

Title here

Author 1 *

Department of YYY, University of XXX
and

Author 2

Department of ZZZ, University of WWW

July 13, 2020

Abstract

Friedman & Tukey commented on their initial paper on projection pursuit in 1974 that “the technique used for maximising the projection index strongly influences both the statistical and the computational aspects of the procedure.” While many projection pursuit indices have been proposed in the literature, few concerns the optimisation procedure. In this paper, we developed a system of diagnostics aiming to visually learn how the optimisation procedures find its way towards the optimum. This diagnostic system can be applied to more general to help practitioners to unveil the black-box in randomised iterative (optimisation) algorithms. An R package, *ferri*, has been created to implement this diagnostic system.

Keywords: optimisation, projection pursuit, guided tour, visual, diagnostics, R

*The authors gratefully acknowledge ...

1 Introduction

In an optimization problem the goal is to find the best solution within the space of all feasible solutions which typically is represented by a set of constraints. The problem consists on optimizing an objective function $f : S \rightarrow \mathbb{R}$ with $S \in \mathbb{R}^n$ in a reduced spaced given by the problem constraints to either minimize or maximize a function. . . . Gradient based optimization has been typically used to solve such problems. However, there are many situations where derivatives of an objective function are unavailable or unreliable and therefore traditional methods based on derivatives are not the best option to solve an optimization problem.

Derivative free methods provides another option to optimise the objective function without evaluating any gradient or Hessian information and a particular class of methods: direct search, has gained its popularity through its conceptual simplicity. However, the whole searching process in the algorithm remains a black-box. Plots are usually used to evaluate and compare the performance of different algorithms but it can easily become tedious because the code will have to be modified significantly when comparing different parameters in the algorithms. For example, a categorical variable with 5 levels can be easily mapped onto color while mapping another categorical variable with 30 levels will not make the plot informative. Thus the plot needs to be re-designed to better suits the characteristics of the parameter (whether the parameter is a scalar or a matrix? whether the parameter is quantitative or categorical? If categorical, how many levels does the parameter have?). This motivates the design of a visual diagnostic framework for optimisation algorithms based on the idea of a *global object*.

The paper is organised as follows. Section 2 gives a general literature review of optimisation, specifically derivative free optimisation. Section 4 presents the new idea of constructing a systematic visual framework that diagnoses the components of an optimisation procedure (parameters, searching path, etc). The rest of the paper serves as a comprehensive example of using the visual diagnostics on one particular problem: *projection pursuit guided tour*. Some background knowledge of projection pursuit guided tour is provided in Section 3. Section ?? applies the concepts proposed in section 4 in the tour problem and sets up the data. The last section, Section ??, presents the visual diagnostic

plots and explains how they can help to understand different aspects of the optimisation in tour.

2 Derivative free optimisation

Given an objective function f , one way of optimising it is to equate its gradient to zero. In modern optimisation problems, gradient information can be hard to evaluate or sometimes even impossible and Derivative-Free Optimisation (DFO) methods can be useful to approach these problems. One common class of methods in DFO is *Direct-search methods*. Coined by Hooke & Jeeves (1961), direct search methods don't require any gradient or Hessian information and has gained its popularity through its simplicity in use and reliability in solving complicated practical problems. Depends on whether a random sample is used in the search, this class of methods can be further classified as *stochastic* or *deterministic*. The stochastic version of the direct-search methods will be the main optimisation procedure analysed in this paper.

[How about adding more details into derivative free methods? ppp]

2.1 Difficulties

[Are we using projection pursuit/guided tour to better understand the convergence of optimization algorithms visually in combination with the algorithms discussed below? Or we are focusing on the optimisation problem only within the project pursuit context? Some of the problems listed below are also applicable to optimization problem in general too. ppp]

Below listed several issues in projection pursuit optimisation. Some are general optimisation problems, while others are more specific for PP optimisation.

- *Finding global maximum*: Although finding local maximum is relatively easy with developed algorithms, it is generally hard to guarantee global maximum in a problem where the objective function is complex or the number of decision variables is large. Also, there are discussions on how to avoid getting trapped in a local optimal in the literature.

- *optimising non-smooth function*: When the objective function is non-differentiable, derivative information can not be obtained, which means traditional gradient- or Hessian- based methods are not feasible. Stochastic optimisation method could be an alternative to solve these problems.
- *computation speed*: The optimisation procedure needs to be fast to compute since tours produces real-time animation of the projected data.
- *consistency result in stochastic optimisation*: In stochastic algorithm, researchers usually set a seed to ensure the algorithm producing the same result for every run. This practice supports reproducibility, while less efforts has been made to guarantee different seeds will provide the same result.
- *high-dimensional decision variable*: In projection pursuit, the decision variable includes all the entries in the projection matrix, which is high-dimensional. Researcher would be better off if they can understand the relative position of different projection matrix in the high-dimensional space.
- *role of interpolation in PP optimisation*: An optimisation procedure usually involves iteratively finding projection bases that maximises the index function, while tour requires geodesic interpolation between these bases to produce a continuous view for the users. It would be interesting to see if the interpolated bases could, in reverse, help the optimisation reach faster convergence.

Think about how does your package help people to understand optimisation

- diagnostic on stochastic optim
- vis the progression of multi-parameter decision variable
- understanding learning rate - neighbourhood parameter
- understand where the local & global maximum is found - trace plot - see if noisy function

3 Projection pursuit guided tour

From Section 3, we presents a comprehensive case study on how to use the visual diagnostics to explore the optimisation in a specific problem: projection pursuit guided tour. Section 3 aims to provide non-experts with an overview of the problem content and the existing optimisation procedures used in projection pursuit guided tour. For those who are already familiar with the techniques, feel free to skip this section.

The optimisation problem we’re interested in is in the context of projection pursuit. Coined by Friedman & Tukey (1974), projection pursuit is a method that detects the interesting structure (i.e. clustering, outliers and skewness) of multivariate data via projecting it in lower dimensions. Let $\mathbf{X}_{n \times p}$ be a data matrix, an n-d projection can be seen as a linear transformation: $T : \mathbb{R}^p \mapsto \mathbb{R}^d$ defined by $\mathbf{P} = \mathbf{X} \cdot \mathbf{A}$, where $\mathbf{P}_{n \times d}$ is the projected data and $\mathbf{A}_{p \times d}$ is the projection basis. Define $f : \mathbb{R}^{p \times d} \mapsto \mathbb{R}$ be an index function that maps the projection basis \mathbf{A} onto a scalar value I , this function is commonly known as the projection pursuit index (PPI) function, or index function and can be used to measure the “interestingness” of the projection. A number of indice functions have been proposed in the literature and this includes Legendre index (Friedman & Tukey 1974), hermite index (Hall et al. 1989), natural hermite index (Cook et al. 1993), chi-square index (Posse 1995), LDA index (Lee et al. 2005) and PDA index (Lee & Cook 2010). [Any literature on holes index? xxx]

As Friedman & Tukey (1974) noted “..., the technique use for maximising the projection index strongly influences both the statistical and the computational aspects of the procedure.” A suitable optimisation procedure is needed to find the projection angle that maximises the PPI and the quality of the optimisation largely affect the interesting projections one could possibly observe.

Projection pursuit is usually used in conjunction with a tour method called *guided tour*. Tour explores the multivariate data *interactively* via playing a series of projections, that form a *tour path* and guided tour uses the path that is geodesically the shortest. Details of the mathematical construction of a tour path can be found in Buja et al. (2005). Modified from Buja et al. (2005), Figure 1 vividly depicts the tour path in guided tour. Guided tour, along with other types of tour, has been implemented in the *tourrr* package in R,

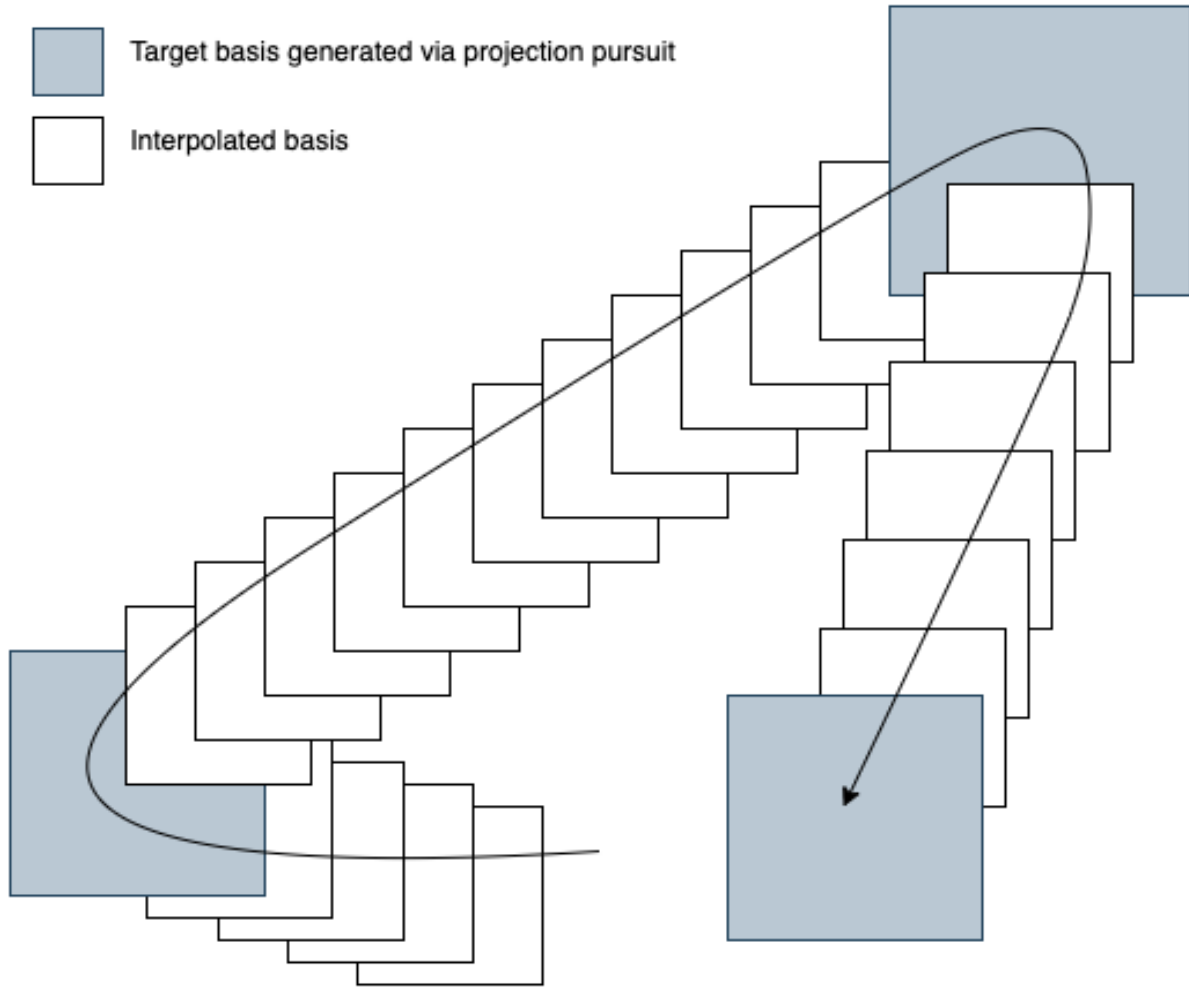


Figure 1: An illustration of the tour path

available on the Comprehensive R Archive Network at <https://cran.r-project.org/web/packages/tourr/> (Wickham et al. 2011).

[doubt if I need the following. xxx - maybe an algorithmic version of guide tour would be helpful for non-expert to understand the method?]

3.1 Optimisation in projection pursuit

[thinking if the following three should be presented as algorithms or plain description is fine. xxx]

Now we formulate the problem in mathematical notation. Projection pursuit finds the

projection basis $\mathbf{A}^J = [\mathbf{a}_1, \dots, \mathbf{a}_d]$, where \mathbf{a}_i is a $p \times 1$ vector, satisfying the following optimisation problem:

$$\arg \max_{\mathbf{A} \in \mathcal{A}} f(\mathbf{A}) \quad (1)$$

$$s.t. \mathcal{A} = \{\forall \mathbf{a}_i, \mathbf{a}_j \in \mathbf{A} : \|\mathbf{a}_i\| = 1 \text{ and } \mathbf{a}_i \mathbf{a}_j = 0\} \quad (2)$$

via an iterative search: $\{\mathbf{A}^1, \mathbf{A}^2, \dots, \mathbf{A}^J\}$. Guided tour interpolates geodesically between each basis found by projection pursuit \mathbf{A}^j to form a tour path:

$$\{\mathbf{A}^0, \mathbf{A}^1, {}_1\mathbf{A}^1, {}_2\mathbf{A}^1, \dots, {}_{k_1}\mathbf{A}^1, \mathbf{A}^2, {}_1\mathbf{A}^2, {}_2\mathbf{A}^2, \dots, {}_{k_2}\mathbf{A}^2, \dots, \mathbf{A}^{J-1}, {}_1\mathbf{A}^{J-1}, \dots, {}_{k_m}\mathbf{A}^{J-1}, \mathbf{A}^J\}$$

Where \mathbf{A}^0 is the starting basis and \mathbf{A}^J is the final basis that has the largest index value in the searching space. Here we use left subscript to denote the index of the interpolation bases in each iteration i.e. ${}_2\mathbf{A}^1$ means the second interpolation basis in the first iteration. The right subscript is reserved for the index of searching points in the optimisation algorithm.

There are three existing methods for optimising PPI function and we review them below. Posse (1995) proposed a stochastic direct search method, a random search algorithm. Given a current basis \mathbf{A}^j , a candidate basis \mathbf{A}_l^{j+1} is sampled in the neighbourhood defined by the radius of the p -dimensional sphere, α of the current basis \mathbf{A}^j by

$$\mathbf{A}_l^{j+1} = \mathbf{A}^j + \alpha \mathbf{A}_{rand}$$

. If the candidate basis has a higher index value than the current basis, it is outputted as the target basis \mathbf{A}^{j+1} along with other metadata and the current iteration stops. If no basis is found to have higher index value after the maximum number of tries l_{\max} , the algorithm stops with nothing outputted. c is the halving parameter and when $c > 30$, the neighbourhood parameter α in the next iteration will be reduced by half. The algorithm for a random search is summarised in Algorithm 1.

Algorithm 1: random search: my first algorithm

```
1  initialisation:  $c = 0, l = 1, l_{\max}, \mathbf{A}^j$ ;  
2  while  $l < l_{\max}$  do  
3       $\mathbf{A}_l^{j+1} = \mathbf{A}^j + \alpha \mathbf{A}_{rand}$ ;  
4       $I_l^{j+1} = f(\mathbf{A}_l^{j+1})$ ;  
5      if  $I_l^{j+1} > I^j$  then  
6           $\mathbf{A}^{j+1} = \mathbf{A}_l^{j+1}$ ;  
7      else  
8           $c = c + 1$ ;  
9      end  
10      $l = l + 1$ ;  
11 end
```

Cook et al. (1995) explained the use of a gradient ascent optimisation with the assumption that the index function is continuous and differentiable. Since some indices could be non-differentiable, the computation of derivative is replaced by a pseudo-derivative of evaluating five randomly generated directions in a tiny nearby neighbourhood. Taking a step on the straight derivative direction has been modified to maximise the projection pursuit index along the geodesic direction.

Algorithm 2: searchgeodesic: this part has problem!

```

1  initialisation:  $l_{\max}, \alpha, l_{\max}, \mathbf{A}^j$ ;
2  while  $l < l_{\max}$  do
3      Generate  $\{A^{(j_1)} : A^{(1)}, \dots, A^{(5)}\}$ ;
4      Find  $A_l^{J_1}$  where  $J_1 \in j_1$  such that  $\{A^{(J_1)} : \arg \max I(A^{(j_1)})\}$ ;
5      Denote  $F_a = A^j$  and  $F_z = A^{(J_1)}$ ;
6      Conduct singular value decomposition (SVD) for  $F_a$  and  $F_z$  such that
           $F_a^T F_z = V_a \Lambda V_z^T$ ;
7      Compute  $G_a = F_a V_a$ ,  $G_z = F_z V_z$  and  $\tau = \arccos(\Lambda)$ ;
8      ...;
9      Find  $A^{(J_2)}$  satisfying  $A^{J_2}$  on the interpolation path such that
           $\{A^{(J_2)} : \arg \max I(A^{(j_2)})\}$ ;
10     Compute  $I^{(J_2)} = f(A^{(J_2)})$ ,  $p_{diff} = (I^{(J_2)} - I^j)/I^j$ ;
11     if  $p_{diff} > 0.001$  then
12          $A^{j+1} = A^{(J_2)}$ ;
13          $I^{j+1} = I^{(J_2)}$ ;
14         Output:  $A^{j+1}, I^{j+1}, i$  and relevant metadata
15     end
16      $i = i + 1$ ;
17 end

```

Simulated annealing (Bertsimas et al. 1993, Kirkpatrick et al. (1983)) is a non-derivative procedure based on a non-increasing cooling scheme $T(l)$. Given an initial T_0 , the temperature at iteration l is defined as $T(l) = \frac{T_0}{\log(l+1)}$. The simulated annealing algorithm works as follows. Given a neighbourhood parameter α and a randomly generated orthonormal basis A_{rand} , a candidate basis is constructed as $A_l^{j+1} = (1 - \alpha)A^j + \alpha A_{rand}$. If the index value of the candidate basis is larger than the one of the current basis, the candidate basis becomes the target basis. If it is smaller, the candidate is accepted with probability $\text{prob} = \min \left\{ \exp \left[\frac{I^j - I_l^{j+1}}{T(l)} \right], 1 \right\}$ where I^j and I_l^{j+1} are the index value of the current and candidate basis respectively. The algorithm can be summarised as in Algorithm 3.

Algorithm 3: simulated annealing

```
1  initialisation:  $T_0, \alpha, l_{\max}, \mathbf{A}^j$ ;  
2  while  $l < l_{\max}$  do  
3       $\mathbf{A}_l^{j+1} = (1 - \alpha)\mathbf{A}^j + \alpha\mathbf{A}_{rand}$  s.t.  $\|\mathbf{a}_i\| = 1$ ;  
4       $I_l^{j+1} = f(\mathbf{A}_l^{j+1})$ ;  
5       $T(l) = \frac{T_0}{\log(l+1)}$ ;  
6      if  $I_l^{j+1} > I^j$  then  
7           $\mathbf{A}^{j+1} = \mathbf{A}_l^{j+1}$ ;  
8           $I^{j+1} = I_l^{j+1}$ ;  
          Output:  $\mathbf{A}^{j+1}, I^{j+1}, l$  and relevant metadata  
9      else  
10          $\text{prob} = \min \left\{ \exp \left[ \frac{I^j - I_l^{j+1}}{T(l)} \right], 1 \right\}$ ;  
11         Draw  $\text{rand} \sim \text{Unif}(0, 1)$ ;  
12         if  $\text{prob} > \text{rand}$  then  
13             | Output:  $\mathbf{A}^{j+1}, I^{j+1}, l$  and relevant metadata  
14         end  
15     end  
16      $l = l + 1$ ;  
17 end
```

Below listed several features characterise the optimisation procedures needed in projection pursuit

- *Being able to handle non-differentiable PPI function:* The PPI function could be non-differentiable, thus derivative free methods are preferred.
- *Being able to optimise with constraints:* The constraint comes from projection matrix being an orthonormal matrix.
- *Being able to find both local and global maximum:* Although the primary interest is to find the global maximum, local maximum could also reveal structures in the data that are of our interest.

4 Visual diagnostic system

[I would expand this section more as the core contribution. ppp]

4.1 Origin idea

Random search methods has a black-box mechanism and focuses solely on finding the global maximum point, while the projection pursuit problem we have aims at *exploring* the data and thus is interested in how the algorithm finds its maximum. This motivates us **to develop a visual diagnostic system for exploring the optimisation searching path.**

The necessity of developing such a system rather than simply producing different diagnostic plots is because the diagnostics of each variable requires a different function and these functions can't be scaled to other problems. Thus we want to establish a set of rules that can generalise the diagnostic of iterative algorithms.

The idea of generalising all the diagnostic plots under one framework is inspired by the concept of *grammar of graphic*(Wickham 2010), which powers the primary graphical system in R, ggplot2 (Wickham 2016). In grammar of graphic, plots are not defined by its appearance (i.e. boxplot, histogram, scatter plot etc) but by “stacked layers”. By this design, ggplot doesn't need to develop a gazillion of functions that each produces a different

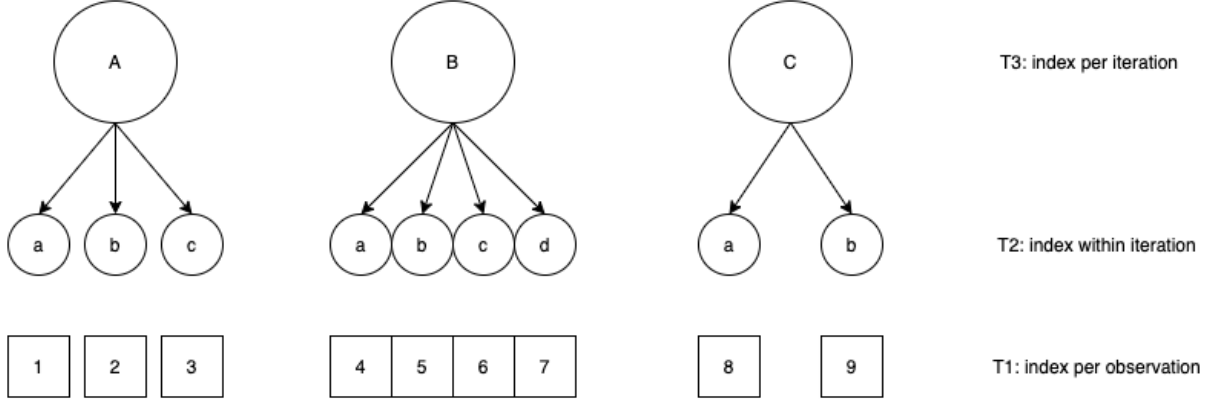


Figure 2: A sketch of the design of iterators in iterative algorithms.

type of plot. Instead, it aesthetically maps the variables to the geometric objects and builds the plot through different layers.

4.2 Global object

Ggplot requires a data frame that contains all the variables to plot and a *global object* is constructed as the data frame supplied to the visual diagnostic plots to better suit the characters of iterative optimisation algorithms. Given an optimisation algorithm, two primary variables of interest are the *decision variable* and the *value of the objective function*. To align with the notation used previously for projection pursuit, we denote the decision variable to be \mathbf{A} , which are matrices in projection pursuit algorithms and value of the objective function to be I , the index function.

Iterators indexes the data collected and has a time series feature that prescribes the its natural ordering. Given the complexity of the algorithm, there could be multiple iterators. The simplest iterator: *index per observation*, is a unique identifier for each observation in the data. For each level of iteration, we design two indices: *index per iteration* is fixed for each observations in one iteration and has an increment of one once a new iteration starts. *index within iteration* has an increment of one for each observation in an iteration and starts over from one once a new iteration starts. A sketch of the difference between these three iterators is provided in Figure 2. Depend on the number of iteration level (L), the total number of iterators will differ and can be calculated as $1 + 2L$.

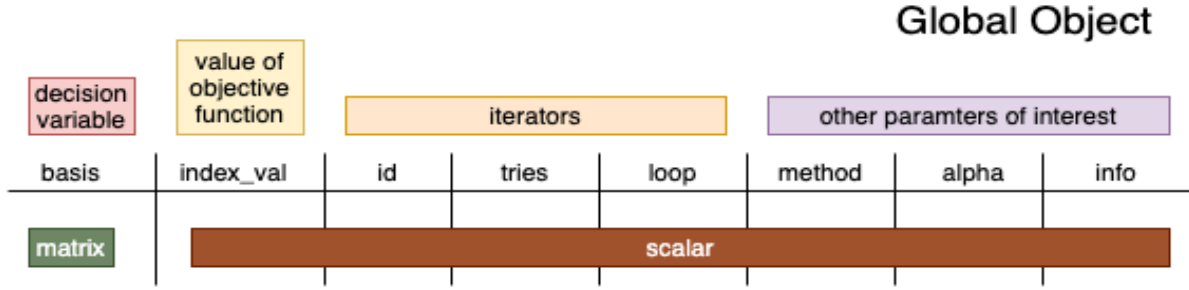


Figure 3: The global object in projection pursuit guided tour.

In projection pursuit optimisation algorithms, there is one level of iteration and hence exists three iterators: `id` indices each observation by a unique number; `tries` is the index per iteration iterator that gets updated once a search-and-interpolate step is finished; and `loop` is the index within iteration iterator and starts over from one at the beginning of a new `tries`.

There could exist other parameters that are also of our interest but can't be classified as one of the three categories. They are defined under the fourth category: *other parameter of interest: O*. In projection pursuit, three other parameters of our interest includes `method`, `alpha` and `info`. `method` identifies the name of the searching method used and we are interested in comparing the performance between different algorithms under direct search. The neighbourhood parameter `alpha` controls the size of the sampling space and we are interested to understand how the searching space shrinks as the algorithm progresses. `info` labels different stages in the searching process. A sketch of the global object for projection pursuit guided tour is presented in Figure 3.

4.3 Notation

Now we formulate the idea of global object using mathematical notation. Let

$$\mathbf{A}_{j,l,k}$$

denotes the projection matrix in iteration j . If it is searching basis produced in the optimisation, then it has a subscript l . If it is an interpolation basis, it has a k subscript. For example $\mathbf{A}_{1,2,/}$ means the second searching basis in the first iteration, $\mathbf{A}_{1,/,2}$ means the

second interpolation basis that on the tour path from the current basis \mathbf{A}^0 to the target basis found in iteration 1: \mathbf{A}^1 . The target basis found in each iteration, $\mathbf{A}_{j,l_j,/}$ can also be simplified using a superscript as \mathbf{A}^j .

Let \mathbf{A}^0 denotes the starting basis. Superscript j is used to denote the number of iteration (index per iteration). Candidate basis are denoted using subscript on the right and interpolation basis are denoted using subscript on the left. For example, \mathbf{A}_2^1 means the second candidate basis in the first iteration and ${}_2\mathbf{A}^1$ means the second interpolation basis that forms the tour path from the current basis \mathbf{A}^0 to the target basis \mathbf{A}^1 . In iteration two, \mathbf{A}^1 becomes the current basis and projection pursuit searches through $\mathbf{A}_1^2, \mathbf{A}_2^2, \dots, \mathbf{A}_{l_2}^2$ to find the target basis \mathbf{A}^2 . The tour path from \mathbf{A}^1 to \mathbf{A}^2 is constructed using the interpolation bases ${}_1\mathbf{A}^2, {}_2\mathbf{A}^2, \dots, {}_{k_2}\mathbf{A}^2$. All the bases together can be arranged in a format that each column represents an iteration as in Equation ??.

$$\begin{bmatrix}
 * & \mathbf{A}_{1,1,*} & \mathbf{A}_{2,l_1+1,*} & \cdots & \mathbf{A}_{j,\sum_{i=1}^{j-1} l_i+1,*} & \cdots & \mathbf{A}_{J,\sum_{i=1}^{J-1} l_i+1,*} \\
 * & \mathbf{A}_{1,2,*} & \mathbf{A}_{2,l_1+2,*} & \cdots & \mathbf{A}_{j,\sum_{i=1}^{j-1} l_i+2,*} & \cdots & \mathbf{A}_{J,\sum_{i=1}^{J-1} l_i+2,*} \\
 \vdots & \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\
 * & \mathbf{A}_{1,l,*} & \mathbf{A}_{2,l_1+l,*} & \cdots & \mathbf{A}_{j,\sum_{i=1}^{j-1} l_i+l,*} & \cdots & \mathbf{A}_{J,\sum_{i=1}^{J-1} l_i+l,*} \\
 \vdots & \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\
 \mathbf{A}_{0,*,*} & \mathbf{A}_{1,l_1,*} & \mathbf{A}_{2,l_1+l_2,*} & \cdots & \mathbf{A}_{j,\sum_{i=1}^j l_i,*} & \cdots & \mathbf{A}_{J,\sum_{i=1}^J l_i,*} \\
 \hline
 * & \mathbf{A}_{1,*,1} & \mathbf{A}_{2,*,k_1+1} & \cdots & \mathbf{A}_{j,*,\sum_{i=1}^{j-1} k_i+1} & \cdots & * \\
 * & \mathbf{A}_{1,*,2} & \mathbf{A}_{2,*,k_1+2} & \cdots & \mathbf{A}_{j,*,\sum_{i=1}^{j-1} k_i+2} & \cdots & * \\
 \vdots & \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\
 * & \mathbf{A}_{1,*,k} & \mathbf{A}_{2,*,k_1+k} & \cdots & \mathbf{A}_{j,*,\sum_{i=1}^{j-1} k_i+k} & \cdots & * \\
 \vdots & \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\
 * & \mathbf{A}_{1,*,k_1} & \mathbf{A}_{2,*,k_1+k_2} & \cdots & \mathbf{A}_{j,*,\sum_{i=1}^j k_i} & \cdots & *
 \end{bmatrix} = \begin{bmatrix}
 \text{candidate basis} \\
 \text{candidate basis} \\
 \vdots \\
 \text{candidate basis} \\
 \vdots \\
 \text{current/target basis} \\
 \hline
 \text{interpolation basis} \\
 \text{interpolation basis} \\
 \vdots \\
 \text{interpolation basis} \\
 \vdots \\
 \text{interpolation basis}
 \end{bmatrix} \quad (3)$$

When aggregated with information of index value and other interesting parameter, it is easier to write the basis as a column vector as in Equation 4 and the full data structure is presneted in Equation 5.

All the basis can also be denoted usign a single subscript by their natural ordering as $\mathbf{A}_0, \mathbf{A}_1, \dots, \mathbf{A}_n$.

$$\begin{bmatrix} \mathbf{A}_{0,*,*} & \mathbf{A}_{1,1,*} & \mathbf{A}_{1,2,*} & \cdots & \mathbf{A}_{1,l_1,*} & \cdots & \mathbf{A}_{1,l_1,*} \\ \mathbf{A}_{1,*,1} & \mathbf{A}_{1,*,2} & \cdots & \mathbf{A}_{1,*,k} & \cdots & \mathbf{A}_{1,*,k_1} & \cdots & \mathbf{A}_{J,l_J,*} \end{bmatrix} \quad (4)$$

$$\begin{bmatrix} \mathbf{A} & I & j & l & k & n & S_1 & \cdots & S_p \\ \mathbf{A}_0 & I_0 & 0 & * & * & 1 & * & \cdots & * \\ \mathbf{A}_1 & I_1 & 1 & 1 & * & 2 & S_{1,1} & \cdots & S_{p,1} \\ \mathbf{A}_2 & I_2 & 1 & 2 & * & 3 & S_{1,2} & \cdots & S_{p,2} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \cdots & \vdots \\ \vdots & \vdots & 1 & l_1 & * & \vdots & \vdots & \cdots & \vdots \\ \vdots & \vdots & 1 & * & 1 & \vdots & \vdots & \cdots & \vdots \\ \vdots & \vdots & 1 & * & 2 & \vdots & \vdots & \cdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \cdots & \vdots \\ \vdots & \vdots & 1 & * & k_1 & \vdots & \vdots & \cdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \cdots & \vdots \\ \mathbf{A}_n & I_n & J & l_J & * & n & S_{1,n} & \cdots & S_{p,n} \end{bmatrix} \quad (5)$$

where $n = 1 + \sum_{i=1}^J l_i + \sum_{i=1}^{J-1} k_i$

4.4 Simulated data

[I would add the maths here to describe the simulations. ppp]

Two set of simulated data are used in the demo of the visual diagnostics in projection pursuit guided tour. A small dataset consists of 1000 randomly simulated observations of five variables ($\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_8, \mathbf{x}_9, \mathbf{x}_{10}$). \mathbf{x}_2 is the informative variable simulated from two bi-modal normal distribution centred at -3 and 3 with variance of one and the other four are simulated from non-informative standard random normal distributions. The data has been scaled to ensure \mathbf{x}_2 has variance of 1. The goal is to find the 1D projection of bi-modal shape, which corresponds to the projection matrix (vector) of $[0, 1, 0, 0, 0]$.

A larger dataset contains more informative variables (\mathbf{x}_3 to \mathbf{x}_7) of different types. The distribution of all the variables except \mathbf{x}_3 is plotted in Figure 4 and \mathbf{x}_3 takes 500 positive one and 500 negative one. [should I introduce the dist for each var? xxx] [also feel like we don't really use \mathbf{x}_3 to \mathbf{x}_6 , should we not mention about these? xxx]

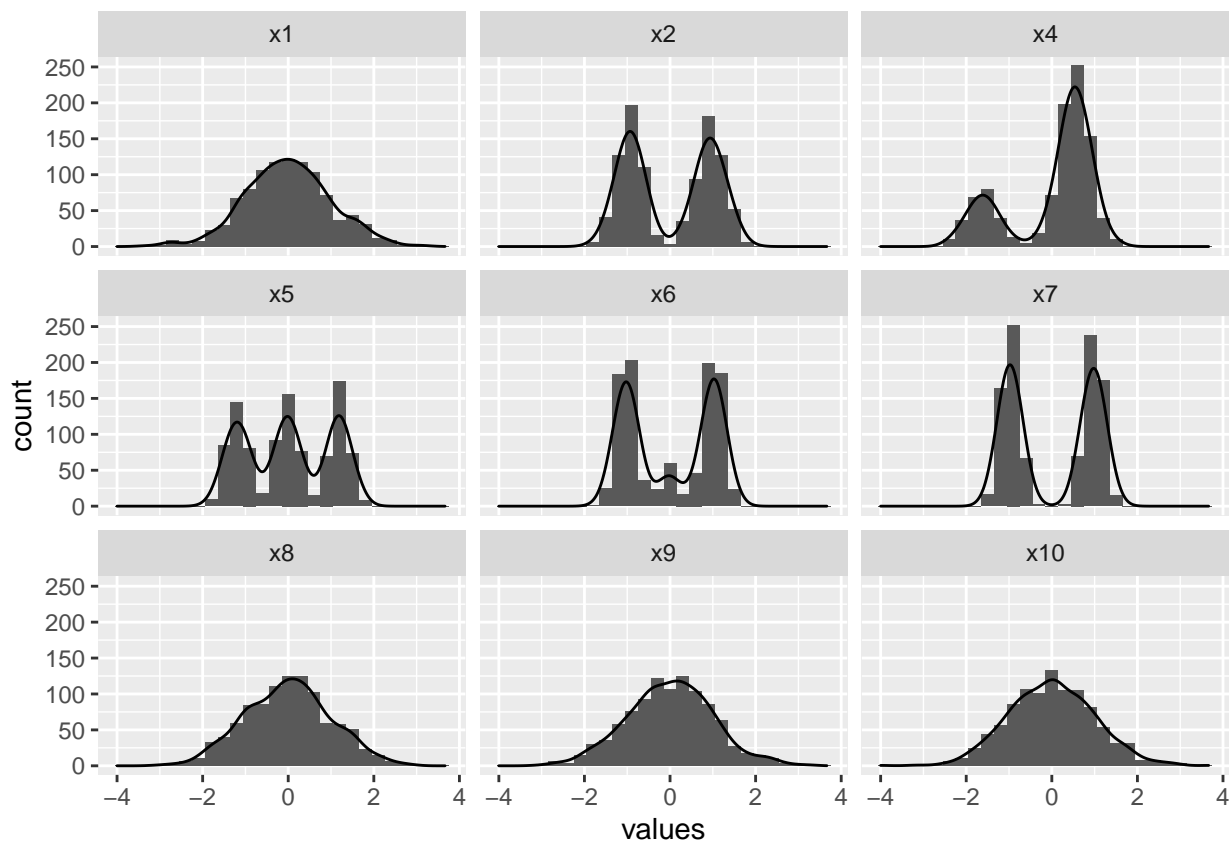


Figure 4: The distribution of simulated data except x3

In tour, when the optimisation ends the global object will be stored and printed and can be turned off via `print = FALSE`. Additional messages during the optimisation can be displayed via `verbose = TRUE`. Below presented the global object of the 1D projection using the small dataset. [the printing is not ideal as it doesn't show all the columns. xxx]

```
## # A tibble: 5 x 8
##   basis          index_val tries info          loop method      alpha    id
##   <list>          <dbl> <dbl> <chr>          <dbl> <chr>          <dbl> <int>
## 1 <dbl[,1] [5 x ~    0.749     1 start          NA <NA>          0.5     1
## 2 <dbl[,1] [5 x ~    0.749     1 direction_sea~    1 search_geode~ NA     2
## 3 <dbl[,1] [5 x ~    0.749     1 direction_sea~    1 search_geode~ NA     3
## 4 <dbl[,1] [5 x ~    0.749     1 direction_sea~    1 search_geode~ NA     4
## 5 <dbl[,1] [5 x ~    0.749     1 direction_sea~    1 search_geode~ NA     5
```

4.5 Visual diagnostics

Below we will present several examples of diagnosing different aspects of the projection pursuit optimisation. We will present:

- 1) static plots to explore the index value,
- 2) animated plots to explore the projection basis and,
- 3) a self-contained example to optimise a complex index function.

In the first two sections, we will first provide a toy example that is easy to grasp and then more examples that can help us to understand the algorithm and parameter choice. Remember the research question we raised earlier, the purpose of visual diagnostics is to understand:

- Whether the algorithm has successfully found the maximum and how the index value changes throughout the algorithm?
- How does the searching space look like, that is, geometrically, where are the projection bases located in the space?

Overall, the first question can be answered using a static plot with x-axis showing the progression of the optimisation and y-axis showing the value of the objective function. The second question can be addressed via visualising the rotating high dimensional space. This seems to be a hard task to see beyond 3D, but the fact is that we never really see all the six faces of a cube at one time. The way humans understand the dimension of a 3D object is either through rotating the physical 3D object or using shade and line type to annotate the 3D object in a paper or electronically. Thus we can do a similar rotation on the screen to perceive even higher dimension. [there is likely to be a learning curve - remember when learning geometry in school, it takes a while to get used to see 3D cube on a paper]. Thus animated visualisation is needed to perceive the optimisation path in the searching space.

The plot design of visual diagnostics depends on the characteristics of the variables to plot. If the searching points are of interest [point geom is not good; need summarise. While updating points are more manageable -> point geom]. However, one must realise that with hundreds or thousands of searching points, exploring the sample space in animation could be slow and this is because of the time it takes to render hundreds point. [strategy may help? - may need more work here.]

4.5.1 Explore the value of objective function

It is a convention to map the time index of a time series on the x-axis and value to the y-axis. In our diagnostic plots, it is to map the iterator on the x-axis and index value to the y-axis. The next question is to ask, given four indices being defined in the global object, which one should be used? *It depends on the .*

As the searching goes towards the end, it is easy to have hundreds of candidate bases in each iteration with very close index value. When plotting, these iteration will occupy the vast majority of the plot space. This motivates to summarise the points in each iteration (**tries**). At each iteration, rather than knowing the index value of *every* points, we are more interested to have a general summary of all the points and more importantly, the point with the largest **index_val** since it prescribes the next geodesic interpolation and future searches.

Using the notation we defined earlier, the plot can be defined as plotting j on the x-axis

and the five point summary of the following groups on the y-axis:

$$\mathbb{1}(I_{j=j^*})$$

Boxplot is a suitable candidate that provides five points summary of the data, however it has one drawback: it doesn't report the number of point in each box. We may risk losing information on how many points it takes to find the target basis by displaying the boxplot alone for all `tries`. Thus, the number of point for each `tries` is displayed at the bottom of each box and we provide options to switch `tries` with small number of points to a point geometry, which is achieved via the `cutoff` argument. A line geometry is also added to link the points with the largest index value in each `tries`. This helps to visualise the improvement made by each `tries`.

Toy example: exploring searching points The larger dataset is used for 2D projection with `search_better` and `max.tries` = 500. In Figure 5, the searching points with `id` and `tries` on the x-axis, colored by `tries`. Label at the bottom indicates the number of observations in each tries and facilitates the choice of cutoff to switch from point geometry to boxplot geometry (`cutoff` = 15). The line geometry suggests the largest improvement happens at `tries` = 5.

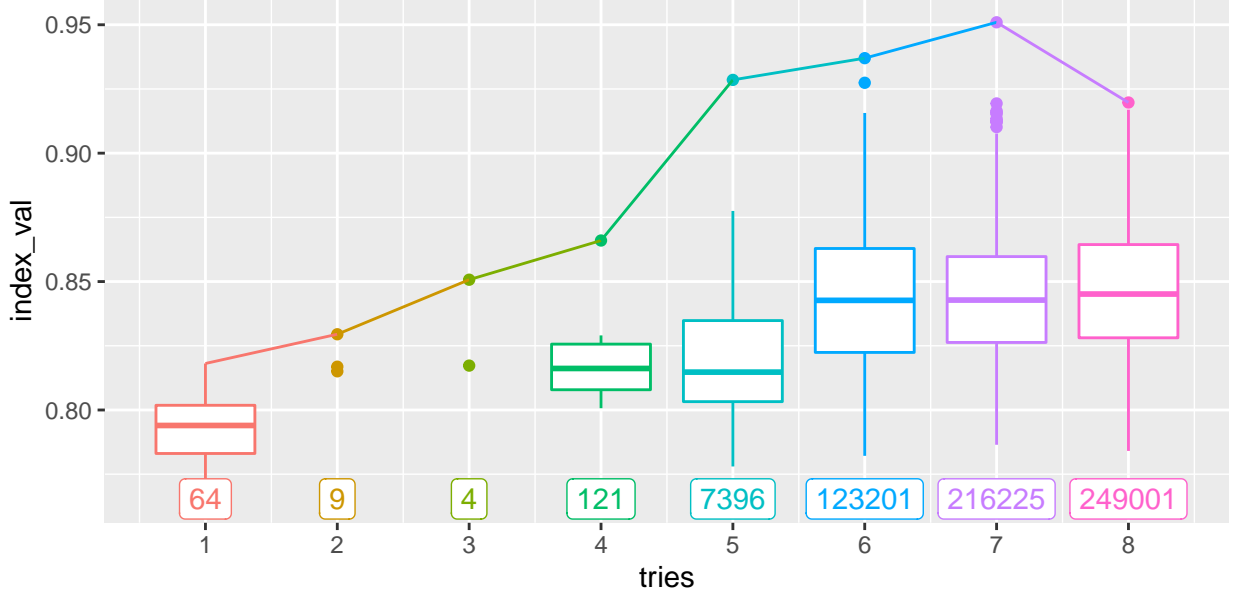


Figure 5: A comparison of plotting the same search points with different plot designs. The left plot doesn't efficiently use the plot space to convey information from the plot while the right plot provides good summarisation of data and number of points in each tries.

The number of interpolation points is much manageable than the searching points and always shows a pattern of gradual increase in the index value. Hence a line plot with interpolating order on the x-axis will be desirable. Using notation, the x-axis is a sequence

$$x = \{1, 2, \dots, k_1, \dots, \sum_{i=1}^J k_i\}$$

and y-axis are the index values defined by $y = I_x$.

A more complex example: Interruption This examples uses `search_better` for a 2D projection on the larger dataset using the `holes` index. In the interpolation stage, continuous frames will be constructed to bridge the project on the current basis to the target basis. The target basis will become the current basis in the next iteration and the searching stage continues to find the next target basis. In the left panel of Figure 6, we observe that there are interpolating bases with higher index values than the target basis and these bases could be used to search for new basis in the next searching stage.

An interruption is then constructed to accept the interpolating bases up to the one

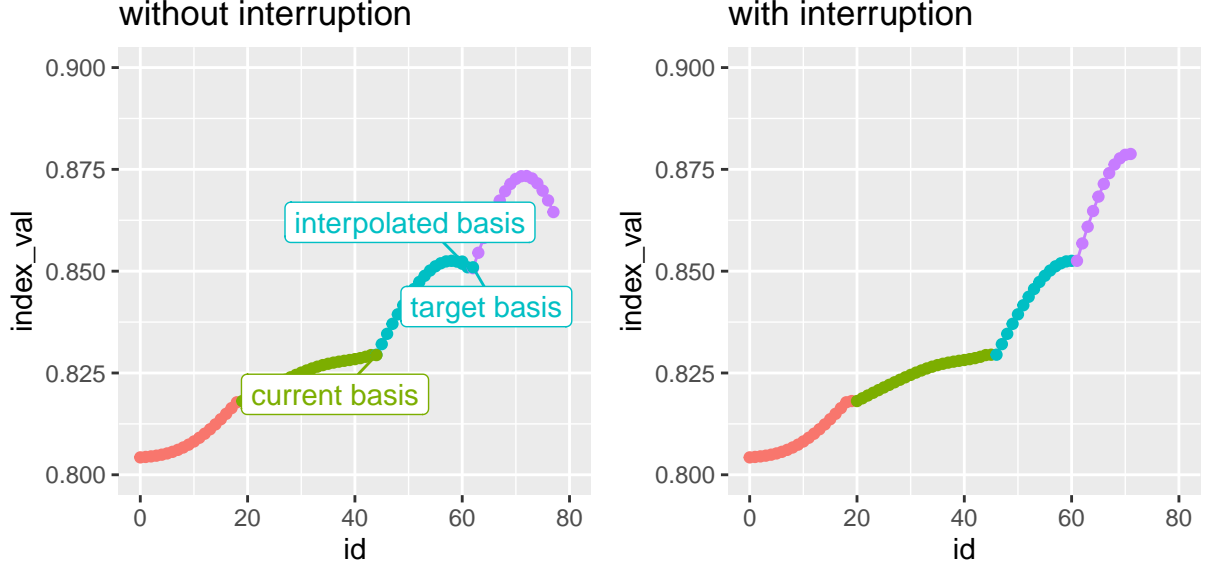


Figure 6: Trace plots of the interpolated basis with and without the interruption. The interruption stops the interpolation when the index value starts to decrease at $\text{id} = 60$. The implementation of the interruption finds an ending basis with higher index value using fewer steps.

with the largest index value and use that basis as the current basis for the next searching stage. After implementing this interruption, the tracing plot with the same configuration is shown on the right panel of Figure 6. In the third `tries`, rather than interpolating fully to the target basis at $\text{id} = 62$, it stops before the index value starts to decrease at $\text{id} = 60$. This implementation results in a higher index value in the end with fewer steps.

Another example: polishing the final projection basis *Example: Polish* In principle, all the optimisation routines should result in the same output on the same problem while this may not be true in real application. This motivates the creation of a polishing search that polishes the ending basis and achieves unity on different methods.

`search_polish` takes the ending basis of a given search as the current basis and uses a brutal-force approach to sample a large number of basis (`n_sample`) in the neighbourhood, whose radius is controlled by `polish_alpha`. Among the `n_sample` basis, the one with the largest index value becomes the candidate basis. If its index value of the candidate basis is

larger than that of the current basis, it becomes the current basis in the next iteration. If no basis is found to have larger index value than the current basis, the searching neighbourhood will be shrunk and the search continues. The polishing search ends when one of the four stopping criteria is satisfied:

- 1) the two basis can't be too close
- 2) the percentage improvement of the index function can't be too small
- 3) the searching neighbourhood can't be too small
- 4) the number of iteration can't exceed the `max.tries`

[should the stopping criteria be more detailed? xxx]

The usage of `search_polish` is as follows. After the first tour, the final basis from the interpolation is extracted and supplied into a new tour with the `start` argument and `search_polish` as the searching function in the `guided_tour`. All the other arguments should remain the same.

```
set.seed(123456)
holes_2d_geo <- animate_xy(data_mult[,c(1,2, 7:10)], tour_path =
  guided_tour(holes(), d = 2,
    search_f = tourr::search_geodesic),
  rescale = FALSE, verbose = TRUE)

last_basis <- holes_2d_geo %>% filter(info == "interpolation") %>%
  tail(1) %>% pull(basis) %>% .[[1]]

set.seed(123456)
holes_2d_geo_polish <- animate_xy(data_mult[,c(1,2, 7:10)], tour_path =
  guided_tour(holes(), d = 2,
    search_f = tourr::search_polish),
  rescale = FALSE, verbose = TRUE,
  start = last_basis)
```

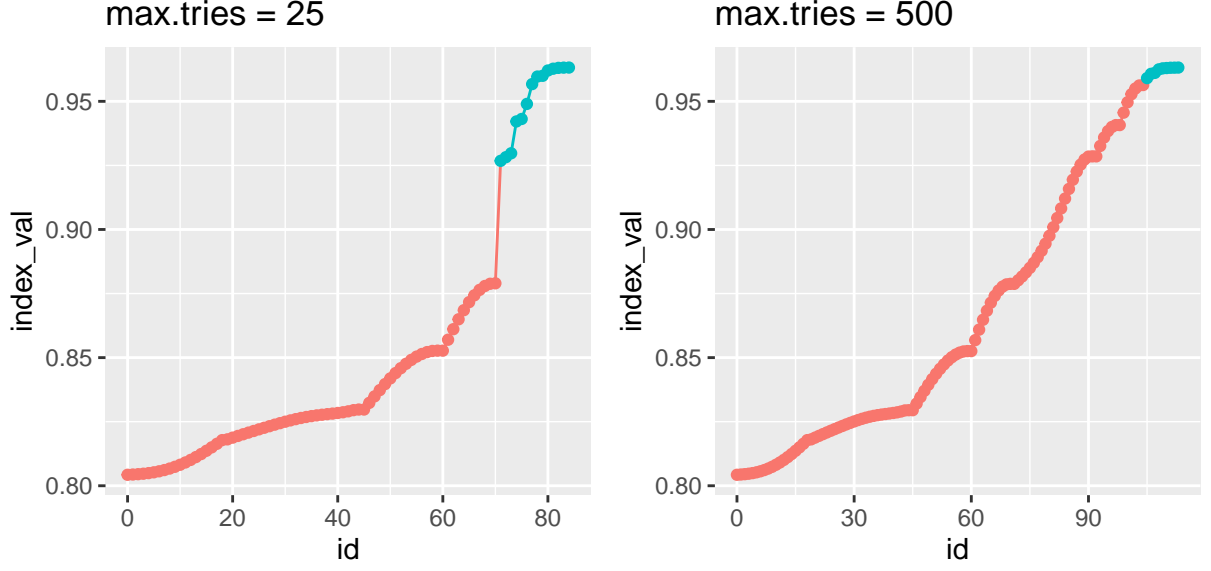


Figure 7: Breakdown of index value when using different `max.tries` in `search_better` in conjunction with `search_polish`. Both attain the same final index value after the polishing while using a `max.tries` 25 is not sufficient to find the true maximum.

The following example conducted a 2D projection on the larger dataset using `search_better` with different configurations. `max.tries` is a hyperparameter that controls the maximum number of try without improvement and its default value is 25. As shown in Figure 7, both trials attain the same index value after polishing while a small `max.tries` of 25 is not sufficient for `search_better` to find the global maximum. The `max.tries` argument needs to be adjusted to ensure it is sufficient to explore the parameter space.

4.5.2 Explore searching space

In projection pursuit, the projection bases $\mathbf{A}_{p \times d}$ are usually of dimension $p \times d$ and hence can't be visualised in a 2D plot. An option to explore the searching space of these bases is to explore a reduced space via principal component analysis (PCA). We first flat each projection bases \mathbf{A} to a row vector and stack them by rows. This gives a matrix of dimension n by $(p * d)$. Centering and scaling are usually done before performing PCA and we denotes the pre-processed matrix as $\tilde{\mathbf{A}}_{n \times (p \times d)}$. By definition, the principal components maximise the variance of the projection bases in the reduced space. Using lagrange multiplier, this is equivalent to find the eigenvector of the variance-covariance matrix: $\Sigma = \tilde{\mathbf{A}}^T \tilde{\mathbf{A}}$. Decomposing $\tilde{\mathbf{A}}$ using singular value decomposition (SVD) gives

$$\tilde{\mathbf{A}} = \mathbf{U} \lambda \mathbf{V}$$

and hence the variance-covariance matrix

$$\Sigma = \tilde{\mathbf{A}}^T \tilde{\mathbf{A}} = (\mathbf{V}^T \lambda^T \mathbf{U}^T)(\mathbf{U} \lambda \mathbf{V}) = \mathbf{V}^T \lambda^2 \mathbf{V}$$

has eigenvector $\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_{p \times d}]$. The projection mapping to the x- and y-axis uses the first two principal components and hence

$$\tilde{\mathbf{A}} \mathbf{v}_1 \text{ and } \tilde{\mathbf{A}} \mathbf{v}_2$$

The resulting two vectors are of dimension $n \times 1$, which is compatible with the dimension of the global object. This means all the iterators and interesting other variables defined in the global object can be binded with the these two projections and can be mapped to other aesthetics (i.e. color).

In this way, the first two principle components will take up the x and y axis and all the other information will be mapped using other aesthetic attribution. [need to think further about how different type of variables can be shown]. While explore the reduced space is an initial attempt to understand the searching space, there are existing technology for rotating a higher dimensional space for visualisation. Geozoo is an option. It generates random points on the high dimensional space and we can overlay it with the points on the optimisation path to visualise the spread of it on the high-D sphere.

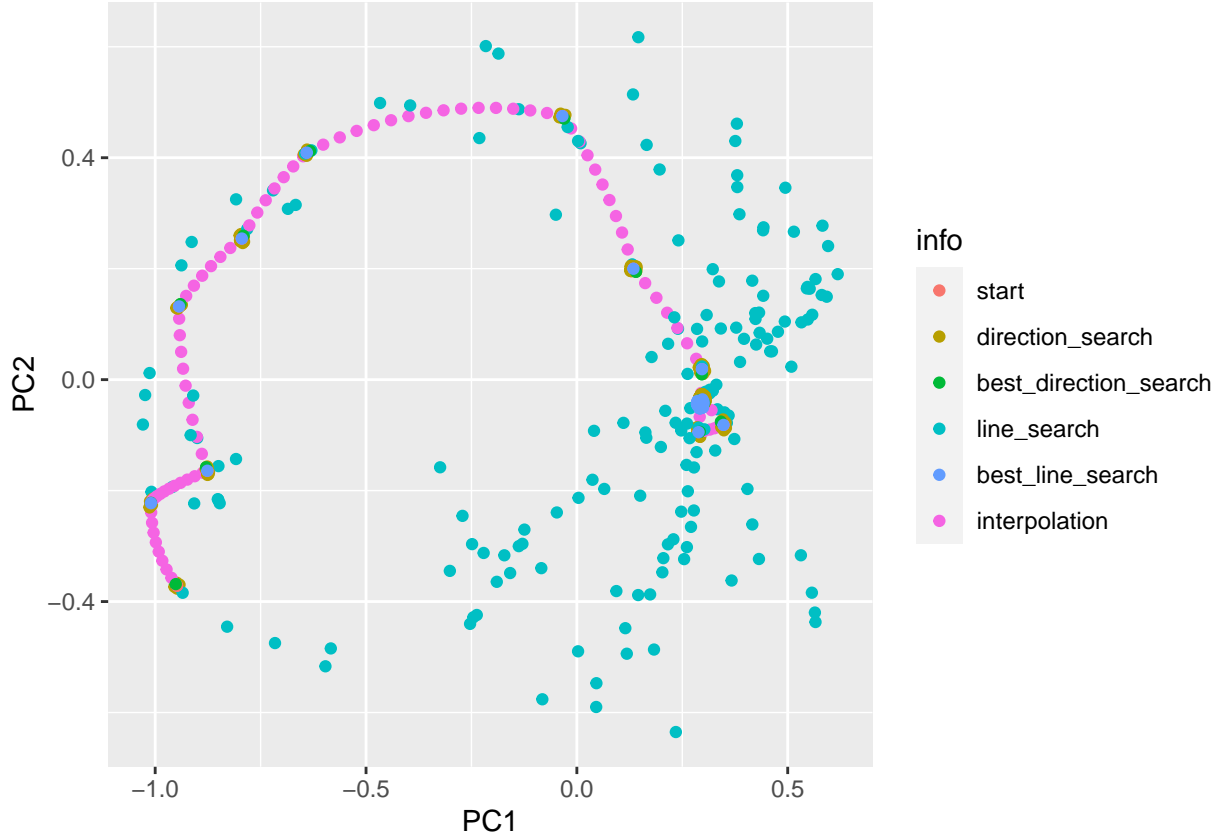


Figure 8: PCA plot of search geodesic Coloring by info allows for better understanding of each stage in the geodesic search

A toy example: understand different stage of search_geodesic *Example: understand search_geodesic* [feel like this example is merely explaining search geodesic algorithm, so maybe introduce the animated plot here? xxx] `search_geodesic` is a two-stage ascending algorithm with four different stages in the search and a PCA plot useful to understand how the algorithm progresses and the relative position of each basis in the PCA projected 2D space. Starting from the start basis, a directional search is conducted in a narrow neighbourhood on five random directions. The best one is picked and a line search is then run on the geodesic direction to find the target basis. The starting and target bases are then interpolated. In the next iteration, the target basis becomes the current basis and then procedures continues.

A more complex example: Choosing the initial value for polishing parameter

Example: initial value for polishing alpha `search_polish` is a brutal-force algorithm that evaluate 1000 points in the neighbourhood at each loop. Setting an appropriate initial value for `polish_alpha` would avoid wasting search on large vector space that are not likely to produce higher index value. The default initial value for polishing step is 0.5 and we are interested in whether this is an appropriate initial value to use after `search_geodesic`. The problem is a 1D projection of the small dataset using `search_geodesic` and followed by `search_polish`. The top-left panel of Figure 9 displays all the projection bases on the first two principal components, colored by the `polish_alpha`. We can observe that rather than concentrating on the ending basis from `search_geodesic` as what polishing step is designed, `search_polish` searches a much larger vector space, which is unnecessary. Thus a customised smaller initial value for `polish_alpha` would be ideal. One way to do this is to initialised `polish_alpha` as the projection distance between the last two target bases. The top-right panel of Figure 9 shows a more desirable concentrated searching space near the ending basis. Both specifications of initial value allow the searches to reach the same ending index values.

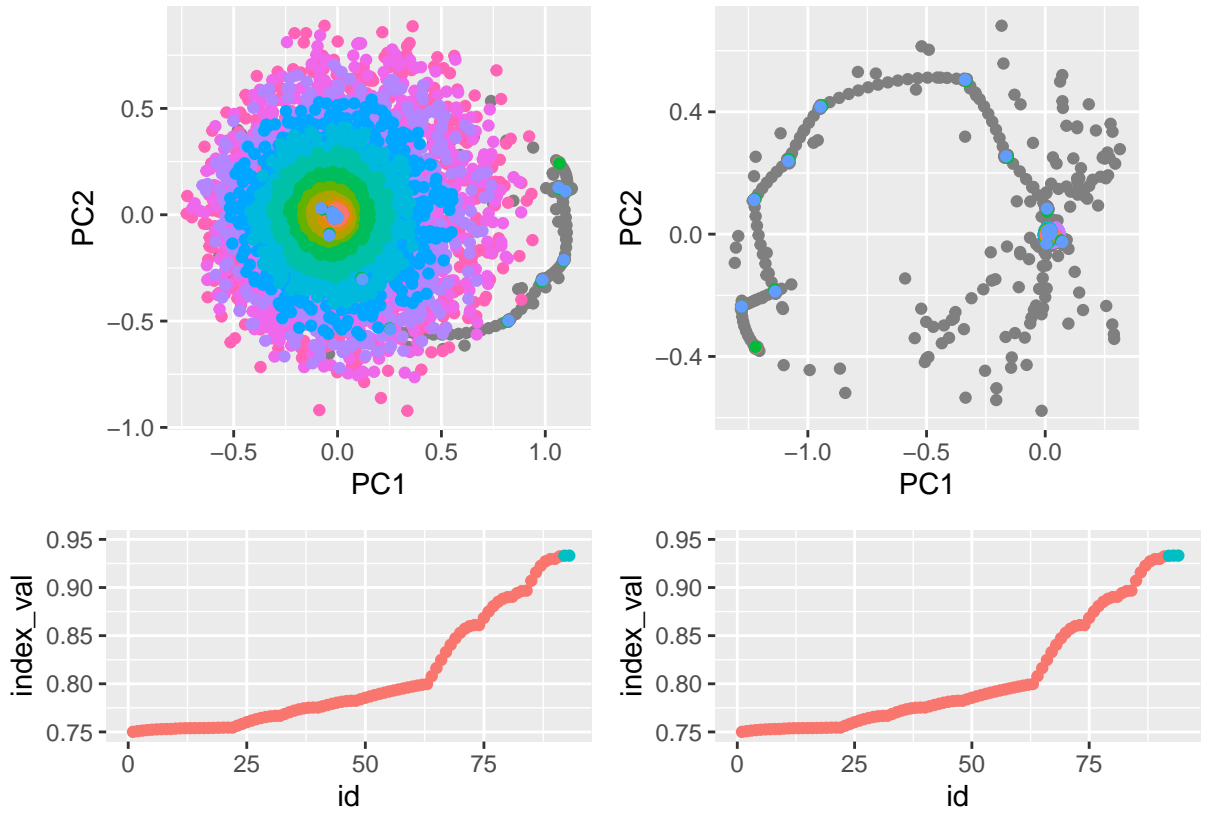


Figure 9: PCA plot of two different polish alpha initialisations. A default polish alpha = 0.5 searches a larger space that is unnecessary while a small customised initial value of polish alpha will search near the ending basis. Both intialisations reach the same ending index values.

4.5.3 A comprehensive example of diagnosing a noisy index function

The interpolation path of holes index, as seen in Figure 6, is smooth, while this may not be the case for more complicated index functions. `kol_cdf` index, an 1D projection index function based on Kolmogorov test, compares the difference between the 1D projected data, $\mathbf{P}_{n \times 1}$ and a randomly generated normal distribution, y_n based on the empirical cumulated distribution function (ECDF). Denotes the ECDF function as $F(u)$ with subscript indicating the variable, the Kolmogorov statistics defined by

$$\max [F_{\mathbf{P}}(u) - F_y(u)]$$

can be seen as a function of the projection matrix $\mathbf{A}_{p \times 1}$ and hence a valid index function.

Explore index value Figure 10 compares the tracing plot of the interpolating points when using different optimisation algorithms: `search_geodesic` and `search_better`. One can observe that

- The index value of `kol_cdf` index is much smaller than that of holes index
- The link of index values from interpolation bases are no longer smooth
- Both algorithms reach a similar final index value after polishing

Polishing step has done much more work to find the final index value ub `search_geodesic` than `search_better` and this indicates `kol_cdf` function favours of a random search method than ascent method.

Now we enlarge the dataset to include two informative variables: `x2` and `x3` and remain 1D projection. In this case, two local maximum appear when projection matrix being $[0, 1, 0, 0, 0, 0]$ and $[0, 0, 1, 0, 0, 0]$.

Using different seeds in `search_better` allows us to find both local maximum as in Figure 11. Comparing the maximum of both, we can see that the global maximum happens when `x2` is found. It is natural to ask then if there is an algorithm that can find the global maximum without trying on different seeds? `search_better_random` manages to do it via a metropolis-hasting random search as shown in Figure 12, although at a higher cost of number of points to evaluate.

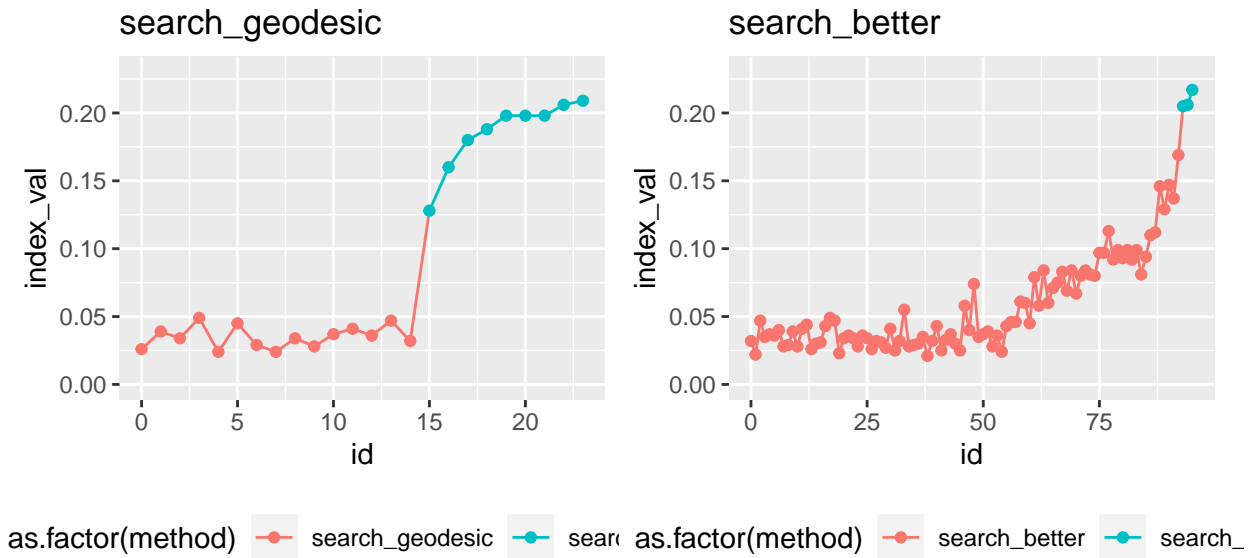


Figure 10: Comparison of two different searching methods: `search_geodesic` and `search_better` on 1D projection problem for a noisier index: `kol.cdf`. The geodesic search rely heavily on the polishing step to find the final index value while search better works well.

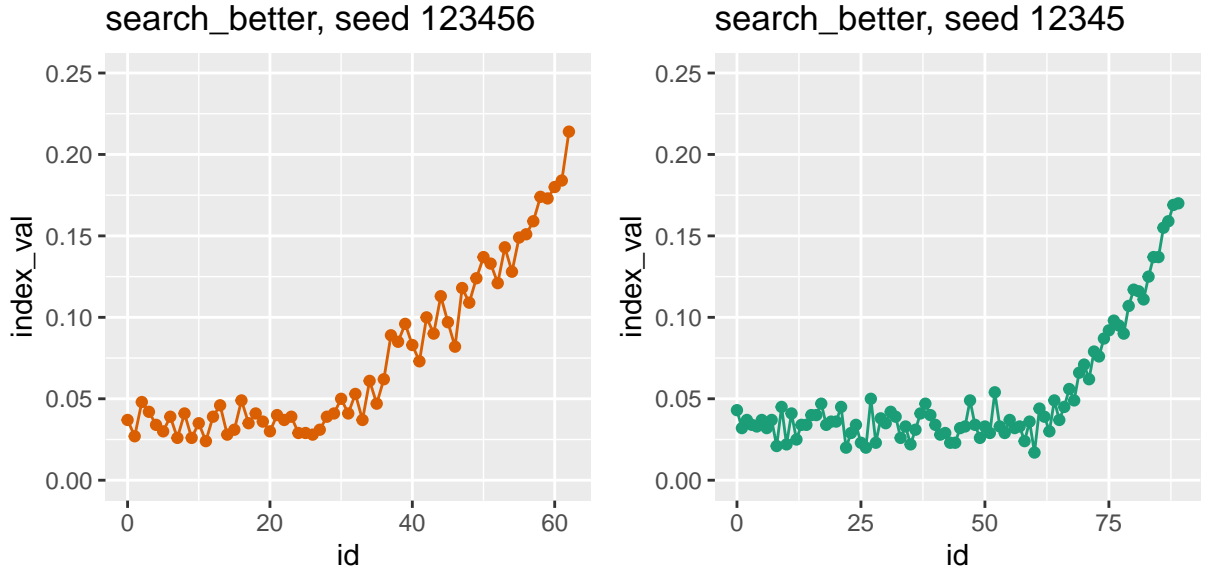


Figure 11: The trace plot search better in a 1D projection problem with two informative variables using different seeds (without polishing). Since there are two informative variables, setting different value for seed will lead search better to find either of the local maximum.

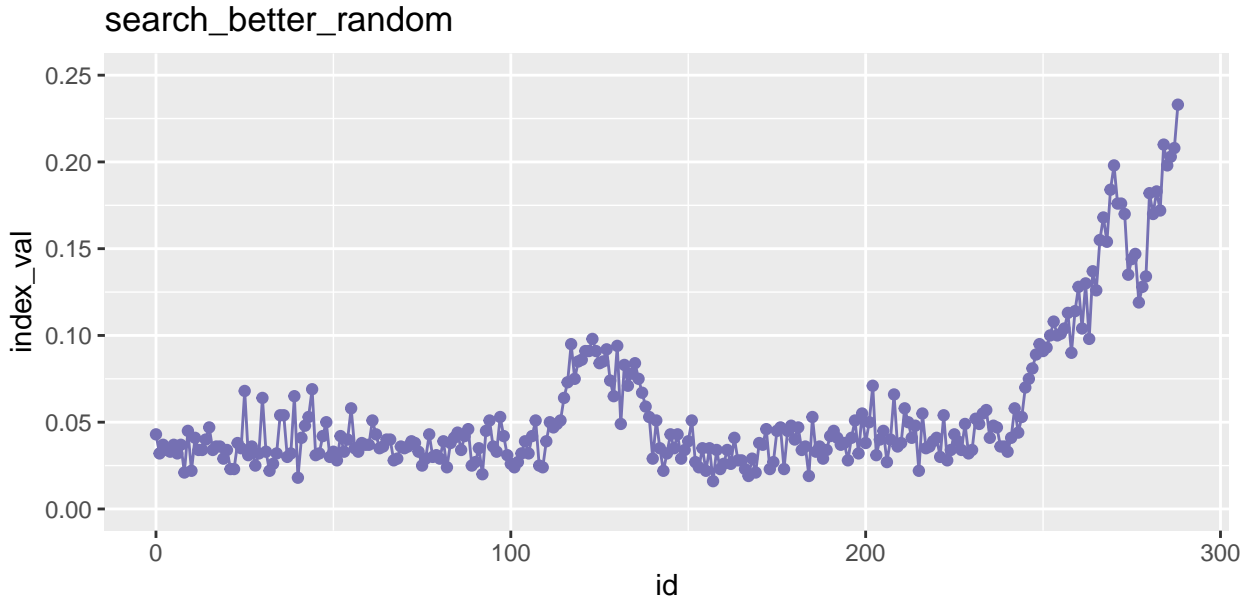


Figure 12: Using search better random for the problem above will result in finding the global maximum but much larger number of iteration is needed.

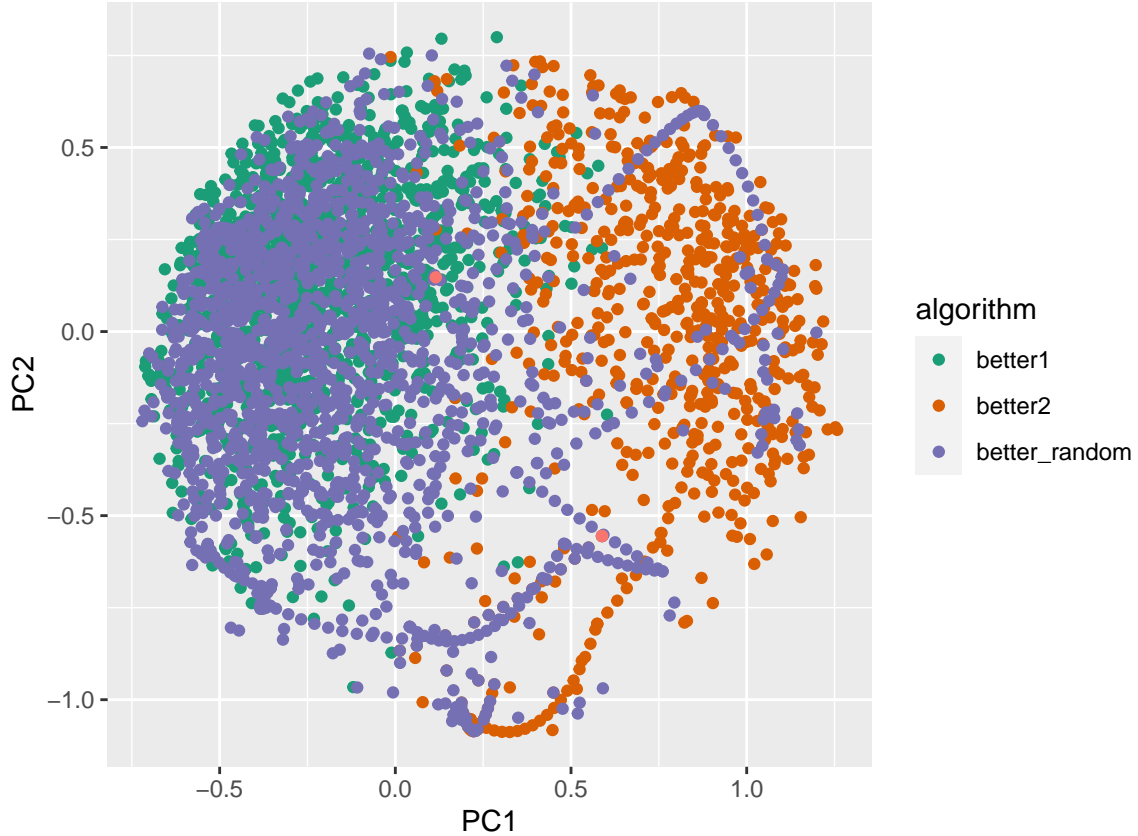


Figure 13: The projected projection basis using principal components. The bases from better1 and better2 only search a proportion of the searching space while better_random produces a more exhaustive search. The large overlapping of better1 and better_random is explained by the fact that both algorithms finds x_2 in the end.

Explore searching space We can also plot the searching points of all the three algorithms in the searching space and explore their relative position against each other using principal components. As shown in Figure 13, the bases from better1 and better2 only search a proportion of the searching space while better_random produces a more exhaustive search. The large overlapping of better1 and better_random is explained by the fact that both algorithms finds x_2 in the end.

5 Implementation: Ferrn pacakge

Everything is coded up in a package. Package structure

6 Conclusion

References

- Bertsimas, D., Tsitsiklis, J. et al. (1993), ‘Simulated annealing’, *Statistical science* **8**(1), 10–15.
- Buja, A., Cook, D., Asimov, D. & Hurley, C. (2005), ‘Computational methods for high-dimensional rotations in data visualization’, *Handbook of statistics* **24**, 391–413.
- Cook, D., Buja, A. & Cabrera, J. (1993), ‘Projection pursuit indexes based on orthonormal function expansions’, *Journal of Computational and Graphical Statistics* **2**(3), 225–250.
- Cook, D., Buja, A., Cabrera, J. & Hurley, C. (1995), ‘Grand tour and projection pursuit’, *Journal of Computational and Graphical Statistics* **4**(3), 155–172.
- Friedman, J. H. & Tukey, J. W. (1974), ‘A projection pursuit algorithm for exploratory data analysis’, *IEEE Transactions on computers* **100**(9), 881–890.
- Hall, P. et al. (1989), ‘On polynomial-based projection indices for exploratory projection pursuit’, *The Annals of Statistics* **17**(2), 589–605.
- Hooke, R. & Jeeves, T. A. (1961), “‘direct search’ solution of numerical and statistical problems’, *Journal of the ACM (JACM)* **8**(2), 212–229.
- Kirkpatrick, S., Gelatt, C. D. & Vecchi, M. P. (1983), ‘Optimization by simulated annealing’, *science* **220**(4598), 671–680.
- Lee, E., Cook, D., Klinke, S. & Lumley, T. (2005), ‘Projection pursuit for exploratory supervised classification’, *Journal of Computational and graphical Statistics* **14**(4), 831–846.
- Lee, E.-K. & Cook, D. (2010), ‘A projection pursuit index for large p small n data’, *Statistics and Computing* **20**(3), 381–392.
- Posse, C. (1995), ‘Projection pursuit exploratory data analysis’, *Computational Statistics & data analysis* **20**(6), 669–687.

Wickham, H. (2010), ‘A layered grammar of graphics’, *Journal of Computational and Graphical Statistics* **19**(1), 3–28.

Wickham, H. (2016), *ggplot2: Elegant Graphics for Data Analysis*, Springer-Verlag New York.

URL: <https://ggplot2.tidyverse.org>

Wickham, H., Cook, D., Hofmann, H. & Buja, A. (2011), ‘tourr: An R package for exploring multivariate data with projections’, *Journal of Statistical Software* **40**(2), 1–18.

URL: <http://www.jstatsoft.org/v40/i02/>