

# Title here

Author 1 \*

Department of YYY, University of XXX  
and

Author 2

Department of ZZZ, University of WWW

May 17, 2020

## Abstract

1. check consistency of using PP for projection pursuit and PPI for projection pursuit index

*Keywords:* 3 to 6 keywords, that do not appear in the title

---

\*The authors gratefully acknowledge ...

# 1 Introduction

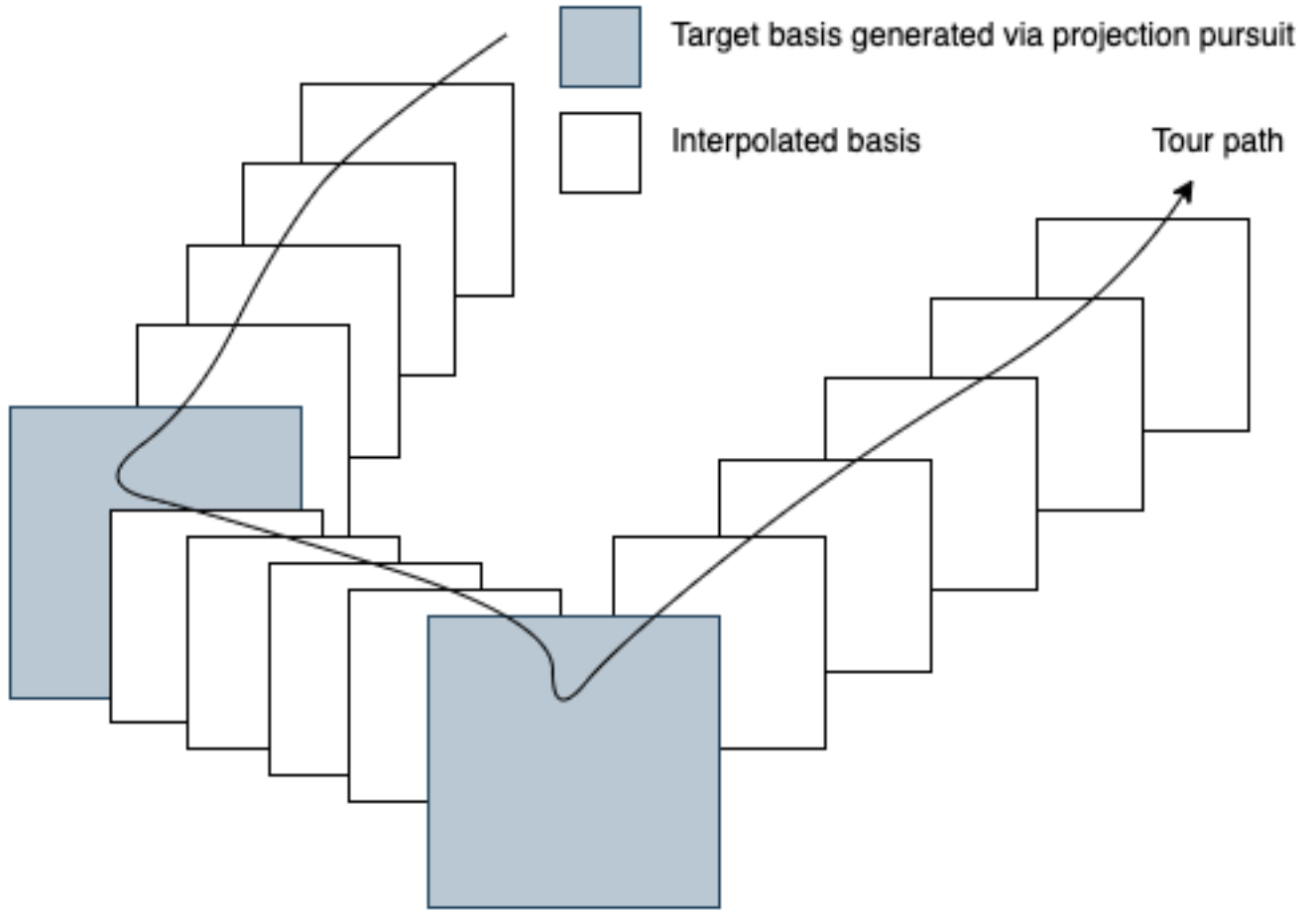
## 1.1 Tour

Tour provides a way to explore multivariate data interactively via an established tour path. A tour path is formed by interpolating between randomly generated plane. Different types of tours are available depends on the purpose of the exploration i.e. a grand tour is suitable for randomly exploring the data from different angles; a guided tour detects a particular structure in the data. a manual tour allows user to manually control the projection (Cook & Buja 1997).

## 1.2 Guided tour

Guided tour is usually used in conjunction projection pursuit, a method coined by Friedman & Tukey (1974) to detect “interesting” low-diemsion projection of multivariate data. A projection pursuit requires the definition of a projection pursuit index function and an optimisation routine. The projection pursuit index measures the “interestingness” of data defined as its departure from normality. Numerous indices have been proposed in the literature, including lengendre 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). An optimisation routine is required to find the projection basis (thus projection) that maximises the projection pursuit index. The discussion of existing optimisation procedure will be discussed in the next section.

Guided tour creates visualisation for the projections found by projection pursuit by constructing a tour path. An illustration modified from (Buja et al. 2005) entails the geodesic interpolation bewteen the plane generated by projection pursuit.



### 1.3 Optimisation methods in projection pursuit literature

As Friedman & Tukey (1974) said “..., the technique use for maximizing the (one- and two-dimensional) projection index strongly influences both the statistical and the computational aspects of the procedure.” The quality of the optimisation procedure largely affect the tour view and thus, the interesting projection one could possibly observe. An ideal optimisation procedure needs to have following characteristics:

- *Being able to handle non-differentiable index function:* the index function could be noisy and non-differentible
- *Being able to optimise with constraints:* the projection matrix is restricted to an orthornormal matrix.

- *Being able to reveal both local and global maximum:* Although the primary interest is to find the global maximum, distinct local structures are also of our interest.

Below present three existing methods in tour.

Posse (1995) presented a random search algorithm that samples new basis in the neighbourhood of the current basis. The neighbourhood is defined via the radius of the  $p$ -dimensional sphere,  $c$ . The new basis is taken as the target basis if it has higher index value, or the sampling continues. If no basis is found to have higher index value after a certain number of tries  $n$ , the radius  $c$  is halved. The algorithm stops when the maximum number of iteration is attained or the radius  $c$  is less than a pre-determined number. [Pursuit package uses this method and it works great! But I don't think we implement this - although don't think it is too hard to do it].

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.

Simulated annealing (Bertsimas et al. 1993, Kirkpatrick et al. (1983)) is a non-derivative procedure based on a non-increasing cooling scheme  $T(i)$ . Given an initial  $T_0$ , the temperature at iteration  $i$  is defined as  $T(i) = \frac{T_0}{\log(i+1)}$ . The simulated annealing algorithm works as follows. Given a neighbourhood parameter  $\alpha$  and a randomly generated orthonormal basis  $B$ , a candidate basis is constructed as  $B_j = (1 - \alpha)B_i + \alpha B$  where  $B_i$  is the current basis. 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  $A = \min\left(\exp\left(-\frac{I(B_j) - I(B_i)}{T(i)}\right), 1\right)$  where  $I(\cdot)$  is the index function.

## 1.4 problems and difficulties in PP optimisation

Below listed several issues in projection pursuit optimisation. Some of them are general problems in optimisation literature, 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 on non-smooth function:* When the objective function is non-differentiable, derivative information can not be obtained. This 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 because 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. While this practice supports reproducibility, less efforts has been made to guarantee different seeds will provide the same result.
- *high-dimensional decision variable:* In projection pursuit, the decision variable is the entry in the projection matrix, which is usually 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 maximise 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 prograssion of multi-parameter decision variable

- understanding learning rate - neighbourhood parameter
- understand where the local & global maximum is found - trace plot - see if noisy function

## 2 Recording the guided tour

### 2.1 Tour components

Guided tour, along with other types of tour, has been implemented in the *tourr* package in R, available on the Comprehensive R Archive Network at <https://cran.r-project.org/web/packages/tourr/> (Wickham et al. 2011). A tour includes two major components: a *generator* that generating the projection basis according to projection pursuit and an *interpolator* that performing geodesic interpolation between the projection basis.

The psudo-code below illustrates the implementation of guided tour in the *tourr* package. Given an projection pursuit index function and a randomly generated projection basis (current basis), the optimisation procedure produces a target basis inside `generator()`. Both the current basis and the target basis will be supplied to `tour_path()` to prepare information needed for constructing a geodesic path. This information is then used to compute a series of interpolating bases inside the `tour()` function. All the basis will be sent to create animation for visualising the tour in the `animate()` function.

```
animation <- function(){

  # compute projection basis
  tour <- function(){

    # construct bases on the tour path
    new_geodesic_path <- function(){
      tour_path <- function(){

        # GENERATOR: generate projection basis via projection pursuit
```

```

guided_tour <- function(){
  generator <- function(){

    # define projection pursuit index
    # generate the target basis from the current basis via optimisation
  }
}

# prepare geodesic information needed for interpolating along the tour path
}

# INTERPOLATOR: interpolate between the current and target basis
function(){
  # generate interpolating bases on the geodesic path
}
}

# animate according to different display methods
}

```

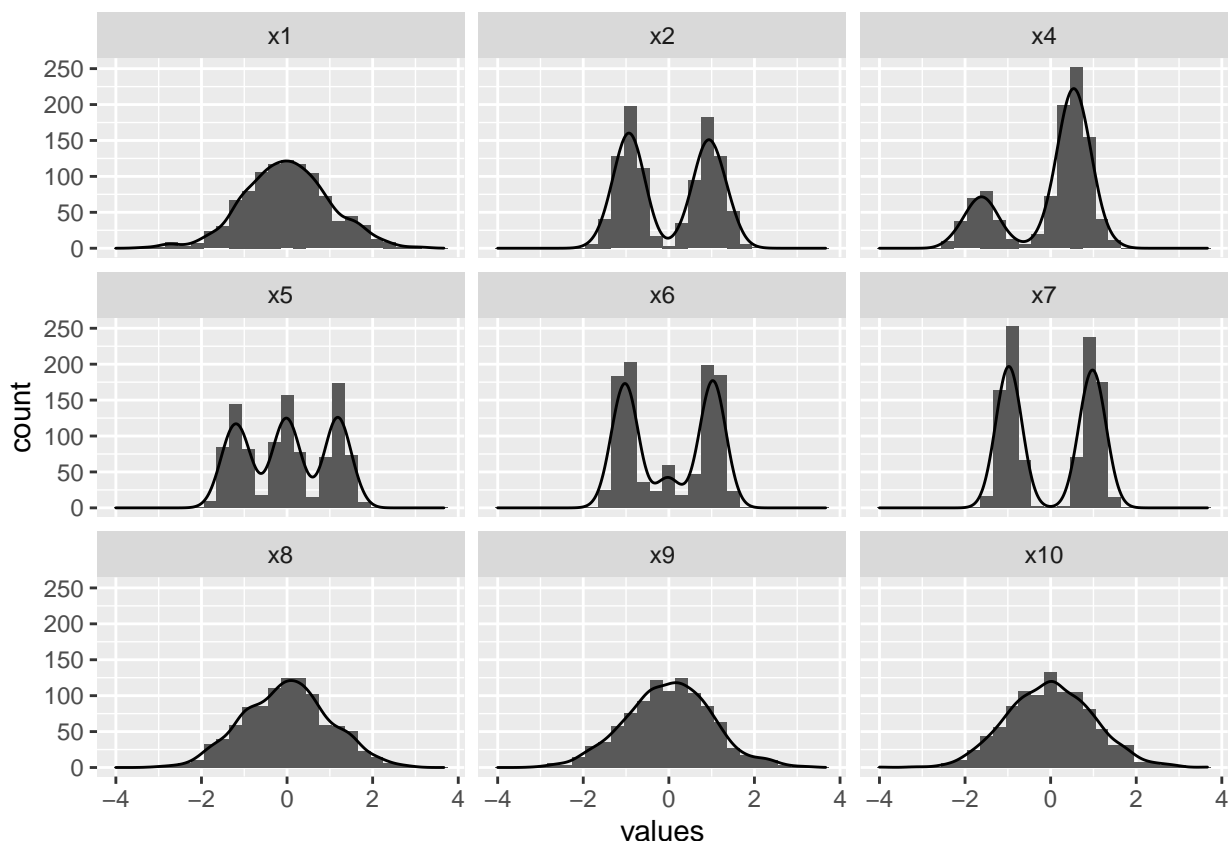
## 2.2 Data Structure

In the current implementation of the `tourr` package, while the target basis generated by the projection pursuit can be accessed later via `save_history()`, interpolating bases and those randomly nearby bases generated in the optimisation are not stored. This creates difficulties for fully understand the behaviour of the optimisation and interpolation of tour in complex scenario **[needa rephrase this part]**.

Two set of simulated data are used in the demonstration of the visualiation and diagnostics of the tour optimisation. A small dataset consists of 1000 randomly simulated observations of five variables (`x1`, `x2`, `x8`, `x9`, `x10`). `x2` is the informative variable simulated

from two bi-modal normal distribution centered at -3 and 3 with variance being 1 and the other four are simulated from  $N(0, 1)$ . The data has been scaled to ensure  $x_2$  has variance of 1.

A larger dataset contains more informative variables ( $x_3$  to  $x_7$ ) of different types.  $x_3$  takes 500 positive one and 500 negative one. The distribution of all the variables except  $x_3$  is plotted below. *[should I introduce the dist for each var?]*



Once the dataset is sent to the `tourr` package, all the information generated will be stored in a global structure. The global structure consists of six columns: `basis`, `index_val`, `tries`, `info`, `loop`, `id` and captured all the basis generated during whole tour process. The example below presents the global object of a 1D projection of the small dataset with geodesic searching method.

```
holes_1d_geo %>% head(5)
```

```
## # A tibble: 5 x 7
```

```
##   basis          index_val tries info          loop method          id
```



##	<list>	<dbl>	<dbl>	<chr>	<dbl>	<chr>	<int>
## 1	<dbl[,1] [5 x 1]>	0.749	1	start	NA	<NA>	1
## 2	<dbl[,1] [5 x 1]>	0.749	1	direction_search	1	search_geodesic	2
## 3	<dbl[,1] [5 x 1]>	0.749	1	direction_search	1	search_geodesic	3
## 4	<dbl[,1] [5 x 1]>	0.749	1	direction_search	1	search_geodesic	4
## 5	<dbl[,1] [5 x 1]>	0.749	1	direction_search	1	search_geodesic	5

`tries` has an increment of one once the generator is called (equivalently a new target basis is generated); `info` records the stage the basis is in. This would include the interpolation stage and the detailed stage in the optimisation i.e. `direction_search`, `best_direction_search`, `line_search` and `best_line_search` for geodesic searching (`search_geodesic`); `random_search` and `new_basis` for simulating annealing (`search_better`). `loop` is the counter used for the optimisation procedure and thus will be `NA` for interpolation steps. `id` creates a sequential order of the basis. This information will be stored and printed when the optimisation ends and can be turned off via `print = FALSE`. Additional messages during the optimisation can be displayed via `verbose = TRUE`. Another examples is a 2D projection of the larger dataset with two informative variable (`x2` and `x7`) using `search_better` method. Notice in this example, the dimension of the bases becomes 6 by 2.

```
holes_2d_better %>% head(5)
```

```
## # A tibble: 5 x 7
```

##	basis	index_val	tries	info	loop	method	id
##	<list>	<dbl>	<dbl>	<chr>	<dbl>	<chr>	<int>
## 1	<dbl[,2] [6 x 2]>	0.804	1	start	NA	<NA>	1
## 2	<dbl[,2] [6 x 2]>	0.793	1	random_search	1	search_better	2
## 3	<dbl[,2] [6 x 2]>	0.784	1	random_search	2	search_better	3
## 4	<dbl[,2] [6 x 2]>	0.773	1	random_search	3	search_better	4
## 5	<dbl[,2] [6 x 2]>	0.795	1	random_search	4	search_better	5

## 3 Visual methods

### 3.1 add something about general vis

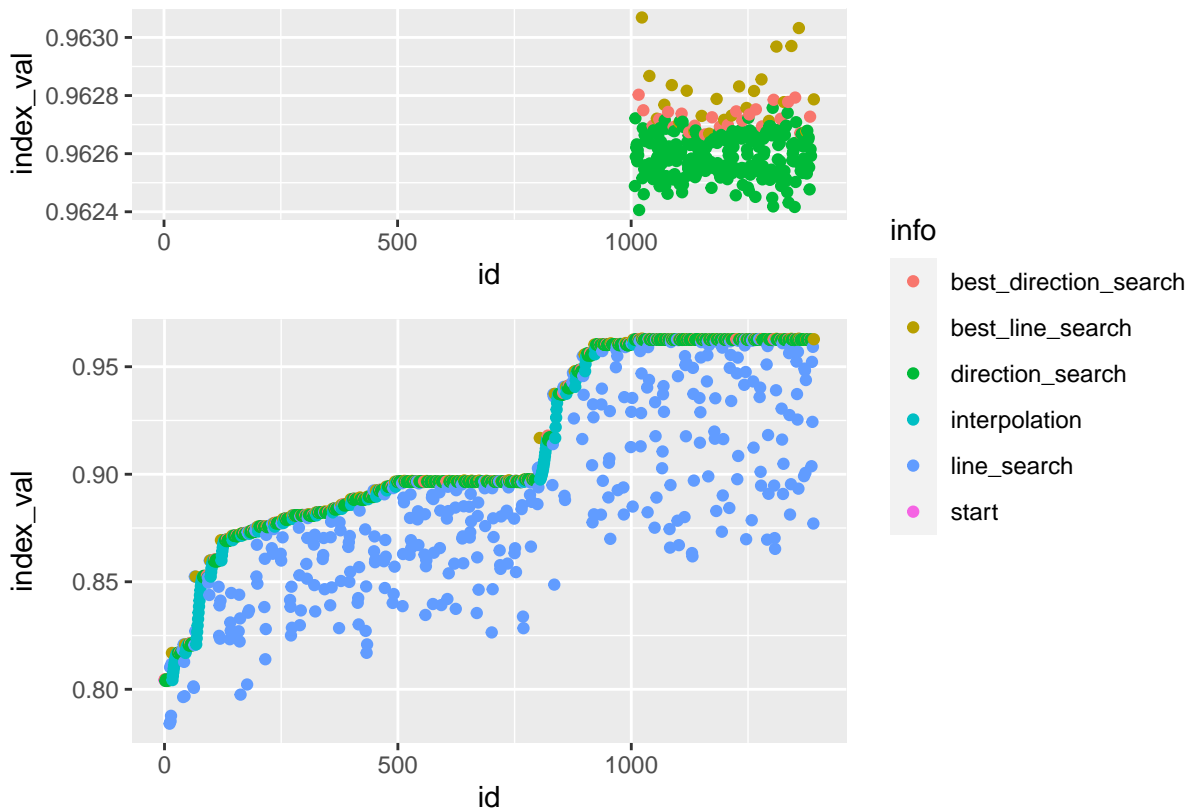
The creation of the global object facilitates the diagnostics of the optimisation in the projection pursuit guided tour. Different diagnostic plots could be made to understand the optimisation [pretty cheesy here]

*The examples and results below still need to be modified as the development / modification of both ferrn and tourr package.*

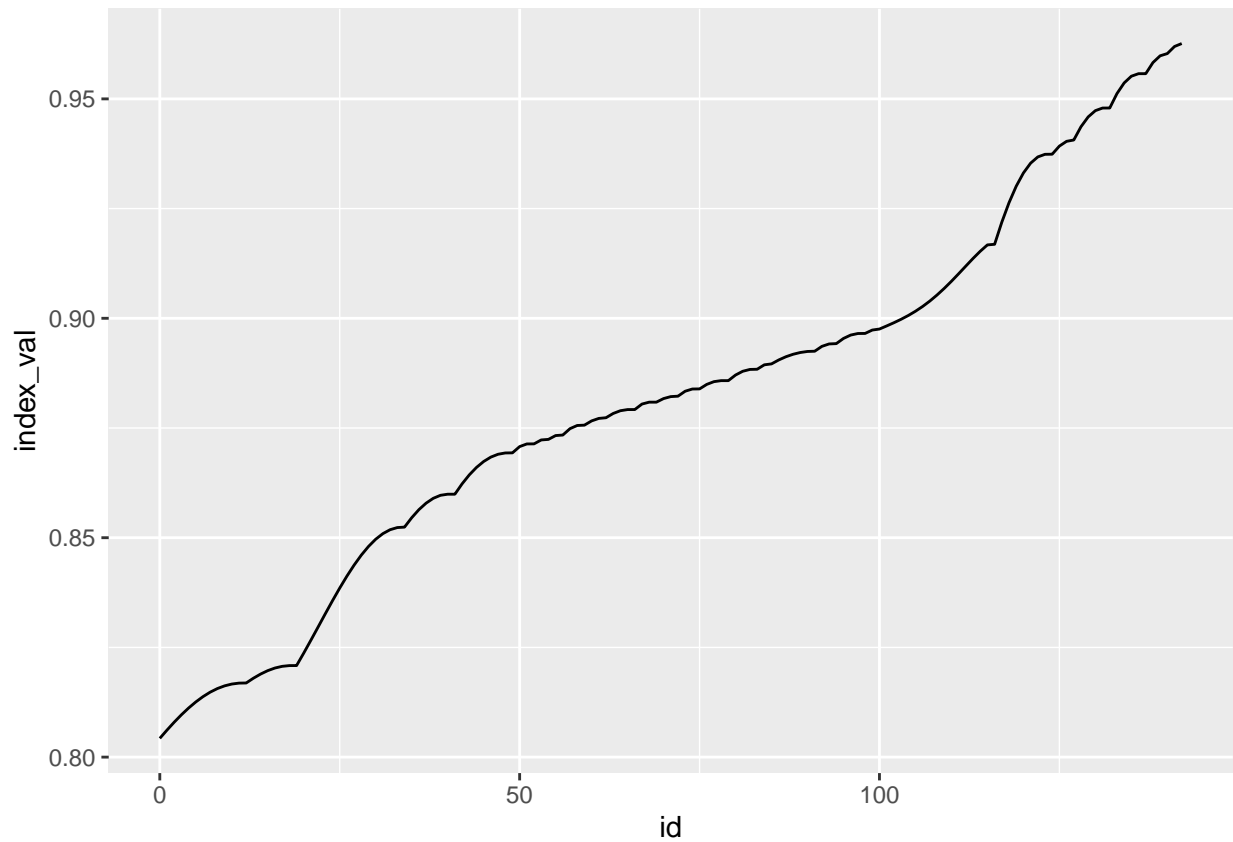
### 3.2 index vs time plots

The tracing plot is useful to explore the index value. Given a global object created, `explore_trace_all()` plots the index value of all the bases.

```
holes_2d_geo %>% explore_trace_all(magnify = TRUE)
```



```
holes_2d_geo %>% explore_trace_interp()
```

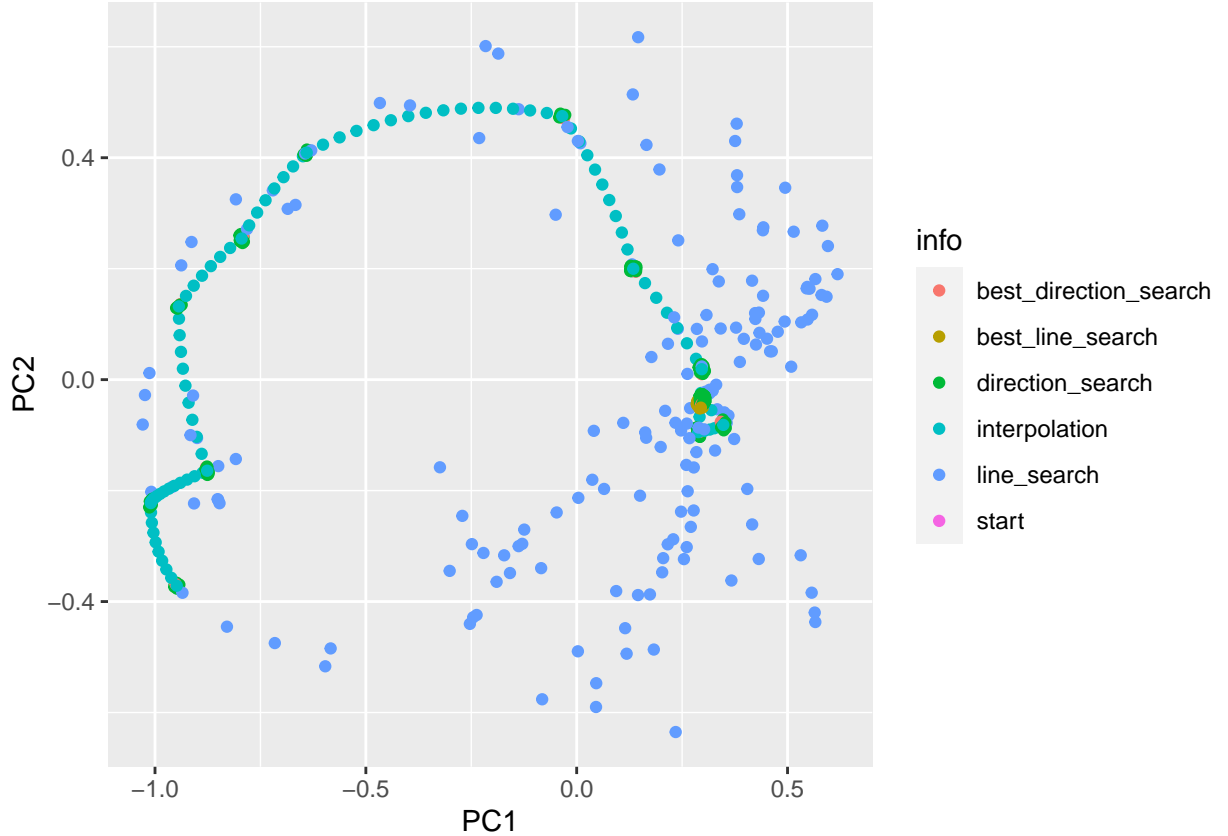


### 3.3 target basis vs interpolation steps, index vs time

### 3.4 low-d representation, PCA

```
holes_1d_geo %>% explore_proj_pca()
```

```
## Warning: The 'x' argument of 'as_tibble.matrix()' must have column names if '.name_repair' is not 'warn'.  
## Using compatibility '.name_repair'.  
## This warning is displayed once every 8 hours.  
## Call 'lifecycle::last_warnings()' to see where this warning was generated.
```



## 4 Animating plots

```
holes_1d_geo %>% explore_proj_pca(animate = TRUE)
```

## 5 Finding errors and developing improvements

The visualisation methods introduced above allow us to assess the performance of each optimisation routine and thus provide improvements.

### 5.1 Interrupt

The left panel of Figure @ref(fig:interruption) presents the trace plot of the interpolated basis using `search_better` to optimise the holes index (2D projection using the larger dataset). Tour interpolates from the current basis to the target basis, which has been

found to have a higher index value by projection pursuit. The target basis will then be sent to the projection pursuit as the current basis to find the next target basis. We can observe from the plot that there are basis with even higher index value on the interpolation path. These bases could be used to search for new basis in the next round. Therefore, an interruption is constructed to only take interpolated basis up to the point where the index value stops increasing and the last basis is taken as the current basis for the next round of search. After implementing the interruption, the tracing plot with the same configuration is shown on the right panel. Rather than interpolate to the target basis, the interruption interrupt the process after  $\text{id} = 60$  and proceeds the next round of search from the interpolated basis. Taking advantage of the interpolated basis results in a higher index value in the end with fewer steps.

Further the case when the index function is not smooth, the interpolation may not attain -> modification

## 5.2 Polish

In principle, all the optimisation routines should present the same output on the same problem. In Figure @ref(fig:trace-compare), red dots shows the trace of the interpolated basis using `search_geodesic` and `search_better` (with `max.tries = 100`), respectively, on the 2D projection problem. We can observe that they attain slightly different ending index value, which is not ideal. This motivates the creation of a polishing search that polishes the ending basis and achieves unity on different methods.

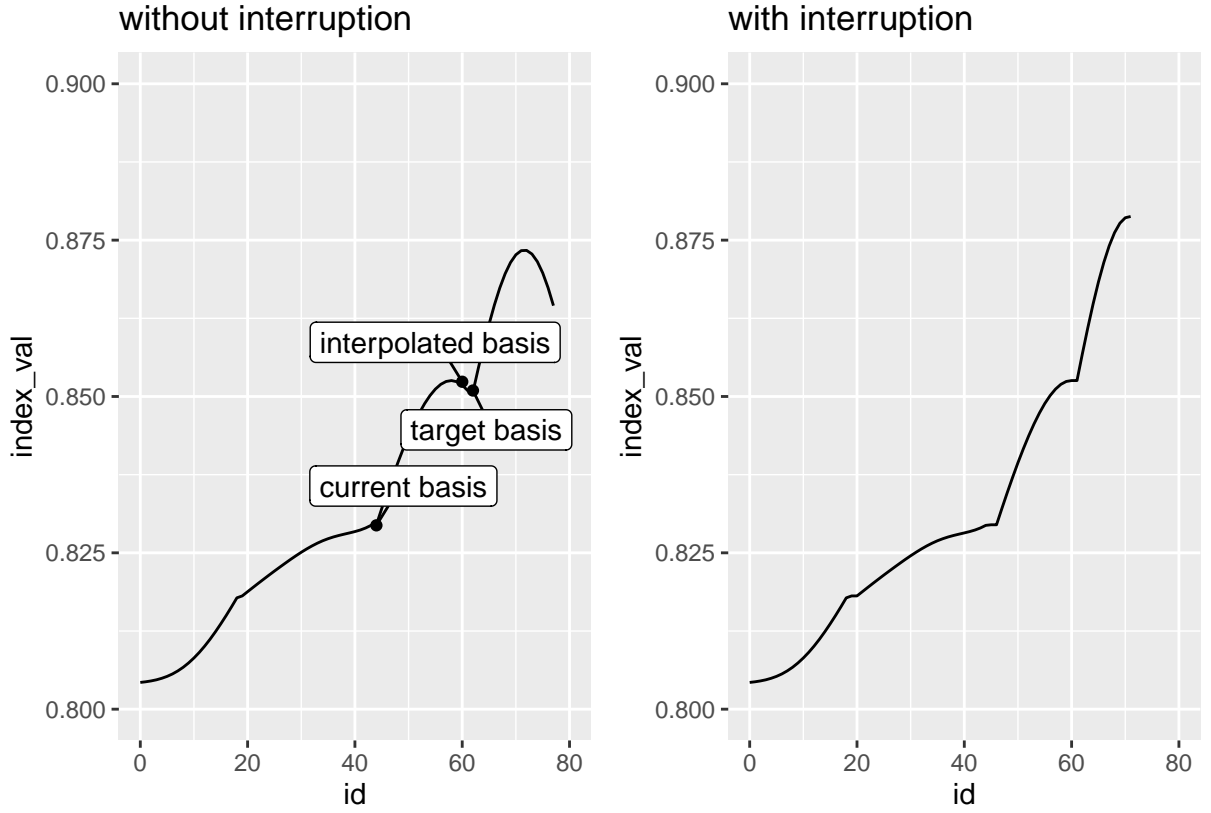
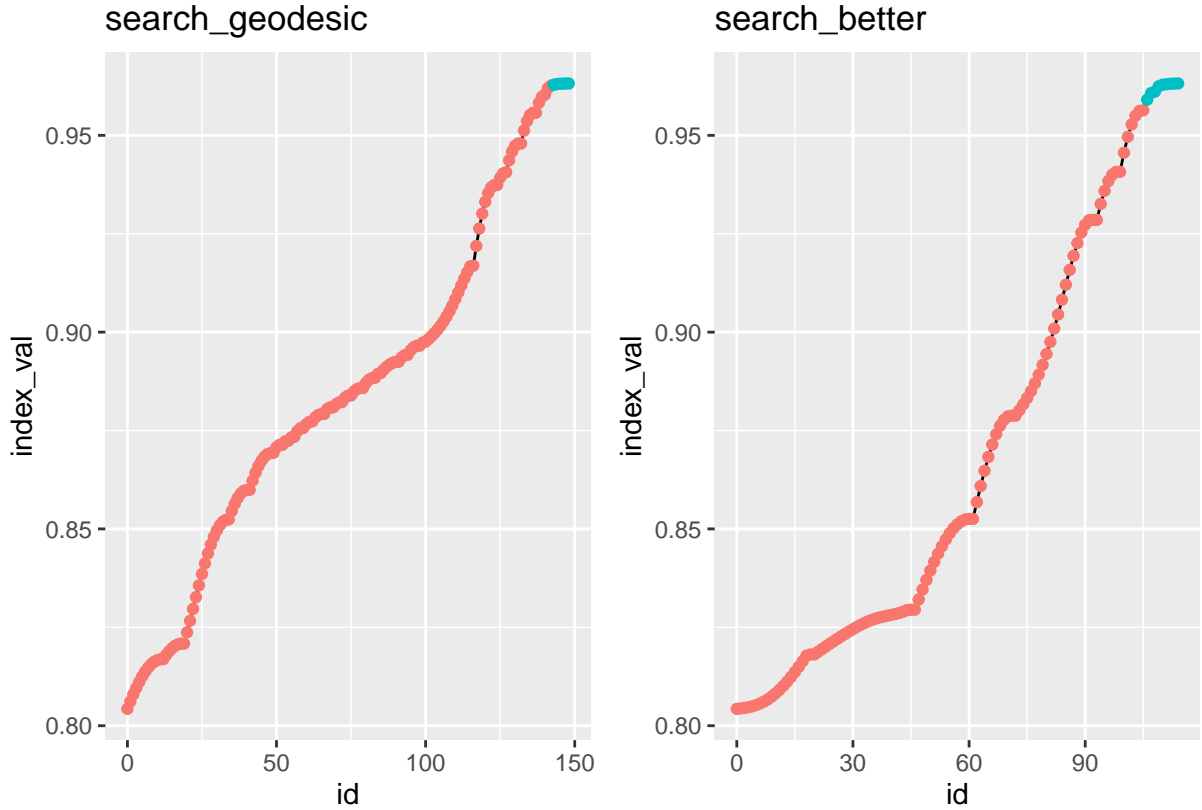


Figure 1: Trace plots of the interpolated basis with and without the interruption. The interruption stops the interpolation when the index value starts to decrease, see the difference at id being around 60. The implementation of the interruption finds an ending basis with higher index value using fewer steps.



`search_polish` takes the ending basis of a given search as the starting basis and uses a brutal-force approach to sample a larger number of basis (`n_sample`) in the neighbourhood, controlled by `polish_alpha`. Among the `n_sample` basis, the one with the largest index value becomes the candidate. If its index value is larger than that of the starting basis, it becomes the center of the searching neighbourhood in the next round. If no basis is found to have larger index value than the starting basis, the searching neighbourhood will shrink 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 improvement 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`

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)
```

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 @ref(fig:trace-compare), after polishing, both trials attain the same index value. However, a small `max.tries` of 25 is not sufficient for the algorithm to find the true maximum. This is because 25 tries is not sufficient for the 2D searching space.

### 5.3 posse?

`search_better` and `search_geodesic` methods are not doing good on 2d projection -> see the simple case above: abs best is 0.96 while they only get to 0.88-ish. Would be worse for higher dimension. thus `search_posse`. get 0.95 and then polish to 0.96 :)



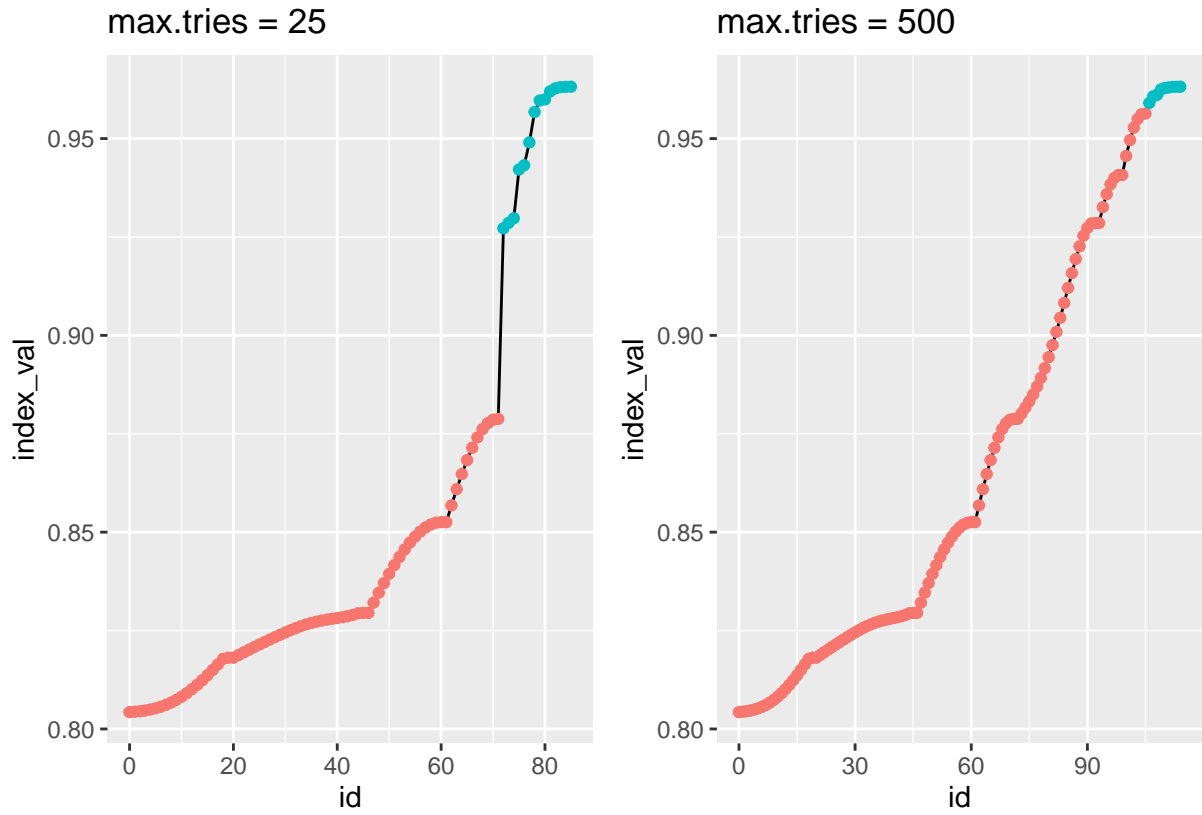


Figure 2: Breakdown of index value when using different max.tries in search better in conjunction with search polish. Both attain the same final index value after the polishing while using a max.tries 25 is not sufficient to find the true maximum.

# 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. (1997), ‘Manual controls for high-dimensional data projections’, *Journal of computational and Graphical Statistics* **6**(4), 464–480.
- 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.
- 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., 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/>