

index

For data $X_{n \times p}$, the projection basis $A_{p \times d}$ gives the projected data, $y = XA$. Projection pursuit finds the projection direction A that maximises the index function f :

$$\max f(XA) \quad A' A = 1.$$

For 2D projections, we can also write $y_{n \times 2} = (y_1, y_2)$, where y_1 and y_2 are the two columns of the projected data.

The following sections summarise the index functions used in the simulation:

hole index

- smooth
- formula:

$$I_{holes}(A) = \frac{1 - 1/n \sum_{i=1}^n \exp(-1/2 y_i y_i')}{1 - \exp(-p/2)}$$

- reference:
 - Cook and Swayne (2007) Interactive and Dynamic Graphics for Data Analysis page 30: [link](#)
 - code from the [tourr package](#)

dcor2d_2 index

- smooth
-
- Reference:
- Note

- From the function manual: “For bivariate data only, these are fast $O(n \log n)$ implementations of distance correlation and distance covariance statistics.”

MIC/ TIC index

- smooth
- formula
 - Let $g(k, l)$ define a partition of the space (y_1, y_2) into $k \times l$ rectangles, for example, $g(2, 3)$ means dividing the data space into 2 rectangles in the y_2 direction and 3 rectangles in the y_1 direction. Let G denotes all the possible partition.
 - MIC finds the maximum mutual information over G where k and l is bounded by the grid size: $k \times l < B(n) = n^\alpha$ where $\alpha = 0.3$ in our simulation

$$I_{MIC}(A) = \max_{g \in G} \frac{I(y_1, y_2 | g)}{\log(\min(k^*, l^*))} = \max_{g \in G} \frac{\sum_{y_1} \sum_{y_2} P(y_1, y_2) \log \frac{P(y_1)}{P(y_1)P(y_2)}}{\log(\min(k^*, l^*))}$$

where k^* and l^* are the number of rectangles in the optimal partition g .

- TIC calculates the sum of mutual information over G

$$I_{TIC}(A) = \sum_{g \in G} \frac{I(y_1, y_2 | g)}{\log(\min(k^*, l^*))} = \sum_{g \in G} \frac{\sum_{y_1} \sum_{y_2} P(y_1, y_2) \log \frac{P(y_1)}{P(y_1)P(y_2)}}{\log(\min(k^*, l^*))}$$

- Reference:
 - Reshef (2011) Detecting Novel Associations in Large Datasets: [link](#) - Figure 1 is usually to understand the general idea
 - Reshef (2016) Measuring Dependence Powerfully and Equitably: [link](#).
- Note:
 - The original Definition 1 on page 6 of Reshef (2016) uses supremum, here I just use the maximum.
 - There are different version of MIC and we use MIC_e and TIC_e (See Sec 4 in the paper in Reshef (2016)), which uses the equicharacteristic matrix rather than the original characteristic matrix in both MIC/ TIC calculation, see page 14 Figure 1 for there difference (equicharacteristic matrix parts the space equally hence faster).
 - Section 4.3/4.4 of Reshef (2016) talks about the time complexity

loess/ spline index

- smooth
- formula:
 - let $e_{y_1 \sim y_2}^{\text{model}}$ denote the residual from model $y_1 \sim y_2$
 - for the loess index, we fit a loess model to $y_1 \sim y_2$ (use $\alpha = 0.05$ for the smoothing parameter in the loess index)

$$I_{\text{loess}}(A) = \max \left(1 - \frac{\text{var}(e_{y_1 \sim y_2}^{\text{loess}})}{\text{var}(y_1)}, 1 - \frac{\text{var}(e_{y_2 \sim y_1}^{\text{loess}})}{\text{var}(y_2)} \right)$$

- for the spline index, a spline model is fitted (use cubic regression spline with $k = 15$ for the dimension of the basis in the spline model).

$$I_{\text{spline}}(A) = \max \left(1 - \frac{\text{var}(e_{y_1 \sim y_2}^{\text{spline}})}{\text{var}(y_1)}, 1 - \frac{\text{var}(e_{y_2 \sim y_1}^{\text{spline}})}{\text{var}(y_2)} \right)$$

- reference:
 - loess: (none) I create it myself, inspired by the spline index
 - spline:
 - * Ursula and Di (2020) Using tours to visually investigate properties of new projection pursuit indexes with application to problems in physics, page 1176: [link](#)
 - * code from the [cassowaryr](#) page

Stringy

- non-smooth
- formula:

$$I_{\text{stringy}}(A) = \frac{\text{number of vertices with 2 edges}}{\text{number of total vertices with more than one edge}}$$

(probably need some graph theory notation to write it mathematically)

- reference:
 - Wilkinson et al (2005) Graph-Theoretic Scagnostics: page 160: [link](#)
 - code from the [cassowaryr](#) page