```python
 1    #   <img src="bingo.png" width=250 style="padding-left:20px;" align=right>
 2    # `bingo.py` reads a  CSV file (`-f`) then
 3    # (1) bins the rows into `-B` Bins along `-d` random projections; then
 4    # (2) actively learns by scoring `-a` random bins, then for `-b` iterations,
 5    # extrapolates from 2 examples to label best `y`-guess; then
 6    # (3) `-c` top bin items are labeled for evaluation.
 7    #
 8    # Success here means that  trees learned from (from `a+b` labels) finds stuff
 9    # as good as anything else (after seeing very few labels).
10    #
11    # ### Coding conventions:
12    # - `the` is for CLI config,
13    # - `row` is a list,
14    # - prefix `_` means "private",
15    # - `col` or `c` means `Num` or `Sym`,
16    # - `i` means "self",
17    # - `d,a,n,s` = dict, array, num, str,
18    # - Structs (no classes), since polymorphic code can be shown together,
19    # - `.it` in structs denotes type
20    # - Uppercase functions are constructors (e.g. `Sym`), and their matching
21    #   lowercase names are variables from that constructor (e.g. `sym`),
22    # - `eg_xxxx` are CLI demos (e.g. `--xxx`),
23    # - In CSV input files, uppercase names on row1 denotes numeric; `+`/`-`
24    # - show `y`-goals (others `x`).
25    """
26    bingo.py: stochastic landscape analysis for multi objective reasoning
27    (c) 2025 Tim Menzies, <timm@ieee.org>. MIT license
28
29    Options, with their (defaults):
30
31      -B Bins   number of bins (10)
32      -a a     rows labelled at random during cold start (4)
33      -b b     rows labelled while reflecting on labels seen so far (30)
34      -c c     rows labels while testing the supposed best bin (5)
35      -d dims   number of dimensions (4)
36      -f file   csv file for data (../moot/optimize/misc/auto93.csv)
37      -G Got    directory to cache downloaded data files (~/tmp/moot)
38      -g get    github repo storing example data files (timm/moot)
39      -K Ksee   sample size, when seeking centroids (256)
40      -k k      Bayes hack for rare classes  (1)
41      -m m      Bayes hack for rare frequencies (2)
42      -p p      minkowski coefficient (2)
43      -r rseed  random number seed (1234567891)
44      -z zero   ignore bins with zero items; 0=auto choose (0)
45      -h        show help
46    """
47    import urllib.request, random, math, sys, re, os
48
49    sys.dont_write_bytecode = True
50    pick = random.choice
51    picks = random.choices
52    BIG = 1E32
53
54    ### Command-line  ------------------------------------------------------------
55
56    # Reset slots from CLI flags, matching on first letter of slot.
57    # e.g. `-f file1` sets `d["file"]` to `file1`. If current value is a bolean then
58    # flags reverse old value. e.g. `-v `negates  (e.g.) `d["verbose"]=False`.
59    def cli(d):
60      for k, v in d.items():
61        for c, arg in enumerate(sys.argv):
62          if arg == "-" + k[0]:
63            d[k] = coerce("False" if str(v) == "True" else (
64                          "True" if str(v) == "False" else (
65                          sys.argv[c + 1] if c < len(sys.argv) - 1 else str(v))))
66
67    # String to thing
68    def coerce(x):
69      for what in (int, float):
70        try: return what(x)
71        except Exception: pass
72      x = x.strip()
73      y = x.lower()
74      return (y == "true") if y in ("true", "false") else x
75
76    def eg_h():
77      "print help text"
78      print(__doc__,"\nExamples:")
79      for s,fun in globals().items():
80        if s.startswith("eg__"):
81          print(f" {re.sub('eg_','--',s):>6}  {fun.__doc__}")
82
83    def eg__all():
84      "run all examples"
85      for s,fun in globals().items():
86        if s.startswith("eg__"):
87          if s != "eg__all":
88            print(f"\n# {s} {'-'*40}\n# {fun.__doc__}\n")
89            random.seed(the.rseed)
90            fun()
91
92    ### Settings  ----------------------------------------------------------------
93
94    # Structs with named fields + pretty print.
95    class o:
96      __init__= lambda i, **d: i.__dict__.update(**d)
97      __repr__= lambda i: \
98                (f.__name__ if (f:=i.__dict__.get("it")) else "")+cat(i.__dict__)
99
100   # Parse the `__doc__` string to generate `the` config variable.
101   the= o(**{m[1]: coerce(m[2])
102            for m in re.finditer(r"-\w+\s+(\w+)[^\(]*\(\s*([^)]+)\s*\)", __doc__)})
103
104   def eg__the() -> None:
105     "Print the configuration."
106     print(the)
107
108
```

```python
108   ### Create -------------------------------------------------------------------
109
110   # Update `i` with  multiple things.
111   def inits(things, i): [add(i,thing) for thing in things]; return i
112
113   # Summarize a stream of numbers
114   def Num(init=[], txt=" ",at=0): # -> Num
115     return inits(init,
116                 o(it=Num,
117                   n=0,          # count of items
118                   at=at,        # column position
119                   txt=txt,      # column name
120                   mu=0,         # mean of what what seen
121                   _m2=0,        # second moment (used to find sd)
122                   lo=  BIG,     # lowest seen
123                   hi= -BIG,     # largest
124                   heaven=(0 if txt[-1]=="-" else 1))) # 0,1 = minimize,maximize
125
126   # Summarize a stream of symbols
127   def Sym(init=[], txt=" ",at=0):  # -> Sym
128     return inits(init, o(it=Sym, n=0,      # count of items
129                          at=at,    # column position
130                          txt=txt, # column name
131                          has={})) # hold symbol counts
132
133   # Turn column names into columns (if upper case, then `Num`. Else `Sym`).
134   def Cols(names): # -> Cols
135     cols,x,y = [],[],[]
136     for c,s in enumerate(names):
137       cols += [(Num if s[0].isupper() else Sym)(txt=s,at=c)]
138       if s[-1] != "X": # what to ignore
139         (y if s[-1] in "+-" else x).append(cols[-1])
140     return o(it=Cols,all=cols, # all the columns
141              x=x,        # just the x columns
142              y=y)        # just the y columns
143
144   # Keep some `rows`, summarize them in the `cols`.
145   def Data(init=[]): # -> Data
146     init = iter(init)
147     names = next(init) # column names
148     return inits(init, o(it=Data, n=0,
149                          rows = [],             # contains the rows
150                          cols = Cols(names))) # summaries of the rows
151
152   # Mimic the structure of an existing `Data`. Optionally, add some rows.
153   def clone(data, rows=[]): # -> Data
154     return inits(Data([[col.txt for col in data.cols.all]]), rows)
155
156
```

```
156   ### Read ----------------------------------------------------------------
157
158   # Iterate over rows in file 's'.
159   def csv(s):
160     with open(webdata(s) or s, 'r', newline='', encoding='utf-8') as f:
161       for line in f:
162         yield [coerce(s) for s in line.strip().split(',')]
163
164   # Get data from repo 'the.get'. Keep a local cache of gotten files at 'the.Got'
165   def webdata(fn):
166     if fn.startswith(the.get):
167       cdir = os.path.expanduser(the.Got)
168       os.makedirs(cdir, exist_ok=True)
169       lfn = fn[len(the.get)+1:]
170       lpath = os.path.join(cdir, lfn)
171       if not os.path.exists(lpath):
172         rurl = f"https://github.com/{the.get}/tree/master/{fn}"
173         urllib.request.urlretrieve(rurl, lpath)
174       return lpath
175
176   def eg__csv():
177     "Print csv data."
178     m = 0
179     for n,row in enumerate(csv(the.file)):
180       if n>0: assert int is type(row[0])
181       m += len(row)
182       if n%50==0: print(n,row)
183     assert m==3192
184
185   def eg__cols():
186     "Print csv data."
187     cols = (lbs,acc,mpg) = Cols( next(csv(the.file))).y
188     assert mpg.heaven==1 and lbs.heaven==0 and acc.at==6
189     [print(cat(col)) for col in cols]
190
191
```

```
191   ### Update --------------------------------------------------------------------
192
193   # 'sub' is just 'add'ing -1.
194   def sub(i,v,purge=False): # -> v
195     return add(i, v, inc= -1, purge=purge)
196
197   # If 'v' is unknown, then ignore. Else, update.
198   def add(i, v, inc=1, purge=False): # -> v
199     def _sym(sym,s): # update symbol counts
200       sym.has[s] = inc + sym.has.get(s,0)
201
202     def _data(data,row): # keep the new row, update the cols summaries.
203       if inc < 0:
204         if purge: data.rows.remove(v)
205         [sub(col, row[col.at], inc) for col in data.cols.all]
206       else:
207         data.rows += [[add(col, row[col.at],inc) for col in data.cols.all]]
208
209     def _num(num,n): # update lo,hi, mean and _m2 (used in sd calculation)
210       num.lo = min(n, num.lo)
211       num.hi = max(n, num.hi)
212       if inc < 0 and num.n < 2:
213         num._m2 = num.mu = num.n = 0
214       else:
215         d        = n - num.mu
216         num.mu  += inc * (d / num.n)
217         num._m2 += inc * (d * (n - num.mu))
218
219     if v != "?":
220       i.n += inc
221       (_num if i.it is Num else (_sym if i.it is Sym else _data))(i,v)
222     return v
223
224   def eg__num():
225     "Demo Numerics."
226     g=lambda: random.gauss(10,2)
227     num = Num(g() for _ in range(256))
228     assert 10 < mid(num) < 10.05 and 2 < div(num) < 2.1
229
230   def eg__sym():
231     "Demo Symbolics."
232     sym = Sym("aaaabbc")
233     assert mid(sym) == "a" and 1,37 < div(sym) < 1.38
234
235   def eg__data():
236     "Read data from disk."
237     model = Data(csv(the.file)).cols.x[2]
238     assert 3.69 <div( model) < 3.7
239     assert model.lo == 70 and model.hi == 82
240
241
```

```python
### Reports ---------------------------------------------------------------

def mids(data): return [mid(col) for col in data.cols.all]
def divs(data): return [div(col) for col in data.cols.all]

def mid(col):
  return col.mu if col.it is Num else max(col.has, key=col.has.get)

def div(col):
  def _num(num):
    return (max(num._m2,0)/(num.n - 1))**0.5

  def _sym(sym):
    return -sum(v/sym.n * math.log(v/sym.n, 2) for v in sym.has.values() if v>0)

  return (_num if col.it is Num else _sym)(col)

def eg__addSub():
  head, *rows = list(csv(the.file))
  data = Data([head])
  for row in rows:
    add(data,row)
    if data.n == 50: m0,d0 = mids(data),divs(data)
  for row in rows[::-1]:
    sub(data,row)
    if data.n == 50:
      m1,d1 = mids(data), divs(data)
      assert all(math.isclose(a,b,rel_tol=0.01) for a,b in zip(m0, m1))
      assert all(math.isclose(a,b,abs_tol=0.01) for a,b in zip(d0, d1))

### Bayes -----------------------------------------------------------------

def like(data, row, nall=2, nh=100):
  prior = (data.n + the.k) / (nall + the.k*nh)
  tmp = [pdf(c,row[c.at],prior)
         for c in data.cols.x if row[c.at] != "?"]
  return sum(math.log(n) for n in tmp + [prior] if n>0)

def pdf(col,v, prior=0):
  def _sym(sym,s):
    return (sym.has.get(s,0) + the.m*prior) / (col.n + the.m + 1/BIG)

  def _num(num,n):
    sd = div(num) or 1 / BIG
    var = 2 * sd * sd
    z = (n - num.mu) ** 2 / var
    return min(1, max(0, math.exp(-z) / (2 * math.pi * var) ** 0.5))

  return (_num if col.it is Num else _sym)(col,v)

def eg__bayes():
  data = Data(csv(the.file))
  L = lambda r: round(like(data,r),2)
  F = lambda a: print(' '.join([f"{x:8}" for x in a]))
  assert all(-20 < L(row) < -9 for row in data.rows)
  rows = [[L(row)] + row for row in sorted(data.rows, key=L)[::30]]
  head = ["Like"] + [col.txt for col in data.cols.all]
  report(rows,head,1)
```

```python
### Distance --------------------------------------------------------------

def norm(i,v):
  return v if (v=="?" or i.it is not Num) else (v - i.lo)/(i.hi - i.lo + 1/BIG)

def dist(col,v,w):
  def _sym(_,s1,s2):
    return s1 != s2

  def _num(num,n1,n2):
    n1,n2 = norm(num,n1), norm(num,n2)
    n1 = n1 if n1 != "?" else (0 if n2 > 0.5 else 1)
    n2 = n2 if n2 != "?" else (0 if n1 > 0.5 else 1)
    return abs(n1 - n2)

  return 1 if v=="?" and w=="?" else (_num if col.it is Num else _sym)(col,v,w)

# Returns the i`p-`th root of sum of the x in a (rarraised to `p`).
def minkowski(a):
  total, n = 0, 1 / BIG
  for x in a:
    n += 1
    total += x**the.p
  return (total / n)**(1 / the.p)

# Distance to ideal, measured across y-columns.
def ydist(data, row):
  return minkowski(abs(norm(c,row[c.at]) - c.heaven) for c in data.cols.y)

# Distance between two rows, measured across x-columns.
def xdist(data, row1, row2):
  return minkowski(dist(c,row1[c.at], row2[c.at]) for c in data.cols.x)

# K-means plus plus: k points, usually D^2 distance from each other.
def kpp(data, k=10, rows=None, few=None):
  def D(x, y):
    key = tuple(sorted((id(x), id(y))))
    if key not in mem: mem[key] = xdist(data,x,y)
    return mem[key]

  few = few or the.Ksee
  row, *rows    = shuffle(rows or data.rows)
  some, rest    = rows[:few], rows[few:]
  centroids, mem = [row], {}
  for _ in range(1, k):
    dists = [min(D(x, y)**2 for y in centroids) for x in some]
    r     = random.random() * sum(dists)
    for j, d in enumerate(dists):
      r -= d
      if r <= 0:
        centroids.append(some.pop(j))
        break
  return centroids, mem, some + rest

def eg__ydist():
  data = Data(csv(the.file))
  L = lambda r: round(like(data,r),2)
  Y = lambda r: round(ydist(data,r),2)
  assert all(0 <= Y(row) <= 1 for row in data.rows)
  rows = [[Y(row),L(row)] + row for row in sorted(data.rows, key=Y)[::30]]
  head = ["Y","Like"] + [col.txt for col in data.cols.all]
  report(rows,head,1)

def eg__kpp():
  "Diversity sample: random vs kpp. Try a few times  with -r $RANDOM --kpp."
  data = Data(csv(the.file))
  repeats=20
  Y = lambda row: ydist(data,row)
  best = lambda rows: Y(sorted(rows, key=Y)[0])
  b4 = Num(Y(row) for row in data.rows)
  print("b4  ", o(Ksee=len(data.rows), repeats=1, lo=b4.lo, mu=b4.mu, hi=b4.hi))
  for k in [10,20,30,40,80,160]:
    print("")
    anys = Num(best(picks(data.rows,k=k))      for _ in range(repeats))
    print("random ", o(Ksee=k, repeats=anys.n, lo=anys.lo, mu=anys.mu, hi=anys.hi
, D=0.35*div(anys)))
    kpps = Num(best(kpp(data,         k=k)[0]) for _ in range(repeats))
    print("kpps  ", o(Ksee=k, repeats=kpps.n, lo=kpps.lo, mu=kpps.mu, hi=kpps.hi,
    D=0.35*div(kpps)))
```

```
### Clustering -----------------------------------------------------------------

def project(data, row, a, b): # -> 0,1,2 .. the.bins-1
  D = lambda row1,row2: xdist(data,row1,row2)
  c = D(a,b)
  if c==0: return 0
  return (D(row, a)**2 + c**2 - D(row, b)**2) / (2 * c *c)

def bucket(data,row,a,b):
  return min(int( project(data,row,a,b) * the.bins), the.bins - 1)

def extrapolate(data,row,a,b):
  ya, yb = ydist(data,a), ydist(data,b)
  return ya + project(data,row,a,b) * (yb - ya)

def poles(data): # -> List[Row]
  r0, *some = picks(data.rows, k=the.some + 1)
  out = [max(some, key=lambda r1: xdist(data.r1, r0))]
  for _ in range(the.dims):
    out += [max(some, key=lambda r2: sum(xdist(data,r1,r2) for r1 in out))]
  return out

def lsh(data, corners): # -> Dict[Tuple, List[Row]]
  buckets = {}
  for row in data.rows:
    k = tuple(bucket(row, a, b) for a, b in zip(corners, corners[1:]))
    buckets[k] = buckets.get(k) or clone(data)
    add(buckets[k], row)
  return buckets

def neighbors(c, hi):
  def go(i, p):
    if i == len(c):
      t = tuple(p)
      if t != c and all(0 <= x < hi for x in t):
        yield t
    else:
      for d in [-1, 0, 1]:
        yield from go(i+1, p + [c[i] + d])
  yield from go(0, [])
```

```
### Tree -------------------------------------------------------------------------

ops = {'<=' : lambda x,y: x <= y,
       "==" : lambda x,y: x == y,
       '>'  : lambda x,y: x >  y}

def selects(row, op, at, y): x=row[at]; return  x=="?" or ops[op](x,y)

def cuts(col,rows,Y,Klass):
  def _sym(sym):
    n,d = 0,{}
    for row in rows:
      if (x := row[sym.at]) != "?":
        n = n + 1
        d[x] = d.get(x) or Klass()
        add(d[x], Y(row))
    return o(div = sum(c.n/n * div(c) for c in d.values()),
             hows = [("==",sym.at,k) for k,v in d.items()])

  def _num(num):
    out, b4, lhs, rhs = None, None, Klass(), Klass()
    xys = [(r[num.at], add(rhs, Y(r))) for r in rows if r[num.at] != "?"]
    xpect = div(rhs)
    for x, y in sorted(xys, key=lambda xy: x[0]):
      if x != b4:
        if the.leaf <= lhs.n <= len(xys) - the.leaf:
          tmp = (lhs.n * div(lhs) + rhs.n * div(rhs)) / len(xys)
          if tmp < xpect:
            xpect, out = tmp, [("<=", num.at, b4), (">", num.at, b4)]
        add(lhs, sub(rhs,y))
        b4 = x
    if out:
      return o(div=xpect, hows=out)

  return (_sym if col.it is Sym else _num)(col)

def tree(data, Klass=Num, Y=None, how=None):
  Y       = Y or (lambda row: ydist(data,row))
  data.kids = []
  data.how  = how
  data.ys   = Num(Y(row) for row in data.rows)
  if data.n >= the.leaf:
    tmp = [x for c in data.cols.x if (x := cuts(c,data.rows,Y,Klass=Klass))]
    if tmp:
      for how1 in sorted(tmp, key=lambda cut: cut.div)[0].hows:
        rows1 = [row for row in data.rows if selects(row, *how1)]
        if the.leaf <= len(rows1) < data.n:
          data.kids += [tree(clone(data,rows1), Klass, Y, how1)]
  return data

def nodes(data1, lvl=0, key=None):
  yield lvl, data1
  for data2 in (sorted(data1.kids, key=key) if key else data1.kids):
    yield from nodes(data2, lvl + 1, key=key)

def leaf(data1,row):
  for data2 in data1.kids or []:
    if selects(row, *data2.decision):
      return leaf(data2, row)
  return data1

def show(data, key=lambda z:z.ys.mu):
  stats = data.ys
  win = lambda x: 100-int(100*(x-stats.lo)/(stats.mu - stats.lo))
  print(f"{'d2h':>4} {'win':>4} {'n':>4} ")
  print(f"{'----':>4} {'----':>4} {'----':>4} ")
  for lvl, node in nodes(data, key=key):
    leafp = len(node.kids)==0
    post = ";" if leafp else ""
    xplain = ""
    if lvl > 0:
      op,at,y = node.decision
      xplain = f"{data.cols.all[at].txt} {op} {y}"
    print(f"{node.ys.mu:4.2f} {win(node.ys.mu):4} {node.n:4}  {(lvl-1)*'| '}{xplain}" + post)
```

```
494 ### Utils -----------------------------------------------------------------
495
496 def cat(v):
497   it = type(v)
498   inf = float('inf')
499   if it is list:  return "{" + ",".join(map(cat, v)) + "}"
500   if it is float: return str(int(v)) if -inf < v < inf and v == int(v) else f"{v
    :.3g}"
501   if it is dict:  return cat([f":{k} {cat(w)}" for k, w in v.items()])
502   if it in [type(abs), type(cat)]: return v.__name__
503   return str(v)
504
505 def report(rows, head, decs=2):
506   w=[0] * len(head)
507   Str  = lambda x   : f"{x:.{decs}f}"      if type(x) is float else str(x)
508   say  = lambda w,x : f"{x:>{w}.{decs}f}" if type(x) is float else f"{x:>{w}}"
509   says = lambda row : ' | '.join([say(w1, x) for w1, x in zip(w, row)])
510   for row in [head]+rows:
511     w = [max(b4, len(Str(x))) for b4,x in zip(w,row)]
512   print(says(head))
513   print(' | '.join('-'*(w1) for w1 in w))
514   for row in rows: print(says(row))
515
516 def shuffle(a):
517   random.shuffle(a)
518   return a
519
520
```

```
520 ### Start-up ---------------------------------------------------------------
521
522 def main():
523   cli(the.__dict__)
524   for s in sys.argv:
525     if fun := globals().get("eg" + s.replace("-", "_")):
526       random.seed(the.rseed)
527       fun()
528
529 if __name__ == "__main__": main()
530
```