

SLAC Computation Research Group  
Stanford, California

CGTM No. 156  
April 1974

**MASTER COPY  
DO NOT REMOVE**

CHAN -- Computer Program for  
One-Dimensional Projection Pursuit

J. H. Friedman

Abstract

A computer program that implements the Projection Pursuit algorithm of Friedman and Tukey is described. This routine seeks to find one-dimensional linear projections of multivariate data that are relatively highly revealing.

## INTRODUCTION

CHAN is a collection of subroutines for finding those linear one-dimensional projections of multivariate data that reveal interesting structure. The algorithm and its properties are detailed elsewhere<sup>1</sup> and are only discussed briefly here as they relate to its implementation.

The Projection Pursuit algorithm assigns to each projection direction,  $\hat{k}$ , in the multidimensional space, a numerical index,  $I(\hat{k})$ , that corresponds to the degree of structuring present in the data as projected onto  $\hat{k}$ . The more structure present in the projection, the larger  $I(\hat{k})$  becomes. The essence of the algorithm is to find those projection directions,  $\hat{k}$ , that maximize  $I(\hat{k})$ . For complex data structures, several solutions may exist, and for each of these, the solution projections can be visually inspected by the researcher for interpretation and judgment as to their usefulness. This multiplicity is often important.

The input to the algorithm consists, principally, of the multivariate data point set and a starting direction,  $\hat{k}_s$ , in the multidimensional data space. Starting at  $\hat{k}_s$ , the algorithm finds the direction  $\hat{k}^*$  that represents the first local maximum of  $I(\hat{k})$  uphill from  $\hat{k}_s$ . Useful starting directions include the larger principal axes of the data set, the original coordinate axes, and even directions chosen at random. From these searches, several quite distinct solutions may result. Each of these projections can then be examined to determine their usefulness in data interpretation.

## THE PROJECTION INDEX

The intent of the projection index is to assign to each direction,  $\hat{k}$ , in the multivariate space, a numerical value that closely corresponds to the degree of data structuring when projected onto  $\hat{k}$ . It was found that those projections that were most interesting to researchers were those that tended to

produce many very small interpoint distances while, at the same time, maintaining the overall spread of the data in the projection. Such projections will, for instance, tend to concentrate the points into clusters while, at the same time, separating the clusters.

The spread of the data as projected onto  $\hat{k}$  is estimated by taking the trimmed standard deviation of the data from the mean:

$$s(\hat{k}) = \left[ \sum_{i=pN}^{(1-p)N} (\vec{X}_i \cdot \hat{k} - \bar{X}_k)^2 / (1-2p)N \right]^{1/2} \quad (1)$$

where

$$\bar{X}_k = \sum_{i=pN}^{(1-p)N} \vec{X}_i \cdot \hat{k} / (1-2p)N$$

Here  $N$  is the total number of data points, and  $\vec{X}_i$  ( $i=1, N$ ) are the multivariate vectors representing each of the data points, ordered according to their projections  $\vec{X}_i \cdot \hat{k}$ . A small fraction,  $p$ , of the points that lie at each of the extremes of the projections are omitted from both sums. Thus, extreme values of  $\vec{X}_i \cdot \hat{k}$  do not contribute to  $s(\hat{k})$ , which is thus robust against extreme outliers.

For a "small interpoint distance" measure, an average nearness function is used which has the form

$$d(\hat{k}) = \sum_{i=1}^N \sum_{j=1}^N (R-r_{ij}) l(R-r_{ij}) \quad (2)$$

where

$$r_{ij} = | \vec{X}_i \cdot \hat{k} - \vec{X}_j \cdot \hat{k} |$$

and  $l(\eta)$  is unity for positive valued arguments and zero for negative values.

Thus, the double sum is confined to pairs with  $0 \leq r_{ij} \leq R$ .

CHAN offers two choices for projection indices. One of these simply takes the product

$$I(\hat{k}) = d(\hat{k}) s(\hat{k}) \quad (3a)$$

keeping the local cutoff radius,  $R$ , in  $d(\hat{k})$  (Eqn. 2) at a constant value,  $R_0$ . (This is the index described in Ref. 1). The other index uses

$$I(\hat{k}) = d(\hat{k}) \quad (3b)$$

but where  $R = f_0 s(\hat{k})$ .

Here  $f_0$  is a constant fraction. The first index (Eqn. 3a) has the advantage that the cutoff radius,  $R$ , (which establishes the distance in the projected subspace over which the local density is averaged and thus establishes the scale of density variation to which the algorithm is sensitive) is a constant completely under user control. The advantage of the second index (Eqn. 3b) is that it is Affine invariant. That is, it is unchanged by scale transformations on the input data. Both projection indices have comparable performance but they do not always find the same solutions.

#### CONSTRAINT DIRECTIONS

In order to encourage the algorithm to find as many distinct solutions as possible, it is useful to be able to reduce the dimensionality of the space to be searched. This can be done by choosing an arbitrary set of directions,  $\{\hat{\sigma}_i\}_{i=1}^m$ ,  $m < n$ , which need not be mutually orthogonal, and applying the constraints

$$\hat{k}^* \cdot \hat{\sigma}_i = 0 \quad i = 1, m \quad (4)$$

on the solution direction  $\hat{k}^*$ . Possible choices for constraint directions might be solution directions found on previous searches, or directions that are known in advance to contain considerable, but well understood, structure. Also, when the choice of scales for the several coordinates is guided by considerations outside the data, one might wish to remove directions with

small variance about the mean since these directions often provide little information about the structure of the data. The introduction of each such constraint direction reduces by one the number of search variables, and thus increases the computational efficiency of the algorithm. CHAN allows for the introduction of an arbitrary number,  $m < \text{NDIM}$  of non-parallel constraint directions.

### IMPLEMENTATION

One-dimensional projection pursuit is invoked from a FORTRAN program by a subroutine call

CALL CHAN(M,RADIUS,NDIM,NPTS,X,KEY,NF,NL,XB,KEYSCL,SCALE,D).

All quantities are input to the routine except XB and SCALE which serve a dual purpose as both input and output arrays. These quantities have the following meanings:

The quantities M and RADIUS refer to the projection index,  $I(\hat{k})$ .

$$M = \begin{cases} 1 & \text{use Affine invariant projection index (Eqn. 3b)} \\ 2 & \text{use constant radius projection index (Eqn. 3a)} \end{cases}$$

$$RADIUS = \begin{cases} \text{for } M=1: & \text{constant fraction, } f_0, \text{ (Eqn. 3b)} \\ \text{for } M=2: & \text{constant local cutoff radius, } R_0, \text{ (Eqn. 3a)} \end{cases}$$

RADIUS = 'DFLT':  $f_0$  and/or R are automatically calculated by the program.

$$f_0 = \text{RADSCAL} / \sqrt{NPTS} \quad (NPTS \leq NPTS0)$$

$$f_0 = \text{RADSCAL} * (\sqrt{NPTS0} / \log(NPTS0)) * (\log(NPTS) / (NPTS)) \quad (NPTS > NPTS0)$$

$$R = f_0 * s(\hat{e}_1)$$

In CHAN, RADSCAL is set to 1.0 and NPTS0 is set to 1000. (See below on how to change these values.) The quantity  $s(\hat{e}_1)$  is the standard deviation about the mean of the data as projected onto the principal axis with the largest eigenvalue. (The trimming factor,  $p$ , (Eqn. 1) is zero for this calculation only; i.e.,  $s(\hat{e}_1)$  is the square root of the largest eigenvalue of the total

sample covariance matrix.) The values for  $f_0$  and R are calculated so that the average number of points within the local cutoff radius grows as the square root of the sample size for  $NPTS \leq NPTS0$ , and grows as the logarithm for  $NPTS > NPTS0$ .

The quantities NDIM, NPTS and X refer to the input data.

NDIM = number of attributes per data point (dimensionality of data space)

NPTS = number of data points (sample size)

X = array, dimensioned X(NDIM,NPTS) in the calling program, that contains the coordinates of the data points.

The quantities KEY,NF,NL and XB refer to the starting directions,  $\hat{k}_s$ , for the search.

KEY = axis flag = {  
  'ORIG': use original axes for starting directions  
  'EIGN': use principal axes as starting directions  
          (ordered in decreasing value of corresponding eigenvalue)  
  'USER': use user supplied vectors as starting directions

NF = axis number of starting direction for first search

NL = axis number of starting direction for last search

There will be  $NL - NF + 1$  different searches whose starting directions will either be the original ('KEY = 'ORIG'), principal (KEY = 'EIGN'), or user supplied (KEY = 'USER') axes NF, NF + 1, ... NL. If  $NF > NL$ , there will be no searches, only a principal axis analysis (see below under output).

XB = array, dimensioned XB(NDIM,NL) in calling program, that contains the user supplied starting directions for KEY = 'USER'. For KEY = 'ORIG' or 'EIGN', XB is ignored on input. Upon return from CHAN to the calling program, XB contains the  $NL - NF + 1$  solution directions,  $\hat{k}^*$ , from the searches starting in XB(1,NF).

The quantities KEYSCL and SCALE refer to the scaling of the input data along the original axes before the projection pursuit. (For the Affine

invariant projection index ( $M = 1$ ), the results are independent of the scaling of the input data.)

KEYSCL = scaling flag =  $\left\{ \begin{array}{l} \text{'NONE'} \text{ do not scale data} \\ \text{'USER'} \text{ use user provided scales} \\ \text{'ORIG'} \text{ employ automatic scaling} \\ \text{'SMAX'} \text{ automatic scaling with user provided} \\ \text{lower limits for scales} \end{array} \right.$

SCALE = array, dimensioned SCALE(NDIM). For KEYSCL = 'NONE' or 'ORIG' this is a scratch array for CHAN. For KEYSCL = 'USER', it contains the user provided scales and for KEYSCL = 'SMAX', it contains the user provided lower limits for the axis scales.

The automatically calculated scales for each original axis are the standard deviations of the data about their means, as projected onto each axis. Under the KEYSCL = 'SMAX' option, the axis scale is set to the maximum of the calculated scale and the user provided minimum scale for that axis. (The standard deviation calculation for the scales is trimmed by the global trimming factor IPER, described below. IPER is set to zero in CHAN but may be changed by the user.) Upon return from CHAN to the calling program, SCALE contains the scales used for each axis.

D = Real\*8 scratch array for CHAN, dimensioned  
Real\*8 D(NDIM,NDIM,3).

#### Labeled Common Scratch Arrays

Three labeled common scratch arrays must be declared in the calling program as working storage for CHAN. These are listed here along with their various dimensions:

```
COMMON // V(240) /POINTS/ Z(NPTS+2) /DATA/P(NDIM*NPTS+1)
```

### Constraint Directions

As discussed above, CHAN allows for the introduction of an arbitrary member (<NDIM) of constraint directions. The solution direction,  $\hat{k}^*$ , will be constrained to be simultaneously perpendicular to each of these directions. The constraint directions need not themselves be mutually orthogonal, but no two may be exactly parallel (to machine accuracy). The constraint directions are entered into a labeled common whose label and dimensions are given below:

```
COMMON / NORMAL / NORM / NORMV / VEC(NDIM,NORM)
```

NORM = number of user supplied constraint directions

VEC = coordinates of constraint vectors

If not referenced by the user, the default value for NORM is zero.

### Alternate CALLS

CHAN may be called repeatedly from the calling program, each call initiating NL-NF+1 projection pursuits. If subsequent calls do not change the input data or its scaling, an alternate shortened calling sequence may be used

```
CALL CHARLY(M,RADIUS,KEY,NF,NL,XB)
```

The arguments in the calling sequence have the same meanings as described above. This call may not be the first call to CHAN. All quantities that are absent from this shortened calling sequence have the values assigned to them on the last call to CHAN, using the full calling sequence.

### Random Direction Generator

To facilitate starting the projection pursuit algorithm at random directions in the multidimensional space, CHAN provides a random direction generator

```
CALL RANDIR(XB,NDIM,NDIR)
```



Upon return to the calling program, XB [dimensioned XB(NDIM,NDIR) in the calling program] contains the coordinates of NDIR random directions in the NDIM-dimensional space. Repeated calls to RANDIR generates new sets of random directions for each call.

### Optional User Control

There are several parameters internal to CHAN that the user may change at his discretion. Some of these parameters control the verbosity of the output. Others are parameters of the projection pursuit algorithm to which it is reasonably insensitive, and thus, they need not be changed frequently. These parameters are stored in the labeled common blocks listed below, along with their default values (set in a BLOCK DATA sub-program internal to CHAN).

```
COMMON /TERSE/ ITERSE
                (0)
```

ITERSE controls the number of pictorial histograms displayed in the output. For ITERSE=1, a pictorial histogram is displayed of the data as projected onto each solution direction,  $\hat{k}^*$ . If KEY = 'ORIG', then, in addition, the data as projected onto each user supplied starting direction,  $\hat{k}_s$ , is also displayed. For ITERSE=0, twice NDIM additional histograms are displayed of the data as projected onto each of the original and principal axes.

```
COMMON /COMST/ IPER,PEROUT
                (0) (.01)
```

IPER is a global trimming parameter for the evaluation of both  $d(\hat{k})$  and  $s(\hat{k})$  (Eqn's 1 and 2). That is, IPER projected points are deleted from each extreme of the projection before both  $d(\hat{k})$  and  $s(\hat{k})$  are evaluated for the projection. PEROUT is the parameter p in Eqn. 1 and is a fractional trimming factor for the evaluation of  $s(\hat{k})$  only.

```
COMMON /EIT/ EITMIN,RADSCL,NEVO
                (.01) (1.0) (1000)
```

If the variance of the data as projected onto a principal axis is sufficiently small (as compared to the largest eigenvalue of the covariance matrix), then CHAN will automatically consider that axis to be a constraint direction. EITMIN is the minimum value of the square root of the ratio of the axis variance to the largest eigenvalue, for the principal axis not to be considered a constraint direction. The parameters RADSCAL and NEVO are discussed above under the default option for the local cutoff radius.

```
COMMON /CNTRL/ NBINS,CONV1,CONV2,NSTPS,NPRNT
              (120) (.05) (-.01) (2) (1)
```

NBINS is the number of bins (channels) in the pictorial histogram displays ( $\leq 120$ ). CONV1 and CONV2 control the convergence criteria of the maximization algorithm used to maximize  $I(\hat{k})$ . The search converges when the change in every parameter is less than  $.1 * CONV1$ , or when  $|[I(\hat{k}) - I(\hat{k}')] / I(\hat{k})| \leq -CONV2$ , where  $I(\hat{k})$  is the projection index value at the current iteration and  $I(\hat{k}')$  is the value at the previous iteration. NSTPS pertains to the maximum number of iterations allowed for the maximizer. This maximum number of iterations is NSTPS, plus the number of search variables. NPRNT controls the verbosity of the printed output from the maximizer.

NPRNT	{	$\leq 0$ : no maximizer printed output
		$> 0$ : printed output for every iteration number that is an integral multiple of NPRNT

#### PRINTED OUTPUT

As indicated above, the nature of the printed output from CHAN depends upon the options chosen by the user. Figure 2 illustrates the printed output from the FORTRAN program listed in Figure 1.

Listed first (SCALE = ) are the scale factors that were used to scale each of the NDIM original axes. For the ITERSE=0 option, this is followed by the 2\*NDIM pictorial histograms (one per page) of the data as projected onto each of the original and principal axes of the data. For ITERSE=1 (the option used in this example), these histograms are skipped. Shown next

(EIGENSTDS = ) are the square roots of the eigenvalues of the total sample (no trimming) covariance matrix, listed in decreasing order. The following line lists the value for  $f_0$  of Eqn. 3b (RADIUS SCALE = ), R of Eqn. 2 (LOCAL RADIUS = ), and EITMIN (see optional user control) times the first (largest) eigen square root value (CUTOFF EIGEN STD = ). This latter value is the smallest allowable eigenvalue square root for the corresponding principal axis not to be considered to be a constraint direction. For KEY = 'USER' (as in this example), a pictorial histogram of the data as projected onto the user supplied starting direction, is shown. For KEY = 'ORIG' or 'EIGN', this histogram does not appear. The printed output from the maximizer follows next (if not suppressed by the user). For each iteration in the search for the maximum of the projection index,  $I(\hat{k})$ , the iteration number, the cumulative number of evaluations of  $I(\hat{k})$ , the value of  $-I(\hat{k})$  (F = ), and the values of the search parameters (X= ), are listed. Iteration zero corresponds to the starting direction. The parameters of the search are those of the solid angle transform,  $SAT(\hat{k})$ , of the corresponding direction  $\hat{k}$  in the NDIM-dimensional space (see Ref. 1, pg. 7 and Appendix), and have no direct interpretation.

After a solution,  $\hat{k}^*$ , is found, its components, expressed in terms of both the original coordinates and the principal axes, are listed. This is followed by a pictorial histogram of the data as projected onto this solution direction.

For each call to CHAN, the output sequence is repeated. If, however, subsequent calls are made to CHARLY using the shortened calling sequence, then the purely data dependent information is not repeated.

REFERENCE

1. J. H. Friedman and J. W. Tukey, "A Projection Pursuit Algorithm for Exploratory Data Analysis," Stanford Linear Accelerator Center, Stanford, California, Report, SLAC PUB-1312, September 1973. (To be published in IEEE Trans. Computers.)

```

C ONE-DIMENSIONAL PROJECTION PURSUIT EXAMPLE. 900 POINTS, 8 DIMENSIONS.
C
C
COMMON /POINTS/ MM(2),Z(900) /DATA/ NM,P(8,900) /TERSE/ ITERSE
REAL SCALE(8),X(8,900),XB(8)
REAL*8 D(8,8,3)
C
C          READ IN DATA
CALL DATAPT(X,900,8)
C
C          GENERATE RANDOM STARTING DIRECTION
CALL RANDIR(XB,8,1)
C
C          SET FLAG FOR TERSE OUTPUT
ITERSE=1
C
C          INITIATE PROJECTION PURSUIT
CALL CHAN(1,'DFLT',8,900,X,'USER',1,1,XB,'NONE',SCALE,D)
C
STOP
END

```

FIGURE 1



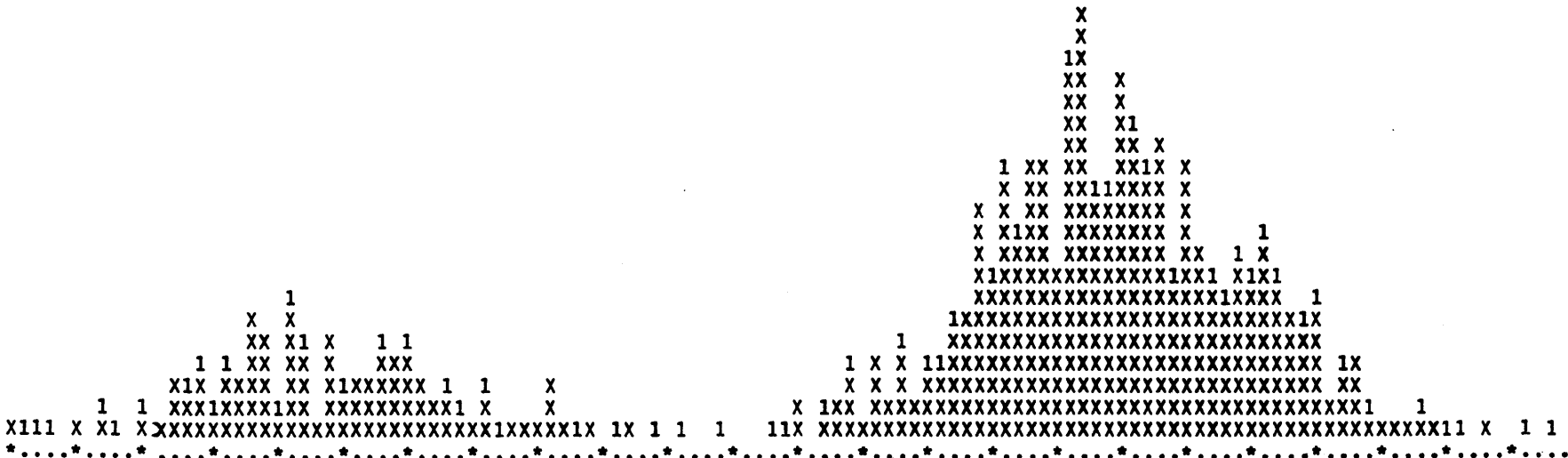
FIGURE 2b

ITER 54 CALLS F=-0.1063106E 04 X= 0.4671E-01 0.4030E 0.9807E 00 0.6242E 00 0.9355E 00 0.3486E 00 0.1316E  
 ITER 81 CALLS F=-0.1079041E 04 X= 0.4780E-01 0.4137E 00 0.9782E 00 0.6236E 00 0.9325E 00 0.3475E 00 0.1299E 00  
 ITER 107 CALLS F=-0.1081021E 04 X= 0.4780E-01 0.4183E 00 0.9779E 00 0.6236E 00 0.9339E 00 0.3474E 00 0.1308E 00

---CHAN HISTOGRAM--- MAXIMUM PROJECTION AXIS-P( 1) =

0.92981E-02 -0.57118E-02 0.42183E-03 -0.75109E 00 0.92701E-01 0.16599E 00 -0.61447E 00 0.14841E 00  
 EIV DECOMP = 0.20489E 00 0.66160E 00 0.51166E-01 -0.90716E-01 -0.91601E-01 -0.93173E-01 0.57440E 00 -0.40305E 00

80.  
78.  
76.  
74.  
72.  
70.  
68.  
66.  
64.  
62.  
60.  
58.  
56.  
54.  
52.  
50.  
48.  
46.  
44.  
42.  
40.  
38.  
36.  
34.  
32.  
30.  
28.  
26.  
24.  
22.  
20.  
18.  
16.  
14.  
12.  
10.  
8.  
6.  
4.  
2.



-----  
 LOWER BIN 8777777777 76666666666655555555554444444444443333333333222222222211111111111000000000000000001111111111222222222233  
 EDGE 09987654321 098765433210987654321098765432109876543210987654321098765432100123445678901234567890112345678901  
 99011223344 56677889001122334455667788900112234455667788900112234455667788900112234455667788900112334455662110099887665544322110099877665544  
 39506173849 5162739405162839406172849506273840516283940617284950617384951627394051728395938271504937261594837160593827160  
 9516284062840628406284951739517395173951728406284062840628406173951739517395173940628404826048260483715937159371504  
 \* \* \* \* \*  
 11 1 1 11212122134223222121111111111

CONTENTS 21110203103 2657376203392056698964532512226212012010100100011403472849477122559666503349585685375952134783222321102001010