

```

130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
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

```

```

228 ---
229 --- UPDATE COLS
230 ---
231 ---
232 local add
233 function add(i,x, inc)
234   inc = inc or 1
235   if x ~= "?" then
236     i.n = i.n + inc
237     i:internalAdd(x,inc) end
238   return x end
239
240 function Sym.internalAdd(i,x,inc)
241   i.all[x] = inc + (i.all[x] or 0)
242   if i.all[x] > i.most then i.most, i.mode = i.all[x], x end end
243
244 function Num.internalAdd(i,x,inc, d)
245   for j=1,inc do
246     d = x - i.mu
247     i.mu = i.mu + d/i.n
248     i.m2 = i.m2 + d*(x - i.mu)
249     i.sd = (i.m2<0 or i.n<2) and 0 or ((i.m2/(i.n-1))^0.5)
250     i.lo = math.min(x, i.lo)
251     i.hi = math.max(x, i.hi)
252     if #i.all < the.keep then push(i.all,x)
253     elseif r() < they.keep/i.n then i.all[r(#i.all)]=x end end end
254
255 ---
256 --- MAKE DATA
257 ---
258 local file2Egs -- not "local data" (since defined above)
259 function data(i,row)
260   push(i.all, row)
261   for _,col in pairs(i.cols) do add(col, row[col.at]) end
262   return i end
263
264 function file2Egs(file, i)
265   for row in file2things(file) do
266     if i then data(i,row) else i = Egs(row) end end
267   return i end
268
269 ---
270 --- SUMMARIZE
271 ---
272 local mids
273 function mids(i,rows,cols) return i:clone(rows):mid(cols) end
274
275 function Egs.mid(i,cols)
276   return map(cols or i.y,function(col) return col:mid() end) end
277
278 function Sym.mid(i) return i.mode end
279 function Num.mid(i) return i.mu end
280
281 function Num.div(i) return i.sd end
282 function Sym.div(i, e)
283   e=0, for _,n in pairs(i.all) do e=e + n/i.n*math.log(n/i.n,2) end
284   return -e end
285
286 ---
287 --- DISTANCE
288 ---
289 local far,furthest,neighbors,dist
290 function far( i,r1,rows,far)
291   return per(neighbors(i,r1,rows),far or the.far)[2] end
292
293 function furthest( i,r1,rows)
294   return last(neighbors(i,r1,rows))[2] end
295
296 function neighbors(i,r1,rows)
297   return sort(map(rows, function(r2) return {dist(i,r1,r2),r2} end),firsts) end
298
299 function dist(i,row1,row2, d,n,a,b,inc)
300   d,n = 0,0
301   for _,col in pairs(i.x) do
302     a,b = row1[col.at], row2[col.at]
303     inc = a=="?" and b=="?" and 1 or col:dist1(a,b)
304     d = d + inc^the.p
305     n = n + 1 end
306   return (d/n)^(1/the.p) end
307
308 function Sym.dist1(i,a,b) return a==b and 0 or 1 end
309
310 function Num.dist1(i,a,b)
311   if a=="?" then b:=norm(b); a=b<.5 and 1 or 0
312   elseif b=="?" then a:=norm(a); b=a<.5 and 1 or 0
313   else a,b = i:norm(a), i:norm(b) end
314   return math.abs(a - b) end
315
316 function Num.norm(i,x)
317   return i.hi - i.lo < 1E-32 and 0 or (x - i.lo)/(i.hi - i.lo) end
318
319 ---
320 --- CLUSTER
321 ---
322 local half, cluster, clusters
323 function half(i, rows, project,row,some,left,right,lefts,rights,c,mid)
324   function project(row,a,b)
325     a = dist(i,left,row)
326     b = dist(i,right,row)
327     return {(a^2 + c^2 - b^2)/(2*c), row}
328   end
329   some = many(rows, the.some)
330   left = furthest(i,any(some), some)
331   right = furthest(i,left, some)
332   c = dist(i,left,right)
333   lefts,rights = {},{}
334   for n, projection in pairs(sort(map(rows,project),firsts)) do
335     if n==#rows//2 then mid=row end
336     push(n <= #rows//2 and lefts or rights, projection[2]) end
337   return lefts, rights, left, right, mid, c end
338
339 function cluster(i,rows, here,lefts,rights)
340   rows = rows or i.all
341   here = {all=rows}
342   if #rows >= 2* (#i.all)^the.leaves then
343     lefts, rights, here.left, here.right, here.mid = half(i, rows)
344     if #lefts < #rows then
345       here.lefts = cluster(i,lefts)
346       here.rights = cluster(i,rights) end end
347   return here end
348
349 function clusters(i,format,t,pre, front)
350   if t then
351     pre=pre or ""
352     front = fmt("%s%s",pre,#t.all)
353     if not t.lefts and not t.rights then
354       print(fmt("%-20s%s",front, o(rnds(mids(i,t.all),format))))
355     else
356       print(front)
357       clusters(i,format,t.lefts,"|.. pre)
358       clusters(i,format,t.rights,"|.. pre) end end end

```

```

359 ---
360 --- DISCRETIZE
361 ---
362 ---
363 local merge,merged,spans,bestSpan
364 function Sym.spans(i, j)
365   local xys,all,one,last,x,y,n = {},{}
366   for x,n in pairs(i.all) do push(xys, {x,"lefts",n}) end
367   for x,n in pairs(j.all) do push(xys, {x,"rights",n}) end
368   for _,tmp in ipairs(sort(xys,firsts)) do
369     x,y,n = unpack(tmp)
370     if x ~= last then
371       last = x
372       one = push(all, {lo=x, hi=x, all=Sym(i.at,i.name)}) end
373     add(one.all, y, n) end
374   return all end
375
376 function Num.spans(i, j)
377   local xys,all,lo,hi,gap,one,x,y,n = {},{}
378   lo,hi = math.min(i.lo, j.lo), math.max(i.hi, j.hi)
379   gap = (hi - lo) / (6/the.cohen)
380   for _,n in pairs(i.all) do push(xys, {n,"lefts",1}) end
381   for _,n in pairs(j.all) do push(xys, {n,"rights",1}) end
382   one = {lo=lo, hi=lo, all=Sym(i.at,i.name)}
383   all = {one}
384   for _,tmp in ipairs(sort(xys,firsts)) do
385     x,y,n = unpack(tmp)
386     if one.hi - one.lo > gap
387     then one = push(all, {lo=one.hi, hi=x, all=one.all:clone()}) end
388     one.hi = x
389     add(one.all, y, n) end
390   all = merge(all)
391   all[1].j.lo = -math.huge
392   all[#all].hi = math.huge
393   return all end
394
395 function merge(b4, j,n,now,a,b,both)
396   j, n, now = 0, #b4, {}
397   while j < #b4 do
398     j = j+1
399     a, b = b4[j], b4[j+1]
400     if b then
401       both = a.all:merge(b.all)
402       if both
403       then a = {lo=a.lo, hi=b.hi, all=both}
404       j = j + 1 end end
405     push(now,a) end
406   return #now == #b4 and b4 or merge(now) end
407
408 function Sym.merge(i,j, k,ei,ej,ek)
409   k = i:clone()
410   for x,n in pairs(i.all) do add(k,x,n) end
411   for x,n in pairs(j.all) do add(k,x,n) end
412   ei, ej, ek = i:div(), j:div(), k:div()
413   if ek*.99 <= (i.n*ei + j.n*ej)/k.n then return k end end
414
415 function spans(egs1,egs2, spans,tmp,coll,col2)
416   spans = {}
417   for c,coll in pairs(egs1.x) do
418     col2 = egs2.x[c]
419     tmp = coll:spans(col2)
420     if #tmp> 1 then
421       for _,one in pairs(tmp) do push(spans,one) end end end
422   return spans end
423
424 function bestSpan(spans)
425   local divs,ns,n,div,stats,dist2heaven = Num(), Num()
426   function dist2heaven(s) return {(1 - n(s))^2 + (0 - div(s))^2^.5,s} end
427   function div(s) return divs:norm( s.all:div() ) end
428   function n(s) return ns:norm( s.all.n ) end
429   for _,s in pairs(spans) do
430     add(divs, s.all:div())
431     add(ns, s.all.n) end
432   return sort(map(spans, dist2heaven), firsts)[1][2] end
433
434 ---
435 --- EXPLAIN
436 ---
437 local xplain,xplains,selects,spanShow
438 function xplain(i,rows,used,
439   stop,here,left,right,lefts0,rights0,lefts1,rights1)
440   used=used or {}
441   rows = rows or i.all
442   here = {all=rows}
443   stop = (#i.all)^the.leaves
444   if #rows >= 2*stop then
445     lefts0, rights0, here.left, here.right, here.mid, here.c = half(i, rows)
446     if #lefts0 < #rows then
447       here.selector = bestSpan(spans(i:clone(lefts0),i:clone(rights0)))
448       push(used, {here.selector.all.name, here.selector.lo, here.selector.hi})
449       lefts1,rights1 = {},{}
450       for _,row in pairs(rows) do
451         push(selects(here.selector, row) and lefts1 or rights1, row) end
452       if #lefts1 > stop then here.lefts = xplain(i,lefts1,used) end
453       if #rights1 > stop then here.rights = xplain(i,rights1,used) end end end
454   return here end
455
456 function xplains(i,format,t,pre,how, sel,front)
457   pre, how = pre or "", how or ""
458   if t then
459     pre=pre or ""
460     front = fmt("%s%s%s",pre,how, #t.all, t.c and rnd(t.c) or "")
461     if t.lefts and t.rights then print(fmt("%-35s",front)) else
462       print(fmt("%-35s%s",front, o(rnds(mids(i,t.all),format))))
463     end
464     sel = t.selector
465     xplains(i,format,t.lefts, "|.. pre, spanShow(sel,.."")
466     xplains(i,format,t.rights, "|.. pre, spanShow(sel,true) ..:") end end
467
468 function selects(span,row, lo,hi,at,x)
469   lo, hi, at = span.lo, span.hi, span.all.at
470   x = row[at]
471   if x=="?" then return true end
472   if lo==hi then return x==lo else return lo <= x and x < hi end end
473
474 function spanShow(span, negative, hi,lo,x,big)
475   if not span then return "" end
476   lo, hi, x, big = span.lo, span.hi, span.all.name, math.huge
477   if not negative
478   then if hi == big then return fmt("%s== %s",x,lo) end
479   if hi == big then return fmt("%s>= %s",x,lo) end
480   if lo == -big then return fmt("%s<= %s",x,hi) end
481   return fmt("%s<= %s< %s",lo,x,hi)
482   else if lo == hi then return fmt("%s!= %s",x,lo) end
483   if hi == big then return fmt("%s< %s",x,lo) end
484   if lo == -big then return fmt("%s>= %s",x,hi) end
485   return fmt("%s< %s and %s>= %s", x,lo,x,hi) end end

```

```

486 ---
487 ---
488 ---
489 ---
490 ---
491
492 function Demo.the() oo(the) end
493
494 function Demo.many(a)
495   a={1,2,3,4,5,6,7,8,9,10}; ok("{1023}" == o(many(a,3)), "manys") end
496
497 function Demo.egs()
498   ok(5140==file2Egs(the.file).y[1].hi,"reading") end
499
500 function Demo.dist(i)
501   i = file2Egs(the.file)
502   for n,row in pairs(i.all) do print(n,dist(i, i.all[1], row)) end end
503
504 function Demo.far( i,j,row1,row2,row3,d3,d9)
505   i = file2Egs(the.file)
506   for j=1,10 do
507     row1 = any(i.all)
508     row2 = far(i,row1, i.all, .9)
509     d9 = dist(i,row1,row2)
510     row3 = far(i,row1, i.all, .3)
511     d3 = dist(i,row1,row3)
512     ok(d3 < d9, "closer far") end end
513
514 function Demo.half( i, lefts, rights)
515   i = file2Egs(the.file)
516   lefts, rights = half(i, i.all)
517   oo(mids(i, lefts))
518   oo(mids(i, rights))
519   end
520
521 function Demo.cluster( i)
522   i = file2Egs(the.file)
523   clusters(i, "%0f", cluster(i)) end
524
525 function Demo.spans( i, lefts, rights)
526   i = file2Egs(the.file)
527   lefts, rights = half(i, i.all)
528   oo(bestSpan(spans(i:clone(lefts), i:clone(rights)))) end
529
530 function Demo.xplain( i,j,tmp,lefts,rights,used)
531   i = file2Egs(the.file)
532   used={}
533   xplains(i, "%0f", xplain(i, i.all, used))
534   map(sort(used, function(a,b)
535     return ((a[1] < b[1]) or
536            (a[1]==b[1] and a[2] < b[2]) or
537            (a[1]==b[1] and a[2]==b[2] and a[3] < b[3]))end), oo) end
538
539
540 -----
541 the = settings(help)
542 Demo.main(the.todo, the.seed)

```