

A modified version of the Procrustes correlation coefficient for high-dimensional data

Received on XXXXX; revised on XXXXX; accepted on XXXXX

Supplementary information: Supplementary data are available at *Bioinformatics* online.

2 Approach

RLs is part of the procruste framework that aims to superimpose a set of points with respect to another through three operations: a translation, a rotation and a scaling. The optimal transfert matrix $Rot_{X \rightarrow Y}$ can be estimated from the singular value decomposition (SVD) of the covariance matrix $X'Y$. SVD factorize any matrix as the product of three matrices (Equation 3).

$$X'Y = U\Sigma V' \quad (3)$$

U and V are two rotation matrices allowing to compute the transfer matrix to superimpose X on Y or reciprocally Y on X . All the elements of Σ except its diagonal are equal to zero. The diagonal elements are the singular values. Singular values are the extension of the eigenvalues to non square matrices. $CovLs(X, Y)$ can also be computed from singular values (Equation 4).

$$CovLs(X, Y) = \frac{\text{Trace}(\Sigma)}{n - 1} \quad (4)$$

This expression illustrates that actually $CovLs(X, Y)$ is the variance of the projections of X on Y or of the reciprocal projection. Therefore $CovLs(X, Y)$ and $Rls(X, Y)$ are always positive and rotation independante. Here we propose to partitionate this variance in two components. A first one corresponding to the actual shared information between X and Y , and a second part that corresponds to that two random matrices of same structure than X and Y are sharing. Two methods are proposed to estimate $\overline{RCovLs(X, Y)}$ the mean the random part of $CovLs(X, Y)$. $ICovLs(X, Y)$ the informative part of $CovLs(X, Y)$ is estimated using Equation (5)

$$ICovLs(X, Y) = \text{Max} \begin{cases} CovLs(X, Y) - \overline{RCovLs(X, Y)} \\ 0 \end{cases} \quad (5)$$

Similarly the informative counter-part of $VarLs(X)$ is defined as $IVarLs(X) = ICovLs(X, X)$.

3 Methods

3.1 Monte-Carlo estimation of $\overline{RCovLs(X, Y)}$

For every values of p and q including 1, $\overline{RCovLs(X, Y)}$ can be estimated using a serie of k random matrices $RX = \{RX_1, RX_2, \dots, RX_k\}$ and $RY = \{RY_1, RY_2, \dots, RY_k\}$ where each RX_i and RY_i have the same structure respectively than X and Y in term of number of columns and of standard deviation of these columns.

$$\overline{RCovLs(X, Y)} = \frac{\sum_{i=1}^k CovLs(RX_i, RY_i)}{k} \quad (6)$$

Even when $X = Y$ to estimate $VarLs(X)$, $\overline{RCovLs(X, Y)}$ is estimated with two independent sets of random matrix RX and RY , both having the same structure than X .

3.2 Empirical assessment of $\overline{RCovLs(X, Y)}$ estimation

3.3 Estimation of $IRLs(X, Y)$

We proposed to define $IRLs(X, Y)$ the informative Procruste correlation coefficient as follow.

$$IRLs(X, Y) = \frac{ICovLs(X, Y)}{IVarLs(X) IVarLs(Y)} \quad (7)$$

Like (X, Y) $IRLs(X, Y) \in [0; 1]$ with the 0 value corresponding to not correlation and the maximum value 1 reached for two strictly homothetic data sets.

3.4 Testing significance of $IRLs(X, Y)$

Significance of $IRLs(X, Y)$ can be tested using permutation test as defined in Jackson (1995) or Peres-Neto and Jackson (2001) and implemented respectively in the `protest` method of the `vegan` R package (Dixon, 2003) or the `procuste.rtest` method of the `ADE4` R package Dray and Dufour (2007).

It is also possible to take advantage of the Monte-Carlo estimation of $\overline{RCovLs(X, Y)}$ to test that $ICovLs(X, Y)$ and therefore $IRLs(X, Y)$ are greater than expected under random hypothesis. Let counting over the k randomization when $\overline{RCovLs(X, Y)}_k$ greater than $CovLs(X, Y)$ name this counts $N_{>CovLs}$. P_{value} of the test can be estimated following Equation (8).

$$P_{\text{value}} = \frac{N_{>CovLs}}{k} \quad (8)$$

3.5 Simulating data for testing sensibility to overfitting

To test sensibility to overfitting correlations were mesured between two random matrices of same dimensions. Each matrix is $n \times p$ with $n = 20$ and $p \in [2, 50]$. Each p variables are drawn from a centered and reduced normal distribution $\mathcal{N}(0, 1)$. Eight correlation coefficients have been tested: RLs the original procruste coefficient, $IRLs$ this work, RV the original R for vector data (Robert and Escoufier, 1976), $RVadjMaye$, $RV2$ and $RVadjGhaziri$ three modified versions of RV (El Ghaziri and Qannari, 2015; Mayer et al., 2011; Smilde et al., 2009), $dCor$ the original distance correlation coefficient (Székely et al., 2007) and $dCor_ttest$ a modified version of $dCor$ not sensible to overfitting (Székely and Rizzo, 2013). For each p value, 100 simulations were run. Computation of $IRLs$ is estimated with 100 randomizations.

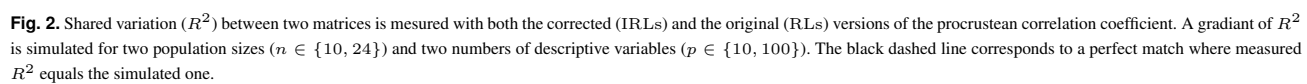
For $p = 1$ random vectors with $n \in [3, 25]$ are generated. As above data are drawn from $\mathcal{N}(0, 1)$ and $k = 100$ simulations are run for each n . The original Pearson correlation coefficient R and the modified version IR are used to estimate correlation between both vectors.

3.6 Empirical assessment of α -risk for the $CovLs$ test

To assess empirically the α -risk of the procruste test based on the randomisations realized during the estimation of $\overline{RCovLs(X, Y)}$, distribution of P_{value} under the H_0 is compared to a uniform distribution between 0 and 1 ($\mathcal{U}(0, 1)$). To estimate such empirical distribution, $k = 1000$ pairs of $n \times p$ random matrices with $n = 20$ and $p \in \{10, 20, 50\}$ are simulated under the null hypothesis of independancy. Procruste correlation between whose matrices is tested based on three tests. Our proposed test ($CovLs.test$), the `protest` method of the `vegan` R package and the `procuste.rtest` method of the `ADE4` R package. Conformance of the distribution of each set of k P_{value} to $\mathcal{U}(0, 1)$ is assessed using the Cramer-Von Mises test (Csörgő and Faraway, 1996) implemented in the `cvm.test` function of the R package `gofTest`.

3.7 Empirical power assessment for the $CovLs$ test

To evaluate relative power of the three considered tests, pairs of to random matrices were produced for various $p \in \{10, 20, 50, 100\}$, $n \in \{10, 15, 20, 25\}$ and two levels of shared variances $R^2 \in \{0.05, 0.1\}$. For each combination of parameters, $k = 1000$ simulations are run. Each



Text Text Text Text Text Text Text Text Text Text Text Text
Text Text Text Text Text Text Text. Figure 2 shows that the above

Text Text Text Text Text Text Text Text Text Text Text Text Text
Text Text Text Text Text Text Text. Figure 2 shows that the above method
Text Text Text Text

Text Text Text Text Text Text Text. Bauer *et al.*, 2007 might want to know about text text text text

This work has been supported by the... Text Text Text Text.

Bauer, M., Klau, G. W., and Reinert, K. (2007). Accurate multiple sequence-structure alignment of RNA sequences using combinatorial optimization. *BMC Bioinformatics*, **8**, 271.

Bravais, A. (1844). *Analyse mathématique sur les probabilités des erreurs de situation d'un point*. Impr. Royale.

Csörgő, S. and Faraway, J. J. (1996). The exact and asymptotic distributions of Cramér-Von mises statistics.

Dixon, P. (2003). VEGAN, a package of R functions for community ecology. *J. Veg. Sci.*, **14**(6), 927–930.

Dray, S. and Dufour, A.-B. (2007). The ade4 package: Implementing the duality diagram for ecologists. *Journal of Statistical Software, Articles*, **22**(4), 1–20.

El Ghaziri, A. and Qannari, E. M. (2015). Measures of association between two datasets; application to sensory data. *Food Qual. Prefer.*, **40**, 116–124.

Escoufier, Y. (1973). Le traitement des variables vectorielles. *Biometrics*, pages 751–760.

- Gower, J. C. (1971). Statistical methods of comparing different multivariate analyses of the same data. *Mathematics in the archaeological and historical sciences*, pages 138–149.
- Jackson, D. A. (1995). PROTEST: A PROcrustean randomization TEST of community environment concordance. *Écoscience*, **2**(3), 297–303.
- Lingoes, J. C. and Schönemann, P. H. (1974). Alternative measures of fit for the schönemann-carroll matrix fitting algorithm. *Psychometrika*, **39**(4), 423–427.
- Mayer, C.-D., Lorent, J., and Horgan, G. W. (2011). Exploratory analysis of multiple omics datasets using the adjusted RV coefficient. *Stat. Appl. Genet. Mol. Biol.*, **10**, Article 14.
- Peres-Neto, P. R. and Jackson, D. A. (2001). How well do multivariate data sets match? the advantages of a procrustean superimposition approach over the mantel test. *Oecologia*, **129**(2), 169–178.
- Ramsay, J. O., ten Berge, J., and Styan, G. P. H. (1984). Matrix correlation. *Psychometrika*, **49**(3), 403–423.
- Robert, P. and Escoufier, Y. (1976). A unifying tool for linear multivariate statistical methods: The RV- coefficient. *J. R. Stat. Soc. Ser. C Appl. Stat.*, **25**(3), 257–265.
- Smilde, A. K., Kiers, H. A. L., Bijlsma, S., Rubingh, C. M., and van Erk, M. J. (2009). Matrix correlations for high-dimensional data: the modified RV-coefficient. *Bioinformatics*, **25**(3), 401–405.
- SzéKely, G. J. and Rizzo, M. L. (2013). The distance correlation t-test of independence in high dimension. *J. Multivar. Anal.*, **117**, 193–213.
- SzéKely, G. J., Rizzo, M. L., and Bakirov, N. K. (2007). Measuring and testing dependence by correlation of distances. *Ann. Stat.*, **35**(6), 2769–2794.