

Discussion of: Visualizing statistical models: Removing the blindfold

Catherine B. Hurley
Department of Mathematics & Statistics
Maynooth University, Ireland
catherine.hurley@nuim.ie

April 28, 2015

Collections, comparisons, graphs and seriation

I enjoyed reading the paper “Visualizing Statistical Models” and thank the authors Wickham, Cook and Hofmann for an informative and stimulating contribution. The notion of using visualization techniques on model results is important as illustrated in the case studies. For those of us interested in taking on the challenge of developing “visualizations for more models in a wider array of statistical and graphical environments”, the current paper (VSM) presents good strategies and a useful distinction between visualising the data in model space and the model in data space. The original work presented here dates from 2007, yet the ideas are fresh.

A lot of careful thought goes in to choices of a visualization, and then its design and implementation. When the presentation of ideas and visualizations are so clear as in the current paper, it all seems so easy, but it is in fact far from easy. The authors hint at this when they say “converting data-vis to model-vis is not always straightforward”. Some new techniques are required, to find prediction boundaries as in Figure 5 of VSM, and to display dendrograms in feature space as in Figure 9. I looked forward to trying out the R packages `classifly` and `clusterfly` but they require `rggobi` which disappointingly no longer has a Mac version. `Meifly` installs but does not implement the techniques of Section 4.1. I think these packages have a lot to offer the statistical community and I’d be delighted if they were updated to work on current hardware, with code for the visualizations of the current paper included as a vignette perhaps.

Collections and comparisons

A theme of this paper is visualising collections of models and iterative model fits. Perhaps there is a need for tools to help organise these collections, that will facilitate comparisons. In their PairViz methodology Hurley and Oldford (2011a,b) proposed the use of a mathematical graph to organise collections of statistical objects, where edges in this graph represent comparisons of interest. Applying this principle to the collections of models proposed in Section 4 VSM, we have graphs whose nodes are models and edges represent model comparisons. For example, to explore the space of all possible regression models we could build a graph where each node represents a regression model, and edges connect models that differ by exactly one predictor.

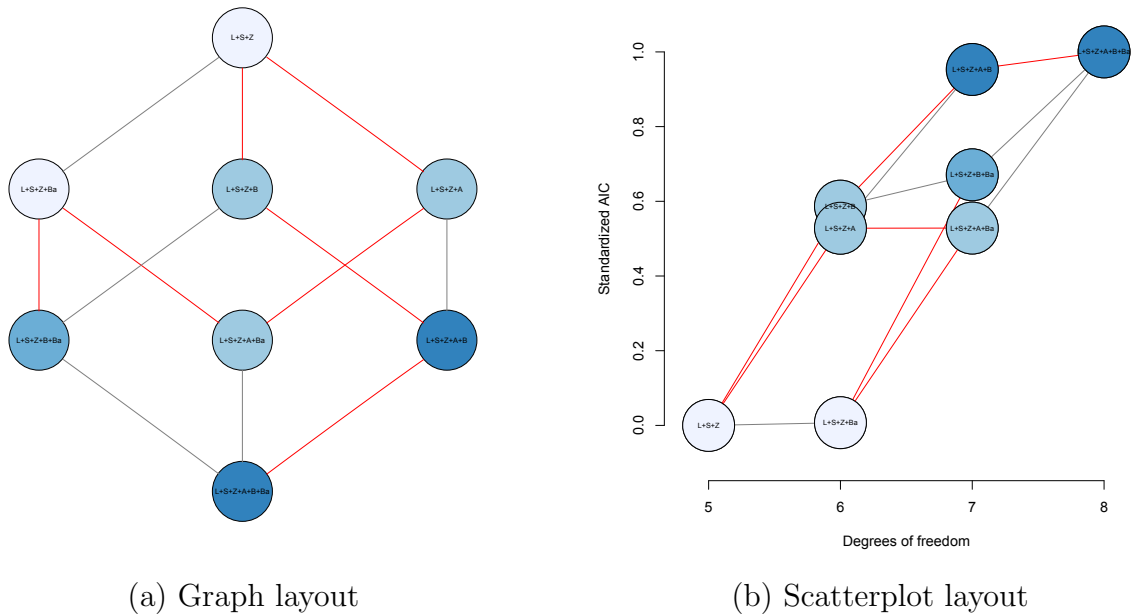


Figure I: Top 8 models according to AIC. Edges connect models that differ by exactly one predictor.

Figure I shows such graphs for the New Haven Housing Data, as in Section 4 of VSM. Only the best 8 models according to AIC are shown. The left hand side figure shows a conventional graph layout, the right hand side shows node position given by (standardised) AIC and degrees of freedom, similar to Figure 11. Paths through the graph such as the one shown in red could be used as a basis for comparing summary measures for collections of models. We could then easily see whether adding or removing a predictor changes other other coefficients, far more effectively than pouring over pages of tables. For such displays to be truly useful they could be supplemented by interactive tools for filtering nodes and edges. Interactive tools for selecting paths through the graph could then be used to drive model comparison visualizations. In their RnavGraph package Waddell and Oldford (2014) (see also Oldford and Wadell (2011)) used similar ideas to drive exploration of high-dimensional data.

Role of seriation

Paths through graphs could be found algorithmically, using eulerian paths that visit all edges or hamiltonian paths that visit all nodes Hurley and Oldford (2011a). Seriation techniques (see for example Earle and Hurley (2015); Hurley and Earle (2013); Hahsler et al. (2008)) produce paths visiting all nodes. Such paths may be chosen specifically to improve comparison of statistical objects (cases, variables and models) through better visualizations.

It is interesting to consider how seriation techniques might be used to improve some of the visualizations presented by the authors. The dendrogram of Figure 8 of VSM is an obvious candidate. The three clusters shown correspond roughly to the three wine classes with just a small number of wines misclassified. This visualization shows the data (leaves) in model space, but potentially in a misleading way. A few of the variety A wines shown in purple are placed at the far end of the predominantly B variety green cluster, giving the impression that these variety A wines are at the edge of the green cluster and are extreme outliers relative to their class. A number of the variety B (green) wines appear within the predominantly C variety orange cluster. Again their position gives the impression that these variety B wines are extreme outliers relative to their class. But, in fact the default arrangement of dendrogram leaves is somewhat arbitrary and these wines may not be extreme outliers. Seriation improves the dendrogram and reduces its potential to mislead. In Figure II we use the DendSer seriation algorithm (Hurley and

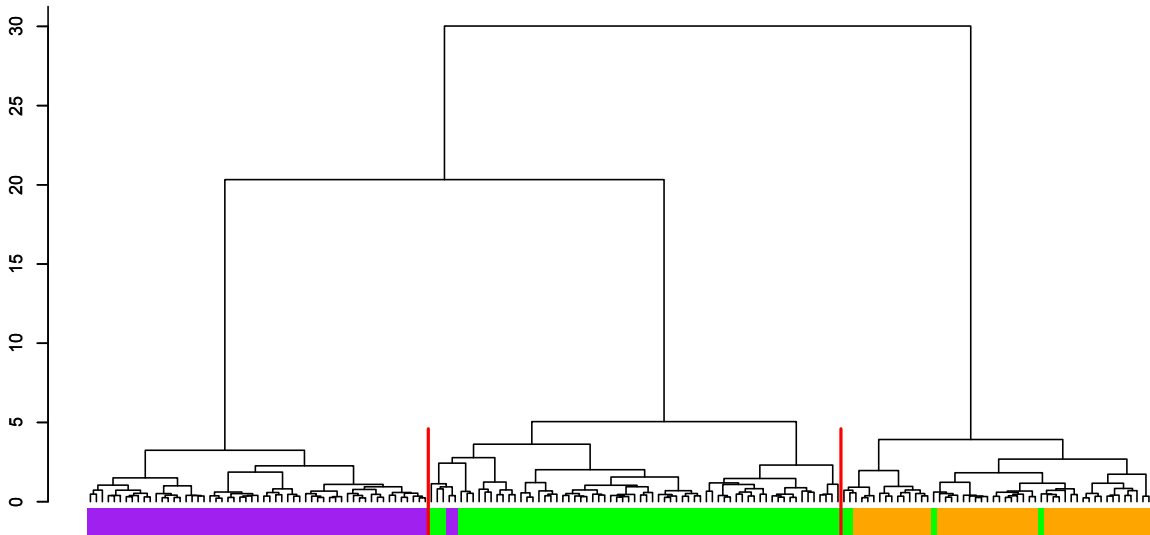


Figure II: Top 8 models according to AIC. Edges connect models that differ by exactly one predictor.

Earle 2013; Earle and Hurley 2015) to re-arrange the leaves of the dendrogram so that leaves close in the dendrogram are close in feature space (as measured by the distance

matrix). The misclassified variety A wines now appear near the cluster with the other variety A wines. Similarly, two of the misclassified variety B wines are moved adjacent to the other B variety wines. This dendrogram gives the impression that perhaps two variety B wines might be outliers relative to their class.

We can check this in feature space. The displays in Figures 9 and 17 of VSM show some overlap between the clusters and that misclassified points are not extreme outliers. Incidentally, as the hierarchy of joins is not evident in Figure 9 even with close scrutiny, I'm not convinced that plotting dendrograms in feature space is a useful m-in-ds visualization.

Enhance teaching

I strongly agree with the authors on the value of model visualization in teaching. As a profession we could make much better use of (preferably interactive) graphics in teaching, even to explain and explore concepts such as tables of numeric regression output and Anova tables. Simple tools such as nomograms allow students to play around with models which should deepen their understanding. The DynNom R package (Jalali et al. 2015) is a nice nomogram implementation based on Shiny (Chang et al. 2015) for general linear models.

I would welcome more and better tools for interactive visualization of models. The diverse case studies and methodologies used by Wickham, Cook and Hofmann present a convincing argument for the value of such tools, and should serve as motivation for those of us hoping to advance research in this area.

References

- Chang, W., Cheng, J., Allaire, J., Xie, Y., and McPherson, J. (2015), *shiny: Web Application Framework for R*, r package version 0.11.1.
- Earle, D. and Hurley, C. B. (2015), "Advances in Dendrogram Seriation for Application to Visualization," *Journal of Computational and Graphical Statistics*, 24, 1–25.
- Hahsler, M., Hornik, K., and Buchta, C. (2008), "Getting Things in Order: An Introduction to the R Package seriation," *Journal of Statistical Software*, 25, 1–34.
- Hurley, C. and Oldford, R. (2011a), "Eulerian tour algorithms for data visualization and the PairViz package," *Computational Statistics*, 26, 613–633.
- (2011b), *PairViz: Visualization using Eulerian tours and Hamiltonian decompositions*, r package version 1.2.1.
- Hurley, C. B. and Earle, D. (2013), "DendSer: Dendrogram seriation," R package version 1.0.1.

- Jalali, A., Alvarez-Iglesias, A., and Newell, J. (2015), *DynNom: A Dynamic Nomogram for Linear and Generalized Linear Models as Shiny Applications*, r package version 1.0.1.
- Oldford, R. and Wadell, A. (2011), “Visual Clustering of High-dimensional Data by Navigating Low-dimensional Spaces,” in *Proceedings of ISI*.
- Waddell, A. R. and Oldford, R. W. (2014), *RnavGraph: Using Graphs as a Navigational Infrastructure*, r package version 0.1.8.