



"This ain't chemistry.  
This is art."

```
function report( p, out, a, b, c, d, p
    p = function(z) return math.floor
```

```

373 -----
374 ---
375 --- SUPER RANGES
376 ---
377
378 local bin={}
379 function bin.new(id,at,name,lo,hi,n,div)
380   return {id=id,at=at,name=name,lo=lo,hi=hi,n=n,div=div} end
381
382 function bin.show(i,negative)
383   local x,lo,hi,big, s = i.name, i.lo, i.hi, math.huge
384   if negative then
385     if lo==hi then s=fmt("%s=%s",x,lo)
386     elseif hi==big then s=fmt("%s< %s",x,lo)
387     elseif lo==big then s=fmt("%s>= %s",x,hi)
388     else s=fmt("%s< %s and %s>= %s",x,lo,x,hi) end
389   else
390     if lo==hi then s=fmt("%s== %s",x,lo)
391     elseif hi==big then s=fmt("%s>= %s",x,lo)
392     elseif lo==big then s=fmt("%s< %s",x,hi)
393     else s=fmt("%s<= %s< %s",lo,x,hi) end end
394   return s end
395
396 function bin.select(i,row)
397   local x, lo, hi = row[i.at], i.lo, i.hi
398   return x=="?" or lo == hi and lo == x or lo <= x and x < hi end
399
400 function bin.Merges(bins)
401   local j,n,new = 0,length(bins),{}
402   while j <= n do
403     j=j+1
404     a=bins[j]
405     if j < n then
406       b = bins[j+1]
407       if a.hi == b.lo then
408         a.hi = b.hi
409         a.div = (a.div*a.n + b.div*b.n)/(a.n+b.n)
410         a.n = a.n + b.n
411         j = j + 1 end end
412     push(new,a) end
413   return #new < #bins and bin.Merges(new) or bins end
414
415 local _argmin
416 function bin.Xys(xys,at,name)
417   xys = sort(xys, upx)
418   local triviallySmall = the.cohen*(per(xys,.9).x - per(xys,.1).x)/2.56
419   local enoughItems = #xys / the.bins
420   local out = {}
421   _argmin(l,xys, xys, triviallySmall, enoughItems, -math.huge, at.name, out)
422   out[#out].hi = math.huge
423   return out end
424
425 function _argmin(lo, hi, xys, triviallySmall, enoughItems, b4, at, name, out)
426   local function add(f,z) f[z] = (f[z] or 0) + 1 end
427   local function sub(f,z) f[z] = f[z] - 1 end
428   local lhs, rhs, cut, div, xpect, xy = {},{}
429   for j=lo,hi do add(rhs, xys[j].y) end
430   div = ent(rhs)
431   if hi-lo+1 > 2*enoughItems then
432     for j=lo,hi - enoughItems do
433       add(lhs, xys[j].y)
434       sub(rhs, xys[j].y)
435       local n1,n2 = j - lo + 1, hi - j
436       if n1 > enoughItems and n2 > enoughItems and
437         xys[j+1].x == xys[j+1].x and -- there is a break here
438         xys[j].x - xys[j+1].x > triviallySmall and
439         xys[j].x - xys[j].x > triviallySmall
440       then xpect = (n1*ent(lhs) + n2*ent(rhs)) / (n1+n2)
441         if xpect < div then -- cutting here simplifies things
442           cut, div = j, xpect end end end
443     end -- end if
444     if cut
445     then b4 = _argmin(lo, cut, xys, triviallySmall, enoughItems, b4, at, name, out)
446     else b4 = _argmin(cut+1, hi, xys, triviallySmall, enoughItems, b4, at, name, out)
447     else -- if no cut then the original div was never updates and is still correct
448       b4 = push(out, bin.new(out+1,at,name,b4,xys[hi].x, hi-lo+1,div)).hi end
449     return b4 end
450
451 function nb3(data, log)
452   local tmp, xnums = {}
453   local function discretize(c,x, col)
454     if x ~= "?" then
455       col = xnums[c]
456       if col then
457         for _,one in pairs(col.bins) do
458           if one.lo <= x and x <= one.hi then return one.id end end end end
459       return x end
460     local function xnum(c,name)
461       if ako.xnum(name) then return {name=name, xys={},bins={}} end end
462       local function train(c,x,row)
463         if xnums[c] and x ~= "?" then push(xnums[c].xys, {x=x,y= row[#row]}) end end
464       -- start
465       for row in items(data) do
466         push(tmp,row)
467         if xnums then collect(row, function(c,x) return train(c,x,row) end)
468         else xnums = collect(row,xnum) end end
469       for where,col in pairs(xnums) do col.bins = bin.Xys(col.xys,where); print(col.name,col.bins) end
470       for j=2,#tmp do tmp[j] = collect(tmp[j], discretize) end
471       return nb1(tmp)
472     end
473   end
474
475 ---
476 --- COLS
477 local num={}
478 function num.new(at,name)
479   local w = ako.weight(name or "")
480   return {nump=true,indep=false,n=0,at=at or 0,name=name or "",
481     ww=w,lo=big,hi=-big,mu=0,m2=0,sd=0,bins={}} end
482
483 local sym={}
484 function sym.new(at,name)
485   return {nump=false,indep=false,n=0,at=at or 0,name=name or "",
486     has={},most=0,mode=nil} end
487
488 -- update to "add" everywhere
489 function num.add(i,x, d)
490   if x == "?" then
491     i.n = i.n + 1
492     i.lo = math.min(x, i.lo)
493     i.hi = math.max(x, i.hi)
494     d = x - i.mu
495     i.mu = i.mu + d/i.n
496     i.m2 = i.m2 + d*(x - i.mu)
497     i.sd = ((i.m2<0 or i.n<2) and 0) or ((i.m2/(i.n - 1))^0.5) end
498   return x end
499
500 function sym.add(i,x)
501   if x == "?" then
502     i.n = i.n + 1
503     i.has[x] = 1 + (i.has[x] or 0)
504     if i.has[x] > i.most then
505       i.mode,i.most = x,i.has[x] end end
506   return x end
507
508 local summary={}
509 function summary.new(names, i)
510   i = {names={},klass=nil,xy={}, x= {}, y={}}
511   i.names = names
512   for at,name in pairs(names) do
513     local now = (ako.num(name) and num.new or sym.new)(at,name)
514     push(i.xy, now)
515     if not ako.ignore(name) then
516       if not ako.goal(name) then now.indep = true end
517       if ako.klass(name) then i.klass=now end
518       push(now.indep and i.x or i.y, now) end end
519   return i end
520
521 function summary.add(i,row)
522   for _,col in pairs(i.xy) do
523     (col.nump and num or sym).add(col, row[col.at]) end
524   return row end
525
526 function summary.better(i,row1,row2)
527   local s1, s2, n, e = 0, 0, #i.y, math.exp(1)
528   for _,col in pairs(i.y) do
529     local a = norm(col.lo, col.hi, row1[col.at])
530     local b = norm(col.lo, col.hi, row2[col.at])
531     s1 = s1 - e^(col.w * (a - b) / n)
532     s2 = s2 - e^(col.w * (b - a) / n) end
533   return s1 / n < s2 / n end
534
535 -----
536 ---
537 --- NUMS and SYMS
538 ---
539 local egs={}
540 function egs.new(data, i)
541   i = {rows={}, cols=nil}
542   for row in items(data) do
543     if not i.cols then i.cols=summary.new(row) else
544       push(i.rows, summary.add(i.cols,row)) end end
545   return i end
546
547 function egs.mid(i,cols)
548   local function mid(col) return col.nump and col.mu or col.mode end
549   return map(cols or i.cols.y, mid) end
550
551 function egs.div(i,cols)
552   local function div(col) return col.nump and col.sd or ent(col.has) end
553   return map(cols or i.cols.y, div) end
554
555 function egs.clone(old,rows)
556   local i={rows={}, cols=summary.new(old.cols.names)}
557   for _,row in pairs(rows or {}) do summary.add(i.cols,row) end
558   return i end
559
560 function egs.bestRest(i)
561   i.rows = sort(i.rows, function(a,b) return summary.better(i.cols,a,b) end)
562   local n = (#i.rows)^the.best
563   return slice(i.rows, 1, n) -- top n things
564   many( i.rows, n,the.rest, n+1) end -- some sample of the rest
565
566 -----
567 ---
568 --- XPLAIN
569 ---
570 local rule={}
571 function rule.new(bins, t)
572   t = {}
573   for _,one in pairs(bins) do t[one.at]=t[one.at] or {}; push(t[one.at],one) end
574   return {bins=t} end
575
576 function rule.selects(i,row)
577   local function ors(bins)
578     for _,x in pairs(bins) do if bin.select(x,row) then return true end end
579     return false end
580   for at,bins in pairs(i.bins) do if not ors(bins) then return false end end
581   return true end
582
583 function rule.show(i,bins)
584   local cat, order, ors
585   cat = function(t,sep)return table.concat(t,sep) end
586   order= function(a,b) return a.lo < b.lo end
587   ors= function(bins)
588     return cat(map(bin.Merges(sort(bins,order)),bin.show)," or ") end
589   return cat(map(i.bins, ors)," and ") end
590
591 function egs.Contrasts(i, rows1, rows2)
592   local function contrast(col)
593     local function asBin(x,ys, n,div)
594       n,div = ent(ys)
595       return bin.new(id, col.at, col.name, x, x, n, div) end
596     local symbols, xys, x = {},{}
597     for klass,rows in pairs{rows1,rows2} do
598       for _,row in pairs(rows) do
599         x = row[col.at]
600         if x ~= "?" then
601           if not col.nump then inc2(symbols,x,klass) end
602           push(xys, {x=x, y=class}) end end
603       return col.nump and bins(xys, col.at) or collect(symbols, asBin) end
604     local out, tmp = {}
605     for _,col in pairs(i.cols.x) do
606       tmp = contrast(col)
607       if #tmp > 1 then
608         for _,one in pairs(tmp) do push(out, one) end end end
609     return out end
610
611 function egs.xplain(i)
612   best, rest = bestRest(i)
613   return egs.contrasts(i, best,rest) end
614
615 function egs.dist(i,row1,row2)
616   local function sym(_,x,y) return x==y and 0 or 1 end
617   local function num(c,x,y)
618     if x=="?" then y = norm(c.lo, c.hi, y); x=y<.5 and 1 or 0
619     elseif y=="?" then x = norm(c.lo, c.hi, x); y=x<.5 and 1 or 0
620     else x,y = norm(c.lo, c.hi, x), norm(c.lo, c.hi, y) end
621     return math.abs(x-y) end
622   local function dist(c,x,y)
623     return x=="?" and y=="?" and 1 or (c.nump and num or sym)(c,x,y) end
624   local d, n = 0, #i.cols.x
625   for _,c in pairs(i.cols.x) do d= d + dist(c, row1[c.at], row2[c.at])^the.e end
626   return (d/n)^(1/the.e) end
627
628 -----
629 ---
630 --- MISC
631 ---
632 ---
633 ---
634 ---
635 ---
636 ---
637 function per(t,p) return t[ (p or .5)*#t//1 ] end
638
639 function ent(t)
640   local n=0; for _,m in pairs(t) do n = n+m end
641   local e=0; for _,m in pairs(t) do if m>0 then e = e+m/n*math.log(m/n,2) end end
642 end

```

```

642     return -e,n end
643
644 function norm(lo,hi,x) return math.abs(hi-lo)<1E-9 and 0 or (x-lo)/(hi - lo) end
645
646 --- c h i o _ c l <
647
648 function ish(x,y,z) return math.abs(x-y) <= (z or 0.001) end
649
650 local fails=0
651 function ok(test,msg)
652     print("", test and "PASS" or "FAIL",msg or "")
653     if not test then
654         fails = fails+1
655         if the and the.dump then assert(test,msg) end end end
656
657 function rogues()
658     for k,v in pairs(_ENV) do if not b4[k] then print("???",k,type(v)) end end end
659
660 --- count
661
662 function inc(f,a,n) f=f or {};f[a]=(f[a] or 0) + (n or 1) return f end
663 function inc2(f,a,b,n) f=f or {};f[a]=inc(f[a] or {},b,n); return f end
664 function inc3(f,a,b,c,n) f=f or {};f[a]=inc2(f[a] or {},b,c,n);return f end
665
666 function has(f,a) return f[a] or 0 end
667 function has2(f,a,b) return f[a] and has(f[a],b) or 0 end
668 function has3(f,a,b,c) return f[a] and has2(f[a],b,c) or 0 end
669
670 --- lists
671
672 unpack = table.unpack
673
674 function push(t,x) t[1 + #t] = x; return x end
675
676 function map(t, f, u) u={};for k,v in pairs(t) do u[1+#u]=f(v) end;return u end
677 function collect(t,f, u) u={};for k,v in pairs(t) do u[k]=f(k,v)end;return u end
678 function copy(t, u)
679     if type(t) ~= "table" then return t end
680     u={}; for k,v in pairs(t) do u[copy(k)] = copy(v) end; return u end
681
682 function powerset(s)
683     local function aux(s)
684         local t = {}
685         for i = 1, #s do
686             for j = 1, #t do
687                 t[#t+1] = {s[i],unpack(t[j])} end end
688             return t end
689         return sort(aux(s), function(a,b) return #a < #b end) end
690
691 function sort(t,f) table.sort(t,f); return t end
692
693 function upx(a,b) return a.x < b.x end
694 function upl(a,b) return a[1] < b[1] end
695 function downl(a,b) return a[1] > b[1] end
696
697 function slots(t, u)
698     local function public(k) return tostring(k):sub(1,1) ~= "." end
699     u={};for k,v in pairs(t) do if public(k) then u[1+#u]=k end end
700     return sort(u) end
701
702 function any(a,lo,hi)
703     lo,hi = lo or 1, hi or #a; return a[ (lo+(hi-lo)*math.random())/1 ] end
704
705 function many(a,n,lo,hi, u)
706     u={}; for j=1,n do push(u,any(a,lo,hi)) end; return u end
707
708 function slice(a,lo,hi, u)
709     u,lo,hi = {},lo or 1,hi or #a; for j=lo,hi do u[1+#u]=a[j] end; return u end
710
711 --- string '2' things
712
713 function words(s,sep, t)
714     sep="([^\s"..(sep or " ").." ")+)"
715     t={}; for y in s:gmatch(sep) do t[1+#t] = y end; return t end
716
717 function things(s) return map(words(s), thing) end
718
719 function thing(x)
720     x = x:gmatch("%s*(-)%s*$")
721     if x=="true" then return true elseif x=="false" then return false end
722     return tonumber(x) or x end
723
724 function items(src,f)
725     local function file()
726         src,f = io.input(src),f or things
727         return function() x=io.read();if x then return f(x) else io.close(src) end e
728     end
729     local function tbl( x)
730         x,f = 0, f or function(z) return z end
731         return function() if x< #src then x=x+1; return f(src[x]) end end end
732     if src then
733         return type(src) == "string" and file() or tbl() end end
734
735 --- things '2' string
736
737 fmt = string.format
738
739 function oo(t) print(o(t)) end
740
741 function o(t, seen, u)
742     if type(t) ~= "table" then return tostring(t) end
743     seen = seen or {}
744     if seen[t] then return "..." end
745     seen[t] = t
746     local function show1(x) return o(x, seen) end
747     local function show2(k) return fmt("%.5s",k, o(t[k],seen)) end
748     u = #t>0 and map(t,show1) or map(slots(t),show2)
749     return (t.s or "").."["..table.concat(u," ").."]" end
750
751 function dent(t, seen,pre)
752     pre,seen = pre or "", seen or {}
753     if seen[t] then t = "..." end
754     if type(t) ~= "table" then return print(pre .. tostring(t)) end
755     seen[t]=t
756     for _,k in pairs(slots(t)) do
757         local v = t[k]
758         local after = type(v)=="table" and "\n" or "\t"
759         io.write(pre,":",k,after)
760         if type(v)=="table" then dent(v,seen,"|" ..pre) else print(v) end end end
761
762 function rnds(t,f) return map(t, function(x) return rnd(x,f) end) end
763 function rnd(x,f)
764     return fmt(type(x)=="number" and (x==x//1 and f or "%.52f") or "%s",x) end
765
766 --- cli
767
768 function cli(help)
769     local d,used = {},{}
770     help:gsub("(\\-|\\[\\^%$%+)|[%s]+(-[%s%+]|\\n)%s+(\\^%$%+)",
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
777 function (long,key,short,x)
778     assert(not used[short], "repeated short flag ["..short.."]")
779     used[short]=short
780     for n,flag in ipairs(arg) do
781         if flag==short or flag==long then
782             x = x=="false" and true or x=="true" and "false" or arg[n+1] end end
783         d[key] = x==true and true or thing(x) end
784     if d.help then os.exit(print(help)) end
785     return d end
786

```



```
905 -- nb1 and nb2 has "?"  
906 -- nb3 needsa new train.  
907
```