

SHORT TECHNICAL NOTE

A new simplified manual tour, with application to exploring classification boundaries in high dimensions

Alex Aumann^a, German Valencia^a, Ursula Laa^b, Dianne Cook^c

^aSchool of Physics and Astronomy, Monash University; ^bInstitute of Statistics, University of Natural Resources and Life Sciences, Vienna; ^cDepartment of Econometrics and Business Statistics, Monash University

ARTICLE HISTORY

Compiled July 7, 2022

ABSTRACT

Something here

KEYWORDS

data visualisation; grand tour; statistical computing; statistical graphics; multivariate data; dynamic graphics

1. Introduction

From a statistical perspective it is rare to have data that are strictly 3D, and so unlike in most computer graphics applications, the more useful methods for data analysis show projections from an arbitrary dimensional space. These are dynamic data visualizations methods and are collected under the term *tours*. Tours involve views of high-dimensional (p) data with low-dimensional (d) projections. In his original paper on the grand tour, Asimov (1985) provided several algorithms for tour paths that could theoretically show the viewer the data *from all sides*. Prior to Asimov's work, there were numerous preparatory developments including Fisherkeller, Friedman, and Tukey (1974)'s PRIM-9. PRIM-9 had user-controlled rotations on coordinate axes, allowing one to manually tour through low-dimensional projections. (A video illustrating the capabilities is available through video library of ASA Statistical Graphics Section (2022).) Steering through all possible projections is impossible, unlike Asimov's tours which allows one to quickly see many, many different projections. After Asimov there have been many, many tour developments, which are summarized in Lee et al. (2021).

One such direction of work develops the ideas from PRIM-9, to provide manual control of a tour. Cook and Buja (1997) describe controls for 1D (or 2D) projections, respectively in a 2D (or 3D) manipulation space, allowing the user to select any variable axis, and rotate it into, or out of, or around the projection through horizontal, vertical, oblique, radial or angular changes in value. Spyrisson and Cook (2020) refined this algorithm and implements them to generate animation sequences.

CONTACT Alex Aumann. Email: aaum0002@student.monash.edu, German Valencia. Email: german.valencia@monash.edu, Ursula Laa. Email: ursula.laa@boku.ac.at, Dianne Cook. Email: dicook@monash.edu

Manual controls are especially useful for assessing sensitivity of structure to particular elements of the projection. There are many places where it is useful. In exploratory data analysis, where one sees clusters in a projection, can some variables be removed from the projection without affecting the clustering. For interpreting models, one can reduce or increase a variable's contribution to examine the variable importance. These controls can also be used to interactively generate faceted plots (?), or spatiotemporal glyphmaps (?). Having the user interact with a projection is extremely valuable for understanding high-dimensional data. However, these algorithms have two problems: (1) the pre-processing of creating a manipulation space overly complicates the algorithm, (2) extending to higher dimensional control is difficult.

Through experiments with the relatively new interactive graphics capabilities in mathematica(?), we have realized that there is a simpler approach, which is more direct, and extensible for generating user interaction. This paper explains this, and is organized as follows. The next section describes the new algorithm for manual control. This is followed by details on implementation. The applications section illustrate how these can be used.

2. What is a manual tour?

An orthonormal basis ($A_{p \times d}$) and a variable id ($m \in \{1, \dots, p\}$) to control are provided to initialize a manual tour. A method to update the values of the component of the controlled variable V_m is needed.

2.1. Existing methods

In the original work, the method for updating component values, for a 2D projection, was built trackball controls in 3D. A 3D manipulation space is created, as illustrated in Figure 1, where the controlled variable has full range of motion from -1 to 1. Movements of a cursor are recorded and converted into changes in the values of V_m to change it's values in the displayed 2D projection. Movement could also be constrained to be only in horizontal, vertical, radial or angular motions.

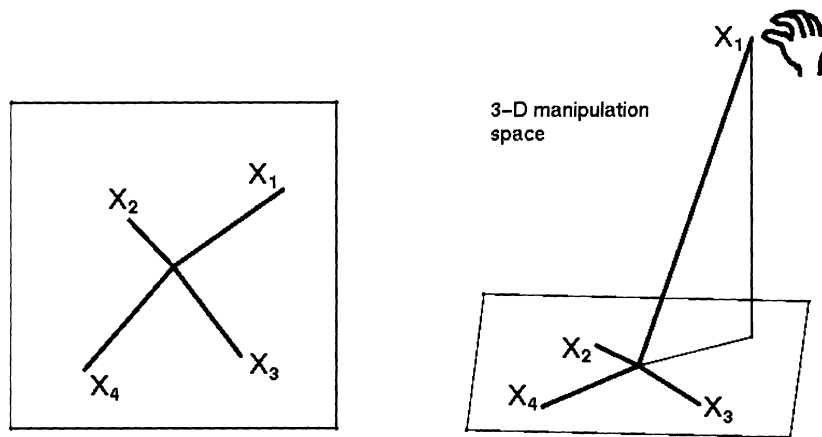


Figure 1. Original construction of the manual tour designed for 2D projections and created a 3D space from which to utilise track ball controls to change it's contribution. (Figure 3 from Cook and Buja (1997).)

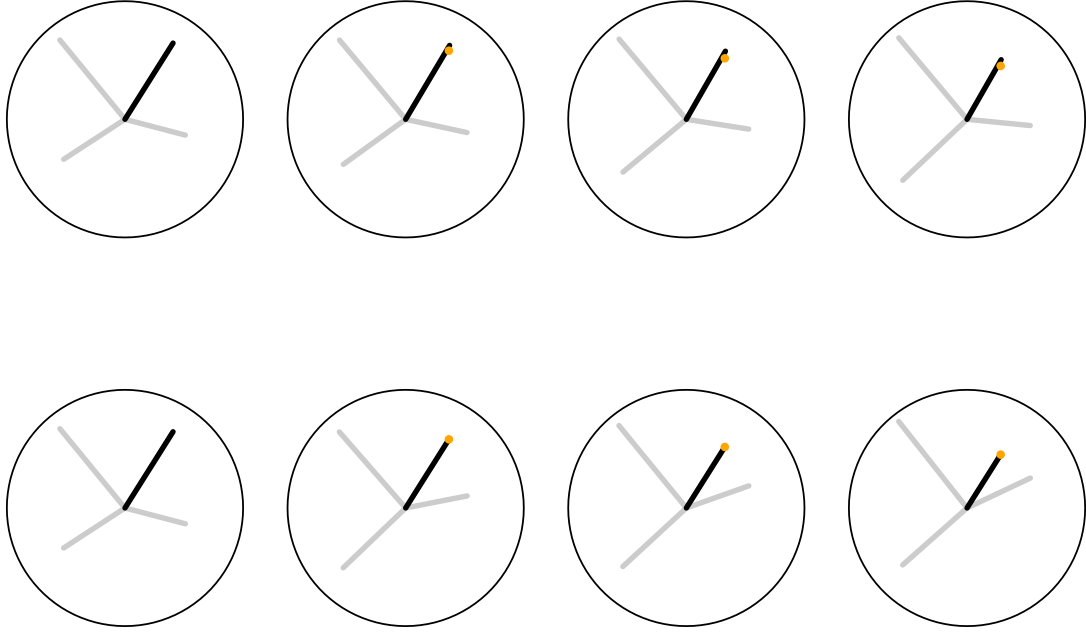


Figure 2. Sequence of projections where contribution of one variable is controlled (black) is changed: (top) unconstrained orthonormalisation, (bottom) constrained as specified. The dot (orange) indicates the chosen values for the controlled variable. For the constrained orthonormalisation it can be seen to precisely match the axis, but not so for the unconstrained orthonormalisation.

2.2. A new simpler and broadly applicable approach

The new approach emerged from experiments in mathematica. The components corresponding to V_m are directly controlled by cursor movement, which updates row m of A . The updated matrix is then orthonormalised.

2.2.1. Algorithm

1. Provide A , and m . (Note that m could also be automatically chosen as the component that is closest to the cursor position.)
2. Change values in row m , giving A^* . A large change in these values would correspond to making a large jump from the current projection. Small changes would correspond to tracking a cursor, making small jumps from the current projection.
3. Orthonormalise A^* , using Gram-Schmidt. For $d = 2$, and $A^* = [a_{.1} \ a_{.2}]$, the steps are:
 - i. Normalise $a_{.1}$, and $a_{.2}$.
 - ii. $a_{.2}^* = a_{.2} - a_{.1}^T a_{.2} a_{.1}$.
 - iii. Normalise $a_{.2}^*$.

This algorithm will produce the changes to a projection as illustrated in Figure 2 (top row). The controlled variable, V_m , corresponds to the black line, and sequential changes to row m of A can be seen to roughly follow a specified position (orange dot). Changes in the other components happen as a result of the orthonormalisation, but are uncontrolled.

2.3. Refinements to enforce exact position

The problem with new simple method (Algorithm 1) is that the precise values for V_m cannot be specified because the orthonormalisation will change them.

2.3.1. Adjustment method 1

A small modification to algorithm 1 will maintain the components of V_m precisely (Figure 2 (bottom row)). It is as follows:

1. Provide A , and m .
2. Change values in row m , giving A^* .
3. Store row m separately, and zero the values of row m in A^* , giving A^{*0} .
4. Orthonormalise A^{*0} , using Gram-Schmidt.
5. Replace row m with the original values, giving A^{**} .
6. For $d = 2$, adjust the values of a_{j2}^{**} using

$$a_{j2}^{**} + \frac{a_{m1}a_{m2}}{p-1}, j = 1, \dots, p, j \neq m$$

which ensures that

$$\sum_{j=1, j \neq m}^p a_{j1}^{**} a_{j2}^{**} + a_{m1}a_{m2} = 0$$

If $d > 2$ the process would be sequentially repeated in the same manner that Gram-Schmidt is applied sequentially to orthonormalise the columns of a matrix. If $d = 1$ no orthonormalisation is needed, and the projection vector would simply need to be normalized after each adjustment.

2.3.2. Adjustment method 2

For $d = 2$ projections, the projection matrix is the sub-matrix of A formed by its first two columns. Whereas orthonormality of the basis for the p -dimensional space is given by $e_i \cdot e_j = \delta_{ij}$, $i, j = 1, \dots, p$, orthonormality of the projection matrix is expressed as $P_i \cdot P_j = \delta_{ij}$, $i, j = 1, 2$. Movement of the cursor takes the two components x_{m1}, x_{m2} into a selected new value a, b . Although the motion is constrained by $a^2 + b^2 \leq 1$, this is not sufficient to guarantee orthonormality of the new projection matrix. One possible algorithm to achieve this is

- (1) Cursor movement takes $x_{m1}, x_{m2} \rightarrow a, b$
- (2) The freedom to change the components $A_{i3} \dots A_{ip}$ (the columns of A not corresponding to the projection matrix) is used to select a new orthonormal basis as follows:
 - (a) For row m one chooses $A_{m3} = \sqrt{1 - a^2 - b^2}$, $A_{m,k>3} = 0$
 - (b) For other rows, $A_{i \neq m, j \geq 3}$ random selections in the range $(-1, 1)$ are made.
 - (c) The Gram-Schmidt algorithm is then used to obtain an orthonormal basis taking e_m as the first vector (which is already normalized), and then

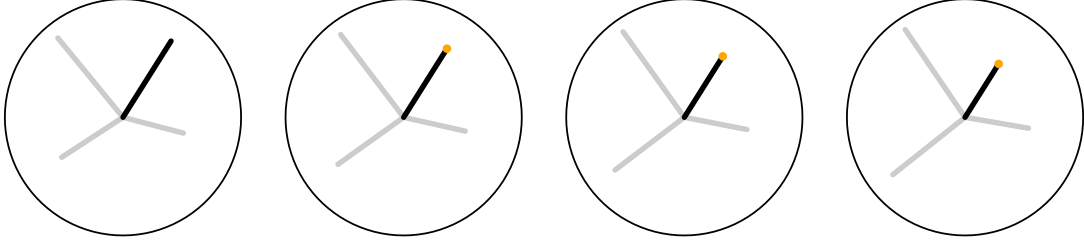


Figure 3. Manual controls with algorithm 2. The precise location of the axis is maintained.

proceeding as usual $e_1 \rightarrow e_1 - (e_1 \cdot e_m)e_m$, $e_1 \rightarrow e_1/(e_1 \cdot e_1)$, etc.

- (3) this results in the orthonormal basis A^* and a new projection matrix with $P_{m1} = a$, $P_{m2} = b$.

The random choice for the components of A not in the projection matrix allows the exploration of dimensions perpendicular to the projection plane.

$$A = \begin{matrix} & \begin{matrix} P_1 \\ \downarrow \\ P_2 \\ \downarrow \end{matrix} \\ \begin{matrix} e_1 \rightarrow \\ e_2 \rightarrow \\ \vdots \\ e_m \rightarrow \\ \vdots \\ e_p \rightarrow \end{matrix} & \begin{pmatrix} x_{11} & x_{12} & x_{13} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdot & \cdots & x_{2p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_{m1} & x_{m2} & \cdot & \cdots & x_{mp} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_{p1} & x_{p2} & \cdot & \cdots & x_{pp} \end{pmatrix} \end{matrix}, \quad (x_{m1}, x_{m2}) \rightarrow (a, b) \quad (1)$$

2.3.3. Adjustment method 3

For now just sketching the idea:

- first step is to click on the axis display to change the contribution of one variable m
- capture that position and replace the corresponding row in the projection matrix
- for 2D projection this gives components m_1 along “x” direction and m_2 along y direction
- next step: rotate basis such that direction of row m now corresponds to the first basis vector, this means we apply 2x2 rotation matrix to each row of the projection matrix, where the rotation angle is $\tan(\theta) = m_2/m_1$ (that matrix can be written just in terms of that ratio, $x = m_2/m_1$, when $m_1 < 0$ I think we just need to translate $\theta \rightarrow \theta + \pi$)
- take the rotated basis, apply Gram-Schmidt, now the direction of m will not change (but the length will change during normalization)
- rotated back to the original xy basis to update the plots

XXX need to check calculation of angle when $m_1 < 0$

3. Experimenting with new techniques using Mathematica

- Why is this a good sandbox

- **Explain the functionality available in the notebooks**

The following functions are implemented in the Mathematica package `mmtour.wl`:

- `ProjectionPlot[data, projmat, xRange(=Automatic), yRange(=Automatic), colour(=Automatic)]` with arguments
 - `data`: A list of data matrices (they must all have the same dimensionality)
 - `projmat`: A projection matrix which describes the projection plane. This is the initial projection and does not need to be orthonormal (XXX is this right?).
 - `xRange`: A list which details the range of the horizontal axis.
 - `yRange`: A list which details the range of the vertical axis.
 - `colour`: A list of graphics directives to be applied to each data matrix.

There are problems with the default values for x and y ranges and colour. It would be better to fix the ranges and colour based on the data as well as the origin of axes and the size of the points.

It would be good if the initial projection matrix (and maybe the center) had default values. Perhaps a random matrix with the correct dimensions read from the data? Fix 1

The format for the printed projection matrix can be improved also. Maybe fix the number of digits and the position of the label.

- `SlicePlot[data, projmat, centrePoint, height, heightRange, ptSize1(=0.005), ptSize2(=0.004), minDist(=0)]` with arguments (used to be called `SliceDynamic` in Alex's first version, the name `SliceDynamic` was taken by something else below)
 - `data`: A data matrix
 - `projmat`: The initial projection matrix
 - `centrepoint`: The initial position of the slice
 - `height`: The initial height of the slice
 - `heightRange`: The range of heights
 - `ptSize1`: The size of the points that exist within the slice
 - `ptSize2`: The size of the points that exist outside of the slice

This function requires some auxiliary functions. It does not accept more than one data set.

It would be good if the initial projection matrix (and maybe the center) had default values. Perhaps a random matrix with the correct dimensions read from the data?

It may be good to fix the size of the points, the center of the axes and the scales of the axes rather than allow them to change as the projection and slicing changes

- `SliceDynamic[data, projmat, centrePoint, height, heightRange, ptSize1 : 0.005, ptSize2 : 0.004]` with arguments
 - `data`: A data matrix with an extra grouping or clustering column (numerical)
 - `projmat`: The initial projection matrix
 - `centrepoint`: The initial position of the slice
 - `height`: The initial height of the slice
 - `heightRange`: The range of heights
 - `ptSize1`: The size of the points that exist within the slice
 - `ptSize2`: The size of the points that exist outside of the slice

Does this need to be a separate function from SlicePlot? Fix 2

It would be good if the initial projection matrix (and maybe the center) had default values. Perhaps a random matrix with the correct dimensions read from the data? Fix 3

It may be good to fix the size of the points, the center of the axes and the scales of the axes rather than allow them to change as the projection and slicing changes. Fix 4

I have found it useful to input different data sets rather than have this function sort them automatically using the last column. One data set with a flag column at the end, is best! Is the code too complicated? Fix 5

- Projected2DSliderPlot[]

This function uses a different display for the projection that I find very useful. However as it stands it is specific to one data set. Needs to be generalised to arbitrary data if we want to keep it.

For Alex to do

- We would need a screenshot (and video) for each of these.
- Fix 1, 2, 3, 4, 5
- Code consistent for applications to work - see slice_tour_compare_sets.nb. See below
 - We need to make sure the drawing range is fixed (no resizing of the axis when changing projection or slicing). It could be nice if there was a manual zoom in/zoom out to change scale occasionally.
 - Should also make sure the size of the points is fixed (maybe this will be automatic if the drawing range is fixed?), bigger points would be preferred, currently the size is sometimes really tiny and it becomes difficult to see anything.
 - Can we have some heuristic to decide the range for the slice thickness slider? In the second graph there is a tiny part of the slider where the changes are interesting, difficult to navigate since this makes it super sensitive to any small movement, while most of the options are not interesting. My suggestion would be to keep (as I understand is the case now) the maximum thickness (all points in) as the upper bound, but instead of zero use a thickness below which slices tend to be empty as the lower bound.
 - For the olives example I think the relation between x2 and the blue region could be interesting.
 - pdfsense data: can we use random samples instead of the first 10k rows? Then we should still be able to see the bigger picture (with picking the first ones we get some artificial structure from how the data was sampled initially).
 - I was playing around with this today. Looks like the axes are fixed when changing slice thickness, but not when rotating?
 - Anyways, some observations that might be useful, but all of this seems difficult for including in a paper:
 - In the first example we can show how the yellow group is associated with the first variable (x1, I guess this is C9?) - when slicing through the center it only appears when x1 has a bigger com-

ponent on the projection, it is associated with low values of x_1 . Since they are not near the average value of x_1 the points do not get captured in a slice where x_1 is not important in the projection. We can move the center point to -1 in the first component to then explore how the yellow group relates to the other three variables. Much less clear, maybe some indication of smaller x_2 and larger x_3 for this.

- The second example is too busy for manual exploration I would say (too many axes to move around separately).
- Third example could be similar to the first: blue region requires small values of x_2 , can be understood from the exploration, then move center point to -0.5 in x_2 to explore further. Now interpretation becomes tricky. Sometimes blue and yellow seem just randomly everywhere, sometimes blue more contained in an ellipsoid around the center (indicating some correlation between variables). For example try zeroing out components 2 and 3 to see this slice in x_1 vs x_4 .
- For the last example x_4 seems to be the one to work with. I feel like this is again very similar, maybe better to instead explore other examples (maybe something without slicing?)

4. Applications

To illustrate the usefulness of the manual tour use the penguins data (Horst, Hill, and Gorman 2020). We will show how classification boundaries can be explored and better understood using a manual tour, on projections and slices through 4D space.

5. Extension added in the R package `tourr`

Explain `radial_tour`

6. What would be desirable for implementations in R?

7. Discussion

Acknowledgements

The authors gratefully acknowledge the support of the Australian Research Council. The paper was written in `rmarkdown` (Xie, Allaire, and Golemund 2018) using `knitr` (Xie 2015).

Supplementary material

The source material and animated gifs for this paper are available at

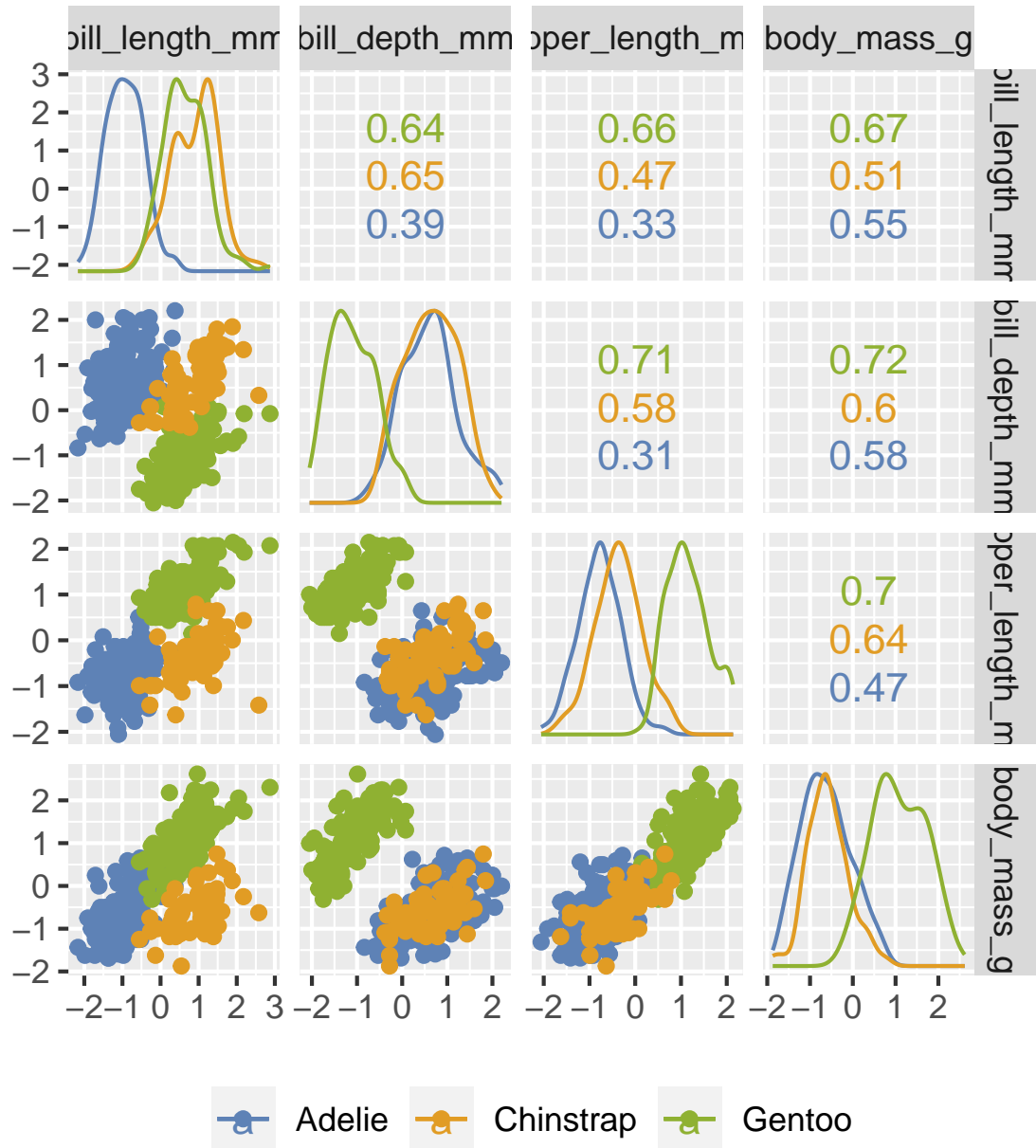


Figure 4. Scatterplot matrix of the (standardised) penguins data.

References

- ASA Statistical Graphics Section. 2022. “Video Library.” <https://community.amstat.org/jointscsg-section/media/videos>.
- Asimov, D. 1985. “The Grand Tour: A Tool for Viewing Multidimensional Data.” *SIAM Journal of Scientific and Statistical Computing* 6 (1): 128–143.
- Cook, Dianne, and Andreas Buja. 1997. “Manual Controls for High-Dimensional Data Projections.” *Journal of Computational and Graphical Statistics* 6 (4): 464–480. <http://www.jstor.org/stable/1390747>.
- Fisherheller, M.A., J.H. Friedman, and J.W. Tukey. 1974. “PRIM-9, an interactive multidimensional data display and analysis system.” In *The Collected Works of John W. Tukey: Graphics 1965-1985, Volume V*, edited by William S. Cleveland, 340–346.
- Horst, Allison Marie, Alison Presmanes Hill, and Kristen B Gorman. 2020. *palmerpenguins: Palmer Archipelago (Antarctica) penguin data*. R package version 0.1.0, <https://allisonhorst.github.io/palmerpenguins/>.
- Lee, Stuart, Dianne Cook, Natalia da Silva, Ursula Laa, Earo Wang, Nick Spyrison, and H. Sherry Zhang. 2021. “Advanced Review: The State-of-the-Art on Tours for Dynamic Visualization of High-dimensional Data.” *arXiv:2104.08016 [cs, stat]* <http://arxiv.org/abs/2104.08016>.
- Spyrison, Nicholas, and Dianne Cook. 2020. “spinifex: an R Package for Creating a Manual Tour of Low-dimensional Projections of Multivariate Data.” *The R Journal* 12 (1): 243. <https://journal.r-project.org/archive/2020/RJ-2020-027/index.html>.
- Xie, Yihui. 2015. *Dynamic Documents with R and knitr*. 2nd ed. Boca Raton, Florida: Chapman and Hall/CRC. <https://yihui.name/knitr/>.
- Xie, Yihui, Joseph J. Allaire, and Garrett Golemund. 2018. *R Markdown: The Definitive Guide*. Chapman and Hall/CRC. <https://bookdown.org/yihui/rmarkdown>.