# molic: An R package for multivariate outlier detection in contingency tables

16 August 2019

## Summary

Outlier detection is an important task in statistical analyses. An outlier is a case-specific unit since it may be interpreted as natural extreme noise in some applications, whereas in other applications it may be the most interesting observation. The **molic** package have been written to facilitate the novel outlier detection method in high-dimensional contingency tables (Lindskou, Svante Eriksen, and Tvedebrink 2019). In other words, the method works for data sets in which all variables are *categorical*, implying that they can only take on a finite set of values (also called *levels*).

The software uses decomposable graphical models (DGMs), where the probability mass function can be associated with an interaction graph, from which conditional independences among the variables can be inferred. This gives a way to investigate the underlying nature of outliers. This is also called *understandability* in the literature. Outlier detection has many applications including areas such as

- Fraud detection
- Medical and public health
- Anomaly detection in text data
- Fault detection (on critical systems)
- Forensic Science

## The Method

The method can be described by the **outlier test** procedure below. Assume we are interested in whether or not a new observation $z$ is an outlier in some data set $D$. First an *interaction graph $G$* is fitted to the variables in $D$; a decomposable undirected graph that describes the association structure between variables in $D$. If the assumption that $z$ belongs to $D$ is true, $z$ should be included in $D$. Denote

by $D_z$ the new data set including $z$. Finally the outlier model $M$ is constructed using $G$ and $D_z$ from which we can query the p-value, $p$, for the test about $z$ belonging to $D$. If $p$ is less than some chosen threshold (significance level), say 0.05, $z$ is declared an outlier in $D$.

```
1: outlier test (D: data, z : new obs.)
2:    G := fit_graph(D)
3:    D_z := D with z included
4:    M := fit_outlier(D_ z, G)
5:    p := pval(M, deviance(M, z))
6:    return p
7: end outlier test
```

The `fit_graph` algorithm has three ways of fitting a graph. The `fwd` type is an implementation of the efficient step-wise selection procedure (Deshpande, Garofalakis, and Jordan 2001) used for model selection in decomposable graphs. There is also a backward, `bwd`, type and finally it is also possible to fit a tree interactions graph (i.e. only first order associations).

The `fit_graph` function can be used to explore dependencies between any kind of discrete variables and make statements about conditional dependencies and independencies. A thorough description of the outlier detection method and how to use the software can be found at

https://mlindsk.github.io/molic/

# Expert Knowledge

If one has prior knowledge of the underlying nature of the association between variables, this can easily be exploited. One can choose to model only the relationship between variables which have no other associations to any of the remaining variables. This will result in a number of interaction graphs which can then be unified as the union of these graphs. This approach was taken in the example below.

# A Use Case in Forensic Science

Recently, advances in DNA sequencing has made it possible to sequence short segments of DNA ($< 200$ basepairs) including two or more SNPs. These are called *microhaplotypes* (or microhaps for short) (K. K. Kidd et al. 2014). They have been demonstrated to be well suited for ancestry assessment in the forensic science community. The short distance between SNPs within a microhap implies that recombination among them rarely occurs. Hence, the methodology of T. Tvedebrink et al. (2018) can not be used as this assumes mutual independence of the SNPs within a population (corresponding the null graph with no edges).

In Lindskou, Svante Eriksen, and Tvedebrink (2019) the **molic** package was used to detect outliers in microhap data from the 1000 Genomes Project (The 1000 Genomes Project Consortium 2015). This data contains DNA profiles from five different continental regions (CRs); Europe (EUR), America (AMR), East Asia (EAS), South Asia (SAS) and Africa (AFR).

Consider for example the region SAS as the hypothesized region and all profiles in AFR as profiles to be tested against the hypothesis that their origin is SAS. Two different interaction graphs are used; $G$ which is the result of using the `efs` algorithm and $G^\emptyset$ where all microhap SNPs are assumed to be independent (a graph with no edges). The proportion of profiles from AFR that are outliers in SAS according to the model, is 1 for $G$ and only 0.834 for $G^\emptyset$, see Table 1. The outlier test was conducted for all pairs of continental regions. It is seen, that $G$ outperforms $G^\emptyset$ in general and the dependency between microhap SNPs cannot be neglected. All tests was conducted on a significance level of 0.05.

| $H_0$ / $z$ | EUR | EAS | AMR | SAS | AFR |
|---|---|---|---|---|---|
| EUR | 0.054 | 1 | 0.191 | 0.509 | 1 |
|     | 0.046 | 1 | 0.145 | 0.231 | 1 |
| EAS | 1 | 0.054 | 0.994 | 0.966 | 1 |
|     | 1 | 0.063 | 0.994 | 0.980 | 1 |
| AMR | 0.778 | 1 | 0.095 | 0.769 | 1 |
|     | 0.697 | 1 | 0.049 | 0.565 | 1 |
| SAS | 0.922 | 1 | 0.710 | 0.037 | 1 |
|     | 0.863 | 1 | 0.620 | 0.047 | 1 |
| AFR | 1 | 1 | 0.997 | 1 | 0.101 |
|     | 0.998 | 1 | 0.918 | 0.834 | 0.106 |

$\blacksquare$ $G$  $\square$ $G^\emptyset$

Table 1: Performance matrix of outlier tests.

Consider the saturated model (a complete graph). This is the equivalent of estimating probabilities using the naive frequency counts in the data. For one, it does not (necessarily) capture the biological association between SNPs and second it would, in general, require an enormous amount of data to obtain valid estimates.

# References

Deshpande, Amol, Minos Garofalakis, and Michael I Jordan. 2001. "Efficient Stepwise Selection in Decomposable Models." In *Proceedings of the Seventeenth Conference in Uncertainty in Artificial Intelligence*, 128–35. Morgan Kaufmann

Publishers Inc.

Kidd, Kenneth K., Andrew J. Pakstis, William C. Speed, Robert Lagacé, Joseph Chang, Sharon Wootton, Eva Haigh, and Judith R. Kidd. 2014. "Current Sequencing Technology Makes Microhaplotypes a Powerful New Type of Genetic Marker for Forensics." *Forensic Science International: Genetics* 12: 215–24. doi:10.1016/j.fsigen.2014.06.014.

Lindskou, Mads, Poul Svante Eriksen, and Torben Tvedebrink. 2019. "Outlier Detection in Contingency Tables Using Decomposable Graphical Models." *Scandinavian Journal of Statistics*. Wiley Online Library. doi:10.1111/sjos.12407.

The 1000 Genomes Project Consortium. 2015. "A Global Reference for Human Genetic Variation." *Nature* 526 (7571). Nature Publishing Group: 68. doi:10.1038/nature15393.

Tvedebrink, T, P S Eriksen, H S Mogensen, and N Morling. 2018. "Weight of the Evidence of Genetic Investigations of Ancestry Informative Markers." *Theoretical Population Biology* 120: 1–10. doi:10.1016/j.tpb.2017.12.004.