

Performing Configural Frequency Analysis by Means of Recurrence Quantification Analysis¹

*Sebastian Wallot*²

Abstract

Recurrence Plots (RPs) were developed at the end of the 1980's as visualization tools for complex dynamics exhibited by time series measures from physical and dynamic systems. At the beginning of the 1990's RPs were further developed into Recurrence Quantification Analysis (RQA), which allowed for numerical characterizations – and analyses – of such time series. In the past couple of years, RQA has been further developed to analyze coupling between two time series, the dynamics of multivariate time series, and can be used to derive correlations between two multivariate time series. The aim of the current paper is to expand on another line of development, which is to extend recurrence-based techniques to the analysis of sample data – here, specifically to use RQA in order to perform the classical or so-called first-order Configural Frequency Analysis (CFA), as developed by Lienert in the 1970's. First, RQA and some of its extensions will be introduced, that allow to properly deal with multidimensional categorical time series or sequences. Then, we will combine these existing techniques to create a framework that can perform a CFA-type of analysis, based on a bootstrap approach.

Keywords: Configural frequency analysis; recurrence quantification analysis; multidimensional data; categorical data

¹ A previous version of this article has been presented in a symposium in honor of Gustav A. Lienert's 100th birthday at the 15th Conference of the Section Methods and Evaluation in the German Psychological Society (DGPs), September 15th to 17th, 2021. This contribution was envisioned for a planned Special Issue on the Gustav A. Lienert Symposium (Guest Editor: Mark Stemmler).

² Department of Psychology, Leuphana University of Lüneburg, Lüneburg, Germany & Max Planck Institute for Empirical Aesthetics, Frankfurt a.M., Germany. *Correspondence concerning this article should be addressed to:* Sebastian Wallot, Leuphana University of Lüneburg, Universitätsallee 1, 21335 Lüneburg, Germany, Sebasitan.wallot@leuphana.de

Introduction

The current paper introduces an extension of Recurrence Quantification Analysis (RQA), which allows to perform the basic Configural Frequency Analysis (CFA) as originally proposed by Lienert (1970). Configural Frequency Analysis provides a mean to detect whether certain patterns (i.e., configurations of categories) in categorical data sets, where each unit of observation is classified in terms of multiple categories, occurs reliably more or less often than what would be expected from the basic occurrences category given the assumption of independence. RQA, on the other hand, was originally proposed as a method to visualize (Eckmann et al., 1987) and quantify (Webber & Zbilut, 1994; Zbilut & Webber, 1992) nonlinear correlation patterns within time series from physics and physiology (Marwan et al., 2007).

RQA has, however, turned out to be a very versatile technique, and applications have been developed to use recurrence quantification analysis for the analysis of multidimensional categorical sequences (Angus et al., 2012; Cox et al., 2016), to test for differences in distributional shapes between multiple samples (Wallot & Leonardi, 2018), and to analyze multidimensional time series (Wallot et al., 2016). The present paper builds on both of these developments to demonstrate how the basic RQA routine can be extended to perform analysis of multidimensional categorical data sets in order to detect the presence of specific patterns of categories.

To do so, I will first briefly introduce the basic recurrence plot for nominal data, the chromatic recurrence plot, and the multidimensional recurrence plot, which are combined to display the distribution of patterns of categories in a data set. Because recurrence-based analyses are model-free, there is no specific test statistic associated with their application. Hence, I will develop a bootstrap procedure that will allow to perform a global test for the distribution of patterns within a data set, in analogy to the χ^2 -omnibus test of the basic or so-called first order Configural Frequency Analysis. Finally, I will apply this procedure to the original data by Lienert (1970) set on syndromes associated with LSD-intake, and compare its results to the classical CFA-procedure.

Recurrence Quantification Analysis (RQA)

Recurrence Quantification Analysis (RQA) is based on the Recurrence Plot (RP), which was first introduced by Eckmann et al. (1987). In subsequent years, Joseph Zbilut and Charles Webber (1992) introduced the quantification of those plots to numerically capture the dynamics and (nonlinear) autocorrelation properties of time series. Accordingly, the Recurrence Plot and Recurrence Quantification Analysis have originally been developed as means to visualize and quantify the dynamics of time series and sequences. These applications were geared towards continuously sampled

signals, but soon it became clear that RQA could also be applied to categorical data (Dale & Spivey, 2006; Dale et al., 2011).

How does a recurrence plot look like? Let us use the example of the following nominal sequence to illustrate RPs: THE RAIN IN SPAIN FALLS MAINLY IN THE PLAIN – the famous phrase sung by the character Eliza Doolittle in the musical “My fair lady”. As the name indicates, recurrence plots are about repetitions in a sequence, and a rhyming sentence like the one above likely contains elements – or whole subsequences – that appear in time and again throughout the sequence.

In the case of such nominal sequences, defining repetitions is simple: If an element “A” occurs on the sixth slot of the sequence, and then occurs again on the 15th slot of that sequence, then it recurs. This is easy because these elements are categorical, and hence we can distinguish between identity relations, like A is A, and A is NOT B. If we want to view such sequences through the lens of a recurrence plot, we first have to chart the sequence against itself to create a two-dimensional plot (Figure 1).

The recurrence plot is a two-dimensional binary matrix, displayed as black and white dots. The black dots on the plot indicate recurrence, the white dots the absence of recurrence. For example, our sequence THE RAIN IN SPAIN FALLS MAINLY IN THE PLAIN, has an “A” on the 6th, 15th, 20th, 26th, and 41st slot. Accordingly, if we look at the recurrence plot in Figure 1, we see black dots on the 6th row, and the 6th, 15th, 20th, 26th, and 41st columns in that plot, indicating the positions where “A” recurs.

Further features of the recurrence plot are the central diagonal, the line of identity (LOI), running from the lower-left to the upper right of the plot. The line of identity simply indicates that the same sequence is perfectly recurrent with itself at lag0, since the same sequence is charted on both, the x- and the y-axis. Moreover, the recurrence plot is symmetrical about the LOI, so we see the same patterns of recurrences in the upper-left and the lower-right triangles of the plot.

If we were interested in analyzing sequences, we could now calculate different measures describing its sequential properties based on different patterns on the plot, such as isolated recurrences, diagonally adjacent recurrences, and vertically/horizontally adjacent recurrence (Marwan et al., 2007). If we were to analyze continuously sampled data, that is, a proper time series, we would furthermore have to set specific parameters – the embedding parameters – to optimize the quantification of recurrences on that plot (e.g., Abarbanel, 1996). For readers interested in such applications, I refer those readers to tutorial introductions describing these applications (Wallot, 2017).

In the current paper, however, we are exclusively interested in dealing with nominal data, hence no embedding parameters are needed. Furthermore, since we are not interested in quantifying sequential properties, but are rather interested in the relative occurrences of different combination of categories, we also do not need to apply recurrence measures that are quantifying sequential properties. Rather, we want to use a recurrence plot quantify the different number of occurrences of categories, as well as their clustering. In order to do so, we need to consider two extensions of the basic

recurrence plot shown in Figure 1, which multidimensional recurrences and chromatic recurrence.

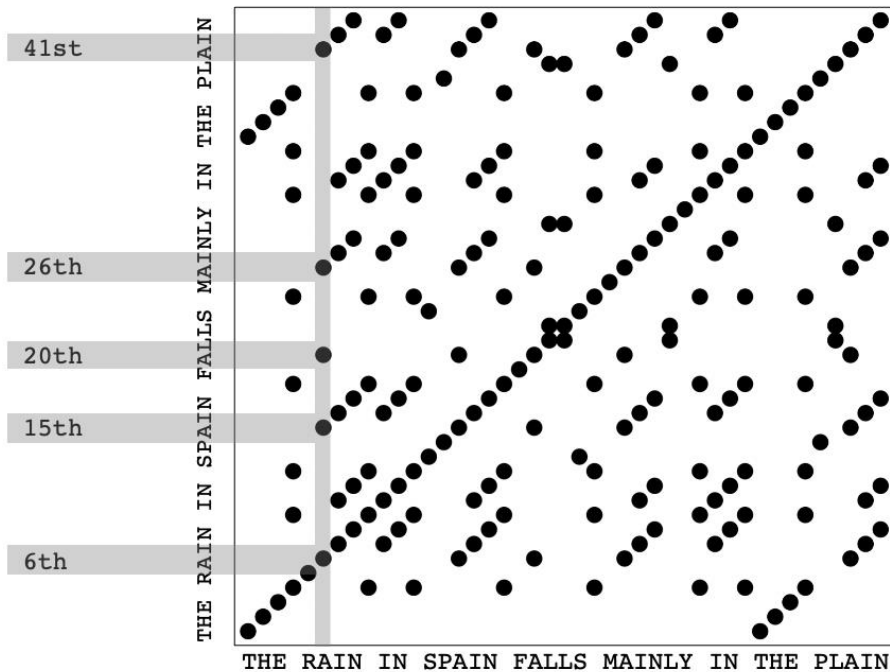


Figure 1. Illustration of a recurrence plot for categorical data, using the text sequence THE RAIN IN SPAIN FALLS MAINLY IN THE PLAIN from the musical my fair lady. For categorical data, recurrences are determined through identity (i.e., “A” equals “A” and only “A”, and is distinct from every other letter). As can be seen in the plot, the latter Appears in the sequence five times in positions 6, 15, 20, 26, and 41, as indicated by the black recurrence points.

First, let us consider multidimensional recurrence quantification analysis. As the name implies, this extension of RQA allows for the analysis of multivariate time series data. To do so, let us expand on our example above. The line of the text taken from “My fair lady” continues in the musical by the exclamation of I THINK SHE HAS GOT IT I THINK SHE HAS GOT IT by Professor Higgins in response to Eliza Doolittle’s perfect pronunciation of THE RAIN IN SPAIN STAYS MAINLY IN THE PLAIN. However, we have now another important sequence of categorical data, which is the

of the speaker variables is also identical (e.g., a 0 for both “A”s). The resulting multidimensional recurrence plot is shown in Figure 3b.

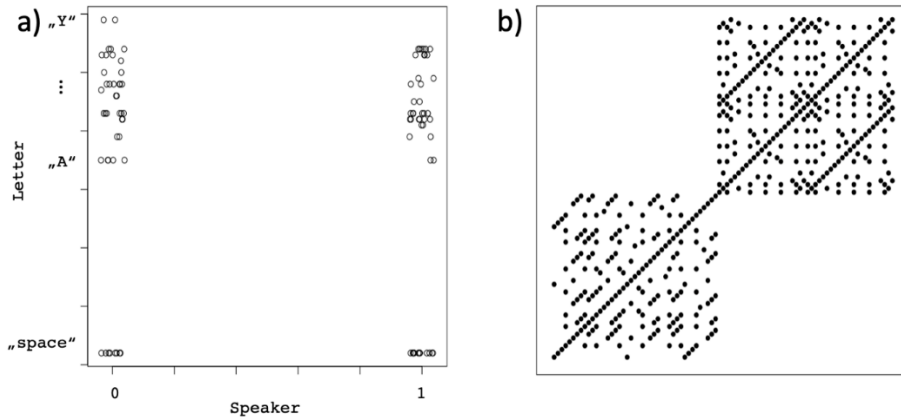


Figure 3. Illustration of a state-space and a multidimensional recurrence plot. a) “State space” of letters in the text (here only the letter “Y”, “A”, and the “space” between words are shown on the y-axis of illustrative purposes) and speakers (0 = Eliza Doolittle, 1 = Professor Higgins). b) Multidimensional recurrence plot that contains the pattern of recurrences when the information from both of the categorical dimensions is integrated.

If we were to calculate recurrence measures from such a plot, however, these measures would still be calculated across the two different categories (i.e., the text by Eliza Doolittle and Professor Higgins), and would reflect their conversation, not their individual contributions to that conversation. In order to treat these combinations of categorical data separately – as we will have to do for a CFA-type of analysis – we will need to add identifiers to this recurrence plot.

Such an analysis of different parts of a recurrence plot has been proposed by Cox and colleagues (2016), who called such a plot a chromatic recurrence plot. Chromatic recurrence plots divide the structures on a plot by a category identifier, and calculate RQA measures separately by the identifier. The chromatic recurrence plot of our categorical data example can be seen in Figure 4.

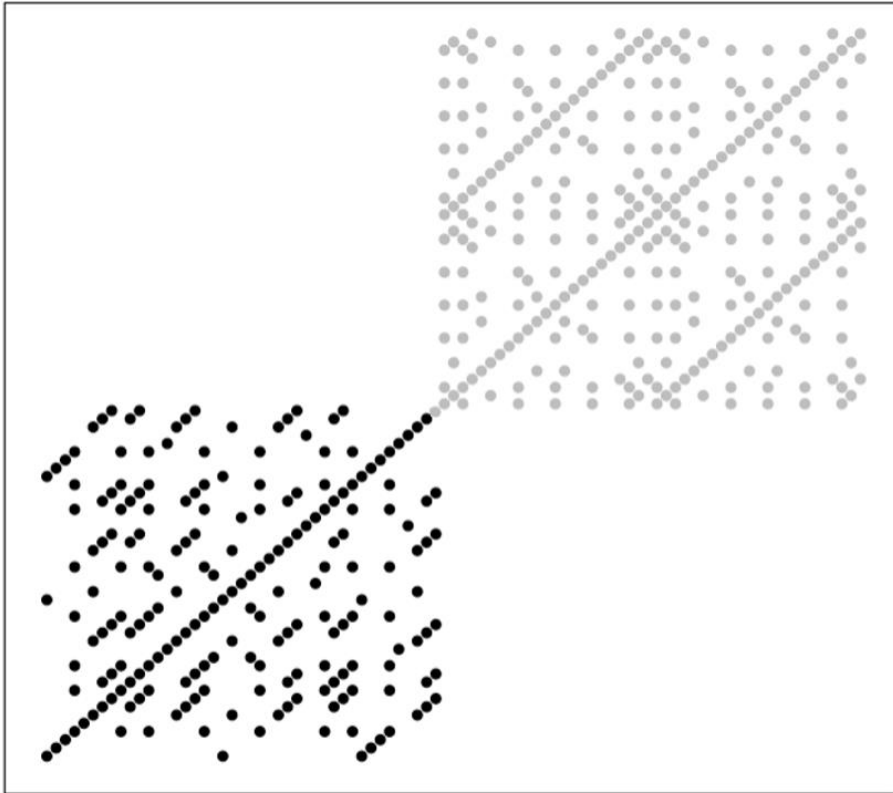


Figure 4. Illustration of a chromatic multidimensional recurrence plot for the two-dimensional categorical data, consisting of text and speaker identity. Here, recurrences for speaker 0 (Eliza Doolittle) are marked in black, while recurrences for speaker 1 (Professor Higgins) are marked in grey.

Now we can compare recurrence measures across the different types of recurrence plots we have defined so far. The most basic recurrence measure is *percent recurrence* (*%REC*), which is simply the sum of all recurrence points in a recurrence plot divided by the size of that plot (i.e., the number of all possible recurrence points). As said above, there are multiple recurrence measures. In the context of the current example, another measure will be of central importance, the *average diagonal length* (*ADL*). This measure is calculated by counting the number of recurrence points that have at least one diagonally adjacent recurrence point, and divide the sum of these points by the number of the so defined diagonals. Table 1 charts the two resulting measures by the different recurrence plots.

Table 1. Recurrence measures for the different types of recurrence plots.

Type of RP	Type of measure	Value
Univariate – text	%REC	11.27%
	ADL	3.72
Univariate – speaker	%REC	50.06%
	ADL	23.16
Multivariate (text & speaker)	%REC	50.91%
	ADL	5.84
Chromatic Multivariate (black)	%REC	10.87%
	ADL	4.15
Chromatic Multivariate (grey)	%REC	13.28%
	ADL	14.43

Note. The black and grey anisotropic recurrences effectively provide RQA measures of the text for the different speakers (black: Elize Doolittle; grey: Professor Higgins).

As we can see, these measures differ substantially between the two univariate and the multivariate recurrence plot. Moreover, we see that we get two different values – one for each category, i.e., the text spoken by Professor Higgins and Eliza Doolittle – for the chromatic multidimensional recurrence plot, which is what we will need to link RQA to CFA.

Using RQA to perform CFA

Our goal here is to perform an analysis akin to the classical Configural Frequency Analysis (CFA). To do this, we need to first compose clusters of combinations of categories, and then extract their frequencies from a recurrence plot. The outlined procedure above delivers this. Let us apply this procedure to the famous data set on LSD-intake-related syndromes, as published by Lienert (1970).

The original question behind this data set was related to findings suggesting that LSD intoxication was characterized by narrowed consciousness, thought disturbances, and affective disturbances. Furthermore, the data seem to suggest that in case of a LSD intake, either all three of these symptoms were present simultaneously, or only one of them, but rarely any combination of only two symptoms (e.g., narrowed consciousness and affective disturbances). To quantitatively test for the presence of such a syndrome-pattern in the data, Lienert (1970) proposed the classical CFA methods to test

whether specific combinations of patterns of symptoms occurred reliably more often than what would be expected from the base frequencies of the individual symptoms in the data set.

Sixty-five students swallowed LSD and the presence or absence of the three symptoms was registered. Accordingly, each participant could be placed in one of possible patterns that could occurred given combinations of the three symptoms (i.e., 2^3). Table 2 below shows the occurrences for each of these patterns. Visual inspection of that table shows that all participants showed at least one of the symptoms. Moreover, the relative distribution of the participants across the categorial combinations seem to speak in favor of the suggestion, that the co-occurrence of the symptoms was interdependent: There were relatively few participants that reported the presence of only two of the three symptoms, and comparatively more that reported either 1 or 3 of these symptoms.

Table 2. Observed and expected frequencies patterns of symptoms for the LSD data set

NC	TD	AD	<i>f</i>	<i>e</i>
0	0	0	0	4.7
1	1	1	20	12.5
1	1	0	1	6.8
1	0	1	4	11.4
1	0	0	12	6.2
0	1	1	3	9.5
0	1	0	10	5.2
0	0	1	15	8.6

Note. The different symptoms are coded for with the letters “NC”, “TD”, and “AD”: “NC” = narrowed consciousness, “TD” = thought disorder, and “AD” = affective disturbances. A “1” marks that the specific symptom was present in participants, a “0” marks that the specific symptom was absent. “*f*” are the observed frequencies of participants that fall into a specific three-dimensional pattern of symptoms, while “*e*” is the expected number of participants for a specific pattern, assuming independent distribution of participants across patterns given the basic number of occurrences of each symptom.

To conduct Configural Frequency Analysis, we use the marginal sums to compute the expected values given the base-occurrences of each symptom and the assumption of independence of the symptom combinations. Now, a χ^2 -value can be deduced from the observed and expected occurrences. This allows the application of individual tests of each of the categories regarding the deviation of its observed and expected values. Moreover, the sum of these deviations allows the application of a χ^2 -omnibus test, testing the null-hypothesis which allows to test deviations from the null-hypothesis that the occurrences of the different symptoms are independently each other.

For the LSD-data set, the χ^2 -value for this omnibus test is $\chi^2(4) = 38.02$, which indicates a significant deviation from the distribution for occurrences over the patterns of categories given the null-hypothesis of independent combinations at $\alpha = 0.05$.

Let us reproduce this test using the recurrence procedure outline above. First, we construct a state-space for our data set as in the My Fair Lady example above, but with our three categories (and some jitter added; see Figure 5a). Now we compute a recurrence plot using multidimensional recurrence quantification procedure, and apply identifiers for each of the patterns of symptoms according to chromatic recurrence analysis (Figure 5b).

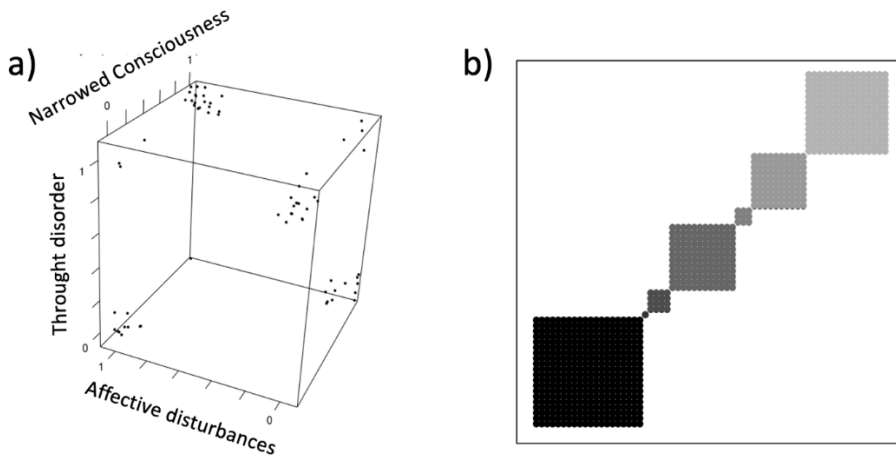


Figure 5. Illustration of the state-space and recurrence plot of the LSD data set. a) State-space of the number of occurrences of the different patterns of symptoms (narrowed consciousness, thought disorder, and affective disturbances). b) Resulting multidimensional chromatic recurrence plot. Note that the square-patches of recurrences are proportional in size to the number of occurrences for each symptom pattern.

After having computed this plot, we can now apply a recurrence measure – the average vertical line length (*AVL*) – to the data. As this measure is applied separately to the different identifiers, we receive a separate *AVL*-value for each of the 8 different patterns of symptoms. As the name implies, *AVL* measures the average vertical line length of a patch of recurrence. As our different categories of occurrences always from squares of recurrence (see Figure 5b), *AVL* equals the square-root of the squares of recurrences associated with each category (e.g., a square-patch of 3-times-3 recurrences has a height of 3 vertically adjacent recurrence points; hence, *AVL* equals 3 for such a patch). This reproduces exactly the observed frequencies that appear in Table 2 above.

However, if we want to stay within the recurrence framework, we cannot have fractions of recurrence points. A data point is either recurrent (= 1) or not (= 0), so we cannot analytically deduce expected frequencies as in the original CFA method, described above. Hence, in order to compose a test of significance, we need to employ a two-step boot-strapping procedure.

In the first step, we create recurrence plot based on a boot-strapped data set of equal sample size to the original sample ($n = 65$, in our case) where the combinations of categories are independent of each other. Then, we run this recurrence plot through the same procedure as the original data, and apply the *AVL* measure. Now we can calculate the difference between the *AVL* measure for each of the categorical patterns of the original data and the *AVL* measure for each of the independent-occurrences boot-strapped data.

If we want to construct an omnibus test to evaluate the overall deviation of the observed data from independently composed combinations of categories – which is our goal – this procedure will not suffice, because sum of the deviations is always zero. Also, taking the absolute value of the deviations is not enough, because this sum will always be positive. Accordingly, we have to add a second step, which is effectively bringing in a baseline deviation given the expected occurrences under the assumption of independence.

In this second step, we draw another bootstrap-sample under the condition of independence and subtract the numbers from the second sample from the one of the first sample, taking the absolute value. This provides us with a measure of the expected deviation of cases under the condition of independence. Now, we can subtract the absolute deviation of the independence-independence case from the absolute deviation of the observed data and the independent data set, which provides us with a measure of how much bigger the deviation between of the observed data is from the independence data compared to the independence case alone. Figure 6 illustrates the procedure. Equation 1 summarized the computation of the difference scores over which a confidence interval can be calculated:

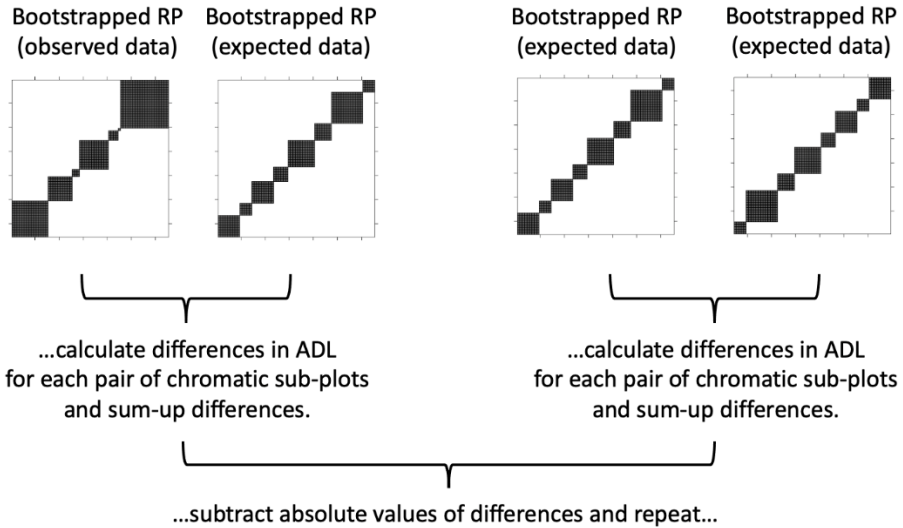


Figure 6. Illustration of the bootstrapping procedure to deduce a confidence interval that test the hypothesis that the deviation of the observed occurrences across the pattern of categories is significantly different from the occurrences expected under the independence assumption. First, a recurrence plot is computed based on the bootstrap of the observed occurrence and another recurrence plot is computed based on the bootstrap of the expected occurrence. For each of the patterns of categories, the *AVL*-measure is computed and pairwise subtracted. Then, the absolute values of these differences are summed. This gives an estimate of the difference between the observed and expected data. Second, two recurrence plots are computed based on bootstraps of the expected occurrences. Likewise, the *AVL*-measure is computed for each of the patterns of categories, which are subtracted from each other, and their absolute differences are summed. Finally, the sum of the differences of the observed-expected pair is subtracted from the sum of the differences of the observed-expected pair. This gives us a measure of how much bigger the deviation between the observed-expected data is compared to the differences one would expect from two expected pairs alone. Then, the process is repeated to generate a distribution of such difference values, over which a confidence interval can be computed.

$$\sum_{i=1}^n (|o_i - e_i|) - \sum_{i=1}^n (|e_i - e_i|) \quad (\text{eq. 1})$$

Where:

n = number of patterns of categories (8 in the present example).

i = the index of the pattern of categories.

o = the average diagonal line measure of a pattern for a bootstrap based on the observed occurrences.

e = the average diagonal line measure of a pattern for a bootstrap based on the expected occurrences.

This procedure can now be used to construct a confidence interval testing the null-hypothesis that the difference between the observed data and the independence-data is zero (i.e., when the values of zero is included in the confidence interval) or bigger than zero (when zero lies outside of the confidence interval), and thus provides an omnibus test for our data. In case of the LSD-dataset, the confidence bounds for the 95%-one-tailed confidence interval are [8, 62], indicating a significant deviation of the observed data from the null-hypothesis of independent composition of the categories of symptoms, analogously to the results obtained from CFA.

Similarly, confidence intervals can be deduced for each of the eight individual cells – the patterns of categories – using the bootstrapped absolute differences between the observed and expected frequencies. Table 3 shows the confidence bounds for the respective 95%-one-tailed intervals. As expected, and similarly to the results obtained by the CFA procedure, these confidence bounds are close to 0 and suggest a less reliable difference, because the individual cells have fewer observations compared to the whole sample.

Table 3. Confidence bounds for the eight individual patterns of categories.

NC	TD	AD	f	95%-CI <i>one-tailed</i>
0	0	0	0	[2, 14]
1	1	1	20	[4, 18]
1	1	0	1	[2, 15]
1	0	1	4	[2, 22]
1	0	0	12	[1, 19]
0	1	1	3	[1, 19]
0	1	0	10	[1, 23]
0	0	1	15	[3, 18]

Note. The different symptoms are coded for with the letters “NC”, “TD”, and “AD”: “NC” = narrowed consciousness, “TD” = thought disorder, and “AD” = affective disturbances. The confidence bounds show the 95%-one-tailed cut-offs.

Conclusion

The current paper briefly introduced recurrence quantification analysis as a time series analysis technique, and proposed an extension of recurrence analysis to conducted analysis of multidimensional categorical sample data in analogy to configural frequency analysis, originally proposed by Lienert (1970). Combining several recurrence analysis methods (particularly multidimensional recurrence quantification analysis and chromatic recurrence quantification analysis) together with a bootstrapping procedure allowed to construct a confidence interval, which can be used as an omnibus test for the null-hypothesis, that the categories in a multidimensional categorical data set cluster independently of each other, in analogy to the χ^2 -omnibus test provided by CFA.

This shows the versatility of the recurrence analysis framework to mimic various statistical procedures testing not only for properties of time series data (Marwan et al, 2007), but extending earlier works that showed that RQA can also be used to analyze properties of sample distributions (Wallot & Leonardi, 2018).

Future developments along this path could potentially include advances methods of configural frequency analysis, such as two-sample CFA or longitudinal CFA (see Stemmler, 2020). However, while the extension of RQA proposed in the present paper is of interest from the perspective of recurrence-based analysis, it remains an open question in how far these methods are superior or even equal to the existing CFA procedures. While bootstrapping procedures as employed in the current method are robust, they are also potentially less sensitive, and further research is needed that provides in-depth head-to-head comparisons between the ability of RQA and CFA to reliably detect differences in the distribution of occurrences across patterns of categorical data.

Acknowledgements

Sebastian Wallot acknowledges funding from the German Research Foundation's Heisenberg programme (DFG; 442405852).

References

- Abarbanel, H. (1996). *Analysis of Observed Chaotic Data*. New York, NY: Springer.
- Angus, D., Watson, B., Smith, A., Gallois, C., & Wiles, J. (2012). Visualising conversation structure across time: Insights into effective doctor-patient consultations. *PLOS ONE*, *7*, e38014.
- Cox, R.F.A., van der Steen, S., Guevara, M., de Jonge-Hoekstra, L., van Dijk, M. (2016). Chromatic and Anisotropic Cross-Recurrence Quantification Analysis of Interpersonal Behavior. In: Webber, Jr., C., Ioana, C., and Marwan, N. (Eds). *Recurrence Plots and Their Quantifications: Expanding Horizons*. Springer Proceedings in Physics (Vol. 180, pp. 209-25). Heidelberg: Springer.
- Dale, R., & Spivey, M. J. (2006). Unraveling the dyad: Using recurrence analysis to explore patterns of syntactic coordination between children and caregivers in conversation. *Language Learning*, *56*, 391-430.
- Dale, R., Warlaumont, A. S., & Richardson, D. C. (2011). Nominal cross recurrence as a generalized lag sequential analysis for behavioral streams. *International Journal of Bifurcation and Chaos*, *21*, 1153-1161.
- Eckmann, J. P., Kamphorst, S. O., & Ruelle, D. (1987). Recurrence plots of dynamical systems. *Europhysics Letter*, *4*, 973-977.
- Lienert, G. A. (1970). Konfigurationsfrequenzanalyse einiger LSD-Wirkungen. *Arzneimittelforschung*, *20*, 912-913.
- Stemmler, M. (2020). *Person-Centered Methods. Configural Frequency Analysis (CFA) and Other Methods for the Analysis of Contingency Tables* (2nd Edition). Cham: Springer.
- Marwan, N., Romano, M. C., Thiel, M., & Kurths, J. (2007). Recurrence plots for the analysis of complex systems. *Physics reports*, *438*(5-6), 237-329.
- Wallot, S. (2017). Recurrence quantification analysis of processes and products of discourse: A tutorial in R. *Discourse Processes*, *54*, 382-405.
- Wallot, S., & Leonardi, G. (2018). Deriving inferential statistics from recurrence plots: A recurrence-based test of differences between sample distributions and its comparison to the two-sample Kolmogorov-Smirnov test. *Chaos: An interdisciplinary journal of nonlinear science*, *28*, 085712.
- Wallot, S., Roepstorff, A., & Mønster, D. (2016). Multidimensional Recurrence Quantification Analysis (MdrQA) for the analysis of multidimensional time-series: A software implementation in MATLAB and its application to group-level data in joint action. *Frontiers in Psychology*, *7*, 1835.
- Webber Jr, C. L., & Zbilut, J. P. (1994). Dynamical assessment of physiological systems and states using recurrence plot strategies. *Journal of Applied Physiology*, *76*, 965-973.
- Zbilut, J. P., & Webber Jr, C. L. (1992). Embeddings and delays as derived from quantification of recurrence plots. *Physics Letters A*, *171*(3-4), 199-203.