowest seen
161     hi= -math.huge, -- highest seen
162     sorted=true, -- no updates since last sort of data
163     w={s or ""}.find"%S" and -1 or 1 -- minimizing if w=-1
164 } end
165
166 -- Holds one record
167 function Row(t) return {cells=t, -- one record
168   cooked=l.copy(t) -- used if we discretize data
169 } end
170
171 ----- Data Functions
172 local add,adds,clone,div,mid,norm,nums,record,records,stats
173
174 ----- Create
175 -- Generate rows from some 'src'. If 'src' is a string, read rows from file;
176 -- else read rows from a 'src' table. When reading, use row1 to define columns.
177 function records(src, data,head,body)
178   function head(sNames)
179     local cols = Cols()
180     cols.names = sNames
181     for c,s in pairs(sNames) do
182       local col = push(cols.all, -- Numerics start with Uppercase.
183         (s:find"%[A-Z]" and Num or Sym)(c,s))
184       if not s:find"%S" then -- some columns are skipped
185         push(find"%[a-z]" and y or cols.x,col) -- some cols are goal cols
186       if s:find"%S" then cols.klass=col end end end
187     return cols
188   end
189   function body(t) -- treat first row differently (defines the columns)
190     if data.cols then record(data,t) else data.cols=head(t) end
191   end
192   data = Data()
193   if type(src)=="string" then l.csv(src, body) else
194     for _,t in pairs(src or {}) do body(t) end end
195   return data end
196
197 -- Return a new data with same structure as 'data'. Optionally, oad in 'rows'.
198 function clone(data1, rows)
199   data2=Data()
200   data2.cols = _head(data1.cols.names)
201   for _,row in pairs(rows or {}) do record(data2,row) end
202   return data2 end
203
204 ----- Update
205 -- Add one thing to 'col'. For Num, keep at most 'nums' items.
206 function add(col,v)
207   if v=="N" then
208     col.n = col.n + 1
209     if not col.isNum then col._has[v] = 1 + (col._has[v] or 0) else
210       col.lo = math.min(v, col.lo)
211       col.hi = math.max(v, col.hi)
212     local pos
213     if #col._has < the.nums then pos = 1 + (#col._has)
214     elseif math.random() < the.nums/col.n then pos = math.random(#col._has) end
215     if pos then col._has[pos] = tonumber(v) end end end end
216
217 -- Add many things to col
218 function adds(col,t) for _,v in pairs(t) do add(col,v) end; return col end
219
220 -- Add a 'row' to 'data'. Calls 'add()' to update the 'cols' with new values.
221 function record(data,xs)
222   local row= push(data.rows, xs.cells and xs or Row(xs)) -- ensure xs is a Row
223   for _,todo in pairs(data.cols.x, data.cols.y) do
224     for _,col in pairs(todo) do
225       add(col, row.cells[col.at]) end end end
226
227

```

```

228 ----- Query
229 -- Return kept numbers, sorted.
230 function nums(num)
231   if not num.sorted then table.sort(num._has); num.sorted=true end
232   return num._has end
233
234 -- Normalized numbers 0..1. Everything else normalizes to itself.
235 function norm(col,n)
236   return x=="N" or not col.isNum and x or (n-col.lo)/(col.hi-col.lo + 1E-32) end
237
238 -- Diversity (standard deviation for Nums, entropy for Syms)
239 function div(col)
240   if col.isNum then local a=nums(col); return (per(a,9)-per(a,.1))/2.58 else
241     local function fun(p) return p*math.log(p,2) end
242     local e=0
243     for _,n in pairs(col._has) do if n>0 then e=e+fun(n/col.n) end end
244     return e end end
245
246 -- Central tendency (median for Nums, mode for Syms)
247 function mid(col)
248   if col.isNum then return per(nums(col),.5) else
249     local most,mode = -1
250     for k,v in pairs(col._has) do if v>most then mode,most=k,v end end
251     return mode end end
252
253 -- For 'showCols' (default='data.cols.x') in 'data', report 'fun' (default='mid').
254 function stats(data, showCols,fun, t)
255   showCols, fun = showCols or data.cols.y, fun or mid
256   t={}; for _,col in pairs(showCols) do t[col.name]=fun(col) end; return t end
257
258 ----- Distance functions
259 local dist
260 -- Distance between rows (returns 0..1). For unknown values, assume max distance.
261 function dist(data,t1,t2)
262   local function fun(col, v1,v2)
263     if v1=="N" and v2=="N" then return 1 end
264     if not col.isNum then return v1==v2 and 0 or 1 end
265     v1,v2 = norm(col,v1), norm(col,v2)
266     if v1=="N" then v1 = v2<.5 and 1 or 0 end
267     if v2=="N" then v2 = v1<.5 and 1 or 0 end
268     return math.abs(v1-v2)
269   end
270   local d = 0
271   for _,col in pairs(data.cols.x) do
272     d = d + fun(col, t1.cells[col.at], t2.cells[col.at])^the.p end
273   return (d/#data.cols.x)^(1/the.p) end
274
275 ----- That's all folks.
276 return {the=the,
277   Data=Data, Cols=Cols, Sym=Sym, Num=Num, Row=Row,
278   add=add, adds=adds, clone=clone, dist=dist, div=div,
279   mid=mid, nums=nums, records=records, record=record, stats=stats}

```

```

210 211 212 213 214 215 216 217 218 219 220 221 222 223 224 225 226 227 228 229 230 231 232 233 234 235 236 237 238 239 240 241 242 243 244 245 246 247 248 249 250 251 252 253 254 255 256 257 258 259 260 261 262 263 264 265 266 267 268 269 270 271 272 273 274 275 276 277 278 279 280 281 282 283 284 285 286 287 288 289 290 291 292 293 294 295 296 297 298 299 300 301 302 303 304 305 306 307 308 309 310 311 312 313 314 315 316 317 318 319 320 321 322 323 324 325 326 327 328 329 330 331 332 333 334 335 336 337 338 339 340 341 342 343 344 345 346 347 348 349 350 351 352 353 354 355 356 357 358 359 360 361 362 363 364 365 366 367 368 369 370 371 372 373 374 375 376 377 378 379 380 381 382 383 384 385 386 387 388 389 390 391 392 393 394 395 396 397 398 399 400 401 402 403 404 405 406 407 408 409 410 411 412 413 414 415 416 417 418 419 420 421 422 423 424 425 426 427 428 429 430 431 432 433 434 435 436 437 438 439 440 441 442 443 444 445 446 447 448 449 450 451 452 453 454 455 456 457 458 459 460 461 462 463 464 465 466 467 468 469 470 471 472 473 474 475 476 477 478 479 480 481 482 483 484 485 486 487 488 489 490 491 492 493 494 495 496 497 498 499 500 501 502 503 504 505 506 507 508 509 510 511 512 513 514 515 516 517 518 519 520 521 522 523 524 525 526 527 528 529 530 531 532 533 534 535 536 537 538 539 540 541 542 543 544 545 546 547 548 549 550 551 552 553 554 555 556 557 558 559 560 561 562 563 564 565 566 567 568 569 570 571 572 573 574 575 576 577 578 579 580 581 582 583 584 585 586 587 588 589 590 591 592 593 594 595 596 597 598 599 600 601 602 603 604 605 606 607 608 609 610 611 612 613 614 615 616 617 618 619 620 621 622 623 624 625 626 627 628 629 630 631 632 633 634 635 636 637 638 639 640 641 642 643 644 645 646 647 648 649 650 651 652 653 654 655 656 657 658 659 660 661 662 663 664 665 666 667 668 669 670 671 672 673 674 675 676 677 678 679 680 681 682 683 684 685 686 687 688 689 690 691 692 693 694 695 696 697 698 699 700 701 702 703 704 705 706 707 708 709 710 711 712 713 714 715 716 717 718 719 720 721 722 723 724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 750 751 752 753 754 755 756 757 758 759 760 761 762 763 764 765 766 767 768 769 770 771 772 773 774 775 776 777 778 779 780 781 782 783 784 785 786 787 788 789 790 791 792 793 794 795 796 797 798 799 800 801 802 803 804 805 806 807 808 809 810 811 812 813 814 815 816 817 818 819 820 821 822 823 824 825 826 827 828 829 830 831 832 833 834 835 836 837 838 839 840 841 842 843 844 845 846 847 848 849 850 851 852 853 854 855 856 857 858 859 860 861 862 863 864 865 866 867 868 869 870 871 872 873 874 875 876 877 878 879 880 881 882 883 884 885 886 887 888 889 890 891 892 893 894 895 896 897 898 899 900 901 902 903 904 905 906 907 908 909 910 911 912 913 914 915 916 917 918 919 920 921 922 923 924 925 926 927 928 929 930 931 932 933 934 935 936 937 938 939 940 941 942 943 944 945 946 947 948 949 950 951 952 953 954 955 956 957 958 959 960 961 962 963 964 965 966 967 968 969 970 971 972 973 974 975 976 977 978 979 980 981 982 983 984 985 986 987 988 989 990 991 992 993 994 995 996 997 998 999 1000

```

```

399     for _,k in pairs(_egs()) do print(string.format(" %-7s",k)) end
400 elseif k=="all" then -- run all
401     for _,k in pairs(_egs()) do
402         fails=fails + (l.run(k,funcs,settings) and 0 or 1) end
403     elseif funcs[k] then -- run one
404         math.randomseed(settings.seed) -- reset seed
405         local b4={}; for k,v in pairs(settings) do b4[k]=v end
406         local out=funcs[k]()
407         for k,v in pairs(b4) do settings[k]=v end -- restore old settings
408         print("!!!!!!", k, out and "PASS" or "FAIL") end
409     l.rogues()
410     return fails end
411
412 -- -----
413 -- That's all folks.
414 return 1

```