

Robust Rank Correlation Coefficients on the Basis of Fuzzy Orderings: Initial Steps

U. Bodenhofer¹ and F. Klawonn²

¹Institute of Bioinformatics, Johannes Kepler University Linz
4040 Linz, Austria

²University of Applied Sciences Braunschweig/Wolfenbüttel
38302 Wolfenbüttel, Germany

bodenhofer@bioinf.jku.at, f.klawonn@fh-wolfenbuettel.de

Abstract

The goal of this paper is to demonstrate that established rank correlation measures are not ideally suited for measuring rank correlation for numerical data that are perturbed by noise. We propose to use robust rank correlation measures based on fuzzy orderings. We demonstrate that the new measures overcome the robustness problems of existing rank correlation coefficients. As a first step, this is accomplished by illustrative examples. The paper closes with an outlook on future research and applications.

1 Introduction

Correlation measures are among the most basic tools in statistical data analysis and machine learning. They are applied to pairs of observations ($n \geq 2$)

$$(x_i, y_i)_{i=1}^n \tag{1}$$

to measure to which extent the two observations comply with a certain model. The most prominent representative is surely *Pearson's product moment coefficient* [1, 18], often nonchalantly called *correlation coefficient* for short. Pearson's product moment coefficient is applicable to numerical data and assumes a linear relationship as the underlying model; therefore, it can be used to detect linear relationships, but no non-linear ones.

Rank correlation measures [11, 13, 16] are intended to measure to which extent a monotonic function is able to model the inherent relationship between the two observables. They neither assume a specific parametric model nor specific distributions of the observables. They can be applied to ordinal data and, if some ordering relation is given, to numerical data too. Therefore, rank correlation measures are ideally suited for detecting monotonic relationships, in particular, if more specific information about the data is not available. The two most common approaches

are *Spearman's rank correlation coefficient* (short *Spearman's rho*) [20, 21] and *Kendall's tau (rank correlation coefficient)* [2, 12, 13].

This paper argues why these well-known rank correlation measures are not ideally suited for measuring rank correlation for numerical data that are perturbed by noise. Consequently, we propose a robust rank correlation measure on the basis of fuzzy orderings. The superiority of the new measure is demonstrated by means of illustrative examples.

2 An Overview of Rank Correlation Measures

Assume that we are given a family of pairs as in (1), where all x_i and y_i are from linearly ordered domains X and Y , respectively. *Spearman's rho* is computed as

$$\rho = 1 - 6 \frac{\sum_{i=1}^n (r(x_i) - r(y_i))^2}{n(n^2 - 1)},$$

where $r(x_i)$ is the rank of value x_i if we sort the list (x_1, \dots, x_n) ; $r(y_i)$ is defined analogously. So, Spearman's rho measures the sum of quadratic distances of ranks and scales this measure to the interval $[-1, 1]$. It can be checked easily that a value of 1 is obtained if the two rankings coincide and that a value of -1 is obtained if one ranking is the reverse of the respective other. Note that the above definition of $r(x_i)$ and $r(y_i)$ was simplified, because it did not take coinciding values, so-called *ties*, into account. In such a case, the values $r(x_i)$ are usually defined as the mean value of all ranks of consecutive coinciding values in the sorted list.

To define the *Kendall tau rank correlation coefficient*, we need to introduce the concepts of concordance, discordance and ties first. For a given pair of indices $(i, j) \in \{1, \dots, n\}^2$, we say that (i, j) is *concordant* if $x_i < x_j$ and $y_i < y_j$; we say that (i, j) is *discordant* if $x_i < x_j$ and $y_i > y_j$. If $x_i = x_j$, we say that (i, j) is a *tie in the first component*. If $y_i = y_j$, we say that (i, j) is a *tie in the second component*. We simply say that (i, j) is a *tie* if (i, j) is a tie in either component.

Let us denote the numbers of concordant, discordant and tied pairs as follows:

$$C = |\{(i, j) \mid x_i < x_j \text{ and } y_i < y_j\}|$$

$$D = |\{(i, j) \mid x_i < x_j \text{ and } y_i > y_j\}|$$

$$T = |\{(i, j) \mid x_i = x_j\}|$$

$$U = |\{(i, j) \mid y_i = y_j\}|$$

Then the basic variant of *Kendall's tau* which we denote with τ_a is computed as the quotient

$$\tau_a = \frac{C - D}{\frac{1}{2}n(n - 1)}.$$

If there are no ties and the two rankings coincide, we have $\frac{1}{2}n(n - 1)$ concordant and no discordant pairs, so $\tau_a = 1$; if we have no ties and one ranking is the reverse of the respective other, we have no concordant and $\frac{1}{2}n(n - 1)$ discordant pairs, so

a value of $\tau_a = -1$ is obtained. So, in these extremal cases, Kendall's tau gives the same results as Spearman's rho.

Ties, no matter whether in the first or in the second list, are not counted in the above definition of τ_a , so they lower the absolute value of τ_a . Therefore, τ_a is best suited for detecting strictly monotonic relationships, but not ideally suited in the presence of ties. A well-established second variant [13] is the following:

$$\tau_b = \frac{C - D}{\sqrt{\frac{1}{2}n(n-1) - T} \sqrt{\frac{1}{2}n(n-1) - U}},$$

It takes ties into account, but is still not fully robust to ties (see next section). A simple and tie-robust rank correlation measure is the *gamma rank correlation measure* according to Goodman and Kruskal [11] that is defined as

$$\gamma = \frac{C - D}{C + D}.$$

Finally, we remark that τ_a , τ_b and γ coincide in all cases where no ties occur in the data.

3 Motivation

Historically, all rank correlation measures highlighted above have been introduced with the aim to measure rank correlation of ordinal data (e.g. natural numbers, marks, quality classes, ranks). The measurement of rank correlation for *real-valued data*, however, is equally important in statistics and machine learning, but raises completely new issues. Depending on the source, numerical data are almost always subject to random perturbations—noise. The concepts introduced above do not take this into account. Pairs are counted as concordant or discordant only on the basis of ordering relations, but without taking into account that only minimal differences may decide whether a pair is concordant or discordant. If one observable depends on the other in a clearly monotonic way and if the level of noise is low, then the rank correlation measures introduced above will still reveal this strictly monotonic relationship and will not be compromised by minor local effects of noise. In the presence of a larger percentage of ties, however, already the slightest perturbations may lead to situations in which the above rank correlation coefficients cannot yield meaningful results anymore.

Consider the data sets in Figure 1. We see a monotonic, yet not strictly monotonic, relationship. The left plot shows data without noise, i.e. $y_i = f(x_i)$ for a non-decreasing function f . For these data, we obtain $\rho = 0.737$, $\tau_b = 0.639$ and $\gamma = 1$ (which confirms that γ is most robust to ties). The middle plot shows the same data, but with additive normally distributed noise with zero mean and $\sigma = 0.001$. Although the noise can hardly be seen from the plot, we obtain $\rho = 0.519$ and $\tau_b = \gamma = 0.387$. These results indicate that none of the three measures can adequately handle a large proportion of ties in the presence of noise. For $\sigma = 0.01$ (right plot), the values are slightly lower, but not significantly: $\rho = 0.456$ and

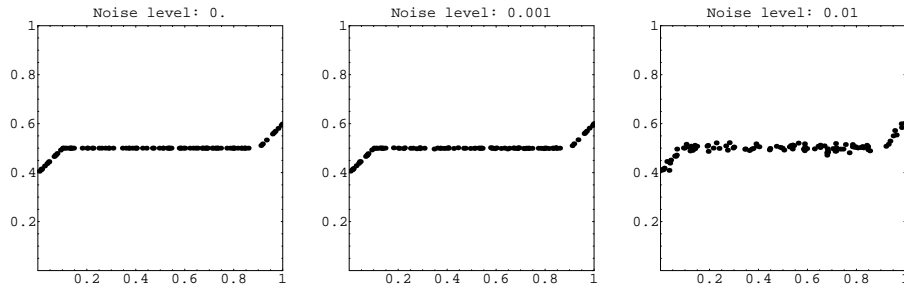


Figure 1: Scatter plots of a simple monotonic relationship with different noise levels.

$\tau_b = \gamma = 0.331$. So we can conclude that it is rather the presence of noise in general than the magnitude of noise that distracts the three rank correlation measures.

The obvious reason for the weakness described above is the fact that all measures only take ordering relationships into account, but neglect similarities/closeness of data points. To illustrate that, consider two pairs (a, c) and (b, c) , where $b > a$. Obviously, this is a tie in the second component. If we add some noise to the second component of the second pair, i.e., if we replace (b, c) by $(b, c + \varepsilon)$, then ε decides whether $((a, c), (b, c + \varepsilon))$ is a tie (for $\varepsilon = 0$), concordant ($\varepsilon > 0$), or discordant ($\varepsilon < 0$), where the magnitude of ε plays no role at all. So we observe a discontinuous behavior. This toy example thereby serves as a proof that all measures introduced above depend on the data in a discontinuous way. Figure 2 illustrates this: it shows graphically how the pairs (i, j) and (j, i) are classified by keeping (x_i, y_i) constant and considering (x_j, y_j) variable. It is obvious that, in close-to-tie situations like the simple example above, little variations of the data can lead to drastic (i.e. discontinuous) changes of the classification of a pair (i, j) .

The question arises how we can define a robust rank correlation measure that depends continuously on the data by taking similarities into account, but still serves as a meaningful measure of rank correlation. Obviously, the measure should be designed such that close-to-tie pairs receive less attention than pairs that are clearly concordant or discordant. A reasonable idea would be to base such a concept on the probabilities to which concordant/discordant pairs are observed as such compared to the probabilities that they are falsely observed as something else. That may be a reasonable approach. Note, however, that such probabilities can only be computed if we know the joint distribution of x and y values or at least if we make distribution assumptions. In practice, such information is most often unavailable and, surely, we do not want to sacrifice the unique feature of rank correlation measures that they are *distribution-free*.

In our opinion, *fuzzy orderings* provide a meaningful way to overcome the difficulties explained above.

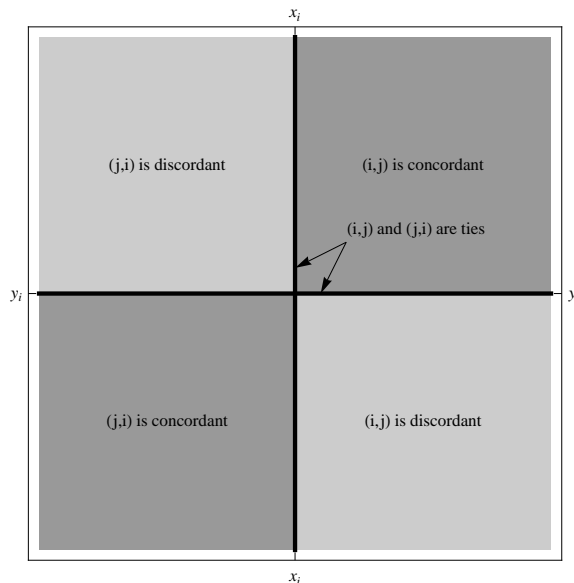


Figure 2: Visualization of the concepts of concordance, discordance and ties for fixed (x_i, y_i) and variable (x_j, y_j) . The horizontal axis corresponds to x_j , while the vertical axis corresponds to y_j .

4 Fuzzy Orderings

Before we can introduce a fuzzy ordering-based rank correlation coefficient, we need to provide some basics of fuzzy orderings. We restrict to an absolutely necessary minimum and refer to literature for details. We assume that the reader is aware of the most basic concepts of triangular norms [15] and fuzzy relations [6, 10, 17].

A fuzzy relation $L : X^2 \rightarrow [0, 1]$ is called *fuzzy ordering* with respect to a t-norm T and a T -equivalence $E : X^2 \rightarrow [0, 1]$, for brevity *T - E -ordering*, if and only if the following three axioms for all $x, y, z \in X$:

- (i) E -Reflexivity: $E(x, y) \leq L(x, y)$
- (ii) T - E -antisymmetry: $T(L(x, y), L(y, x)) \leq E(x, y)$
- (iii) T -transitivity: $T(L(x, y), L(y, z)) \leq L(x, z)$

Moreover, we call a T - E -ordering L *strongly complete* if $\max(L(x, y), L(y, x)) = 1$ for all $x, y \in X$ [4].

Several correspondences between distances and fuzzy equivalence relations are available [7, 8, 14, 23]. From these results, we can easily infer that (assume $r > 0$ in the following)

$$E_r(x, y) = \max(0, 1 - \frac{1}{r}|x - y|)$$

is a $T_{\mathbf{L}}$ -equivalence on \mathbb{R} , where $T_{\mathbf{L}}(x, y) = \max(0, x + y - 1)$ denotes the Łukasiewicz

t-norm. Analogously,

$$E'_r(x, y) = \exp(-\frac{1}{r}|x - y|)$$

is a $T_{\mathbf{P}}$ -equivalence on \mathbb{R} , where $T_{\mathbf{P}}(x, y) = xy$ denotes the product t-norm.¹

Based on a general representation theorem for strongly complete fuzzy orderings [4, Theorem 4.2], we can further prove that

$$L_r(x, y) = \min(1, \max(0, 1 - \frac{1}{r}(x - y)))$$

is a strongly complete $T_{\mathbf{L}}$ - E_r -ordering on \mathbb{R} and that

$$L'_r(x, y) = \min(1, \exp(-\frac{1}{r}(x - y)))$$

is a strongly complete $T_{\mathbf{P}}$ - E'_r -ordering on \mathbb{R} . As $T_{\mathbf{L}} \leq T_{\mathbf{P}}$, we can trivially conclude that L'_r is also a strongly complete $T_{\mathbf{L}}$ - E'_r -ordering.

In order to generalize the notion of concordant and discordant pairs, we need the notion of a strict fuzzy ordering. We call a binary fuzzy relation $R : X^2 \rightarrow [0, 1]$ a *strict fuzzy ordering* with respect to T and a T -equivalence $E : X^2 \rightarrow [0, 1]$, for brevity *strict T - E -ordering*, if it is irreflexive (i.e. $R(x, x) = 0$ for all $x \in X$), T -transitive, and E -extensional, that is,

$$T(E(x, x'), E(y, y'), R(x, y)) \leq R(x', y')$$

for all $x, x', y, y', z \in X$ [5].

Given a T - E -ordering $L : X^2 \rightarrow [0, 1]$,

$$R(x, y) = \min(L(x, y), N_T(L(y, x))), \quad (2)$$

where $N_T(x) = \sup\{y \in [0, 1] \mid T(x, y) = 0\}$ is the residual negation of T , is the most appropriate choice for extracting a strict fuzzy ordering from a given fuzzy ordering L (for a detailed argumentation, see [5]). From this construction, we can infer that the fuzzy relation

$$R_r(x, y) = \min(1, \max(0, \frac{1}{r}(y - x)))$$

is a strict $T_{\mathbf{L}}$ - E_r -ordering and that

$$R'_r(x, y) = \max(0, 1 - \exp(-\frac{1}{r}(y - x)))$$

is a strict $T_{\mathbf{L}}$ - E'_r -ordering.

If a given $T_{\mathbf{L}}$ - E -ordering $L : X^2 \rightarrow [0, 1]$ is strongly complete, it can be proved that the fuzzy relation R defined as in (2) simplifies to

$$R(x, y) = 1 - L(y, x)$$

and that the following holds:

$$R(x, y) + E(x, y) + R(y, x) = 1 \quad (3)$$

$$\min(R(x, y), R(y, x)) = 0 \quad (4)$$

¹In the following, we will further use the well-known *minimum t-norm* $T_{\mathbf{M}}(x, y) = \min(x, y)$.

5 A Fuzzy Ordering-Based Rank Correlation Coefficient

The previous section has provided us with the apparatus that is necessary to define a generalized rank correlation measure. Assume that the data are given as in (1) again (with $x_i \in X$ and $y_i \in Y$ for all $i = 1, \dots, n$). Further assume that we are given two $T_{\mathbf{L}}$ -equivalences $E_X : X^2 \rightarrow [0, 1]$ and $E_Y : Y^2 \rightarrow [0, 1]$, a strongly complete $T_{\mathbf{L}}$ - E_X -ordering $L_X : X^2 \rightarrow [0, 1]$ and a strongly complete $T_{\mathbf{L}}$ - E_Y -ordering $L_Y : Y^2 \rightarrow [0, 1]$. Then we can define a strict $T_{\mathbf{L}}$ - E_X -ordering on X as $R_X(x_1, x_2) = 1 - L_X(x_2, x_1)$ and a strict $T_{\mathbf{L}}$ - E_Y -ordering on Y as $R_Y(y_1, y_2) = 1 - L_Y(y_2, y_1)$.

Spearman's rho is based on rankings. Rankings are crisp concepts in which it is not easy to accommodate degrees of relationship in a straightforward way. Thus it is more meaningful to use pairwise comparisons to define a concept of rank correlation, just like Kendall's tau and the gamma measure do.

Given an index pair (i, j) , we can compute the degree to which (i, j) is a concordant pair as

$$\tilde{C}(i, j) = \bar{T}(R_X(x_i, x_j), R_Y(y_i, y_j))$$

and the degree to which (i, j) is a discordant pair as

$$\tilde{D}(i, j) = \bar{T}(R_X(x_i, x_j), R_Y(y_j, y_i)),$$

where \bar{T} is some t-norm to aggregate the relationships of x and y components.

It is easy to prove that, for all index pairs (i, j) , the equality

$$\tilde{C}(i, j) + \tilde{C}(j, i) + \tilde{D}(i, j) + \tilde{D}(j, i) + \tilde{T}(i, j) = 1 \quad (5)$$

holds. In this equation, $\tilde{T}(i, j)$ denotes the degree to which (i, j) is a tie in either variable

$$\tilde{T}(i, j) = \bar{S}(E_X(x_i, x_j), E_Y(y_i, y_j)),$$

where \bar{S} is the dual t-conorm of \bar{T} (i.e. $\bar{S}(x, y) = 1 - \bar{T}(1 - x, 1 - y)$). Note that (5) does not hold in general. The properties (3) and (4), however, are sufficient conditions for the fulfillment of (5).

If we adopt the simple sigma count idea to measure the cardinality of a fuzzy set [9], we can compute the numbers of concordant pairs \tilde{C} and discordant pairs \tilde{D} , respectively, as

$$\tilde{C} = \sum_{i=1}^n \sum_{j \neq i} \tilde{C}(i, j), \quad \tilde{D} = \sum_{i=1}^n \sum_{j \neq i} \tilde{D}(i, j).$$

The question arises whether we should attempt to generalize τ_a , τ_b or γ . As the main motivation is to get rid of the influence of close-to-ties pairs in the presence of noise, it is immediate that the idea behind γ is the most promising one. So, with the assumptions from above, we define our *fuzzy ordering-based rank correlation measure* $\tilde{\gamma}$ as

$$\tilde{\gamma} = \frac{\tilde{C} - \tilde{D}}{\tilde{C} + \tilde{D}}$$

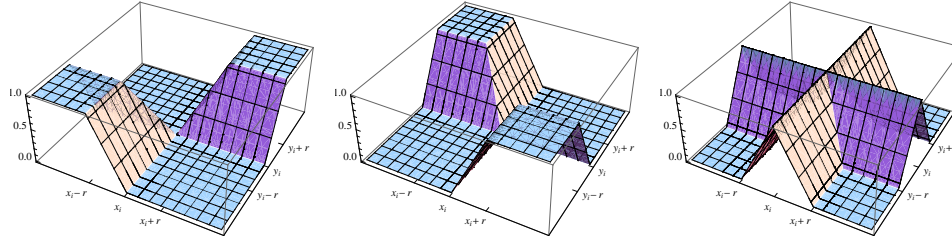


Figure 3: $\tilde{C}(i, j) + \tilde{C}(j, i)$ (left), $\tilde{D}(i, j) + \tilde{D}(j, i)$ (middle), and $\tilde{T}(i, j)$ (right) plotted as functions of x_j and y_j for fixed x_i and y_i (using the relations E_r and R_r and the minimum t-norm $\tilde{T} = T_M$ for aggregation).

Then we can also infer the following:

$$\tilde{C} = \sum_{i=1}^n \sum_{j>i} (\tilde{C}(i, j) + \tilde{C}(j, i)) \quad \tilde{D} = \sum_{i=1}^n \sum_{j>i} (\tilde{D}(i, j) + \tilde{D}(i, j))$$

Thus, by (5), $\tilde{C} + \tilde{D}$ equals the number of non-tie pairs if we consider each choice of indices i, j only once (in contrast to considering (i, j) and (j, i) independently for each i and j). So $\tilde{\gamma}$ measures the difference of concordant and discordant pairs relative to the number of non-tie pairs; the concept of “tiedness” is a fuzzy one, however.

It is obvious that, in case that E_X and E_Y are crisp equalities and that R_X and R_Y are crisp linear strict orderings, $\tilde{\gamma}$ coincides with γ . So what is the difference if R_X and R_Y are non-trivial fuzzy relations? We will see shortly that concordant/discordant pairs are counted more if they are dissimilar and less if they are similar—which perfectly corresponds to our intention. Let us demonstrate this fact with an example.

Assume $X = Y = \mathbb{R}$, $E_X = E_Y = E_r$, and $R_X = R_Y = R_r$ for some $r > 0$. Fixing some x_i and y_i and considering $\tilde{C}(i, j) + \tilde{C}(j, i)$, $\tilde{D}(i, j) + \tilde{D}(j, i)$, and $\tilde{T}(i, j)$ as functions of the two variables x_j and y_j , the graphs shown in Figure 3 can be obtained. It can be seen that pairs are counted fully if $|x_i - x_j| > r$ and $|y_i - y_j| > r$ (i.e. like in the classical γ measure). If one of the two distances is smaller than r , the pair is considered as a tie to the corresponding degree $\tilde{T}(i, j)$ and only counted to a degree of $1 - \tilde{T}(i, j)$. If the relations $E_X = E_Y = E'_r$, and $R_X = R_Y = R'_r$ are used, the effect is qualitatively similar, r also controls to which degree a close-to-tie pair is counted, also in a monotonic, yet asymptotic fashion (see Figure 4). It is obvious that the fuzzy sets depicted in Figures 3 and 4 are nothing else but continuous fuzzifications of the crisp partition depicted in Figure 2.

It is clear from the above examples that, the smaller r , the more $\tilde{\gamma}$ resembles to γ . For both, the variant based on E_r/R_r and the variant based on E'_r/R'_r , it can be proved that $\tilde{\gamma}$ converges to γ for $r \rightarrow 0$. One also sees that, if r is chosen so large that $|x_i - x_j| \leq r$ and $|y_i - y_j| \leq r$ for all pairs, $\tilde{\gamma}$ based on R_r with $\tilde{T} = T_M$ counts all pairs to a degree proportionally to the minimum of these two distances—which is a meaningful rank correlation coefficient, too. For $\tilde{T} = T_P$, $\tilde{\gamma}$ based on R_r counts

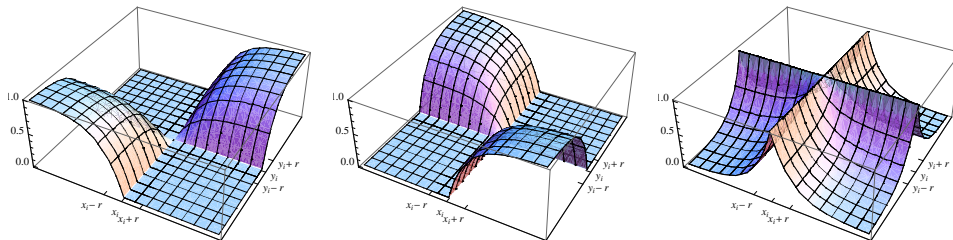


Figure 4: $\tilde{C}(i, j) + \tilde{C}(j, i)$ (left), $\tilde{D}(i, j) + \tilde{D}(j, i)$ (middle), and $\tilde{T}(i, j)$ (right) plotted as functions of x_j and y_j for fixed x_i and y_i (using the relations E'_r and R'_r and the product t-norm $\tilde{T} = T_{\mathbf{P}}$ for aggregation).

pairs to a degree proportionally to the product of the two distances. With more effort, it is possible to prove that the R'_r -based variants of $\tilde{\gamma}$, with $r \rightarrow \infty$, converge to the same limit values as the R_r -based variants.

Another property of $\tilde{\gamma}$ is immediate to see: if the fuzzy relations R_X and R_Y are continuous and if \tilde{T} is a continuous t-norm, then $\tilde{\gamma}$ depends continuously on the data set $(x_i, y_i)_{i=1}^n$.

6 Experiments

Let us first reconsider the example from Section 3. More specifically, we are given 100 uniformly distributed random values (x_1, \dots, x_{100}) from the unit interval. The list (y_1, \dots, y_{100}) is computed as $y_i = f(x_i)$, where f is a simple, piecewise linear, non-decreasing function that has a relatively large flat area. In order to study how different rank correlation measures react to noise, we contaminated the data points with additive, independent, normally distributed noise with 0 mean and standard deviation σ . Figure 5 shows these data sets. Figure 6 displays the results that we obtained for different rank correlation measures. We compared ρ , τ_b , γ and different variants of $\tilde{\gamma}$. Every line in Figure 6 corresponds to the results obtained by one rank correlation measure depending on the noise level σ . The two lines for τ_b (dotted, black) and γ (dotted, light gray) coincide except for no noise ($\sigma = 0$). Both lines reveal that these two measures react to noise in a non-robust way. More or less the same is true for ρ (dotted, medium gray). The other lines correspond to different variants of $\tilde{\gamma}$. Solid lines correspond to $\tilde{\gamma}$ using R_r and dashed lines denote the results for $\tilde{\gamma}$ using R'_r (where we use the same r for both components and $\tilde{T} = T_{\mathbf{M}}$). We used $r = 0.05$ (black), $r = 0.2$ (medium gray), and $r = 0.5$ (light gray). We see that all six different variants react to the noise in a more robust way than the three crisp measures. Clearly, the higher r , the more noise is neglected. Note, however, that, the larger r , the more difficult it is for $\tilde{\gamma}$ to find out whether there are slightly non-monotonic parts in the data.

So let us consider a different setting. Now we fix the noise level $\sigma = 0.01$ and use different functions to create the second list (y_1, \dots, y_{100}) . Right of $x = 0.5$, we use $f(x) = \frac{x}{2} + \frac{1}{4}$ and, to the left of $x = 0.5$, we linearly interpolate between $(0, q)$

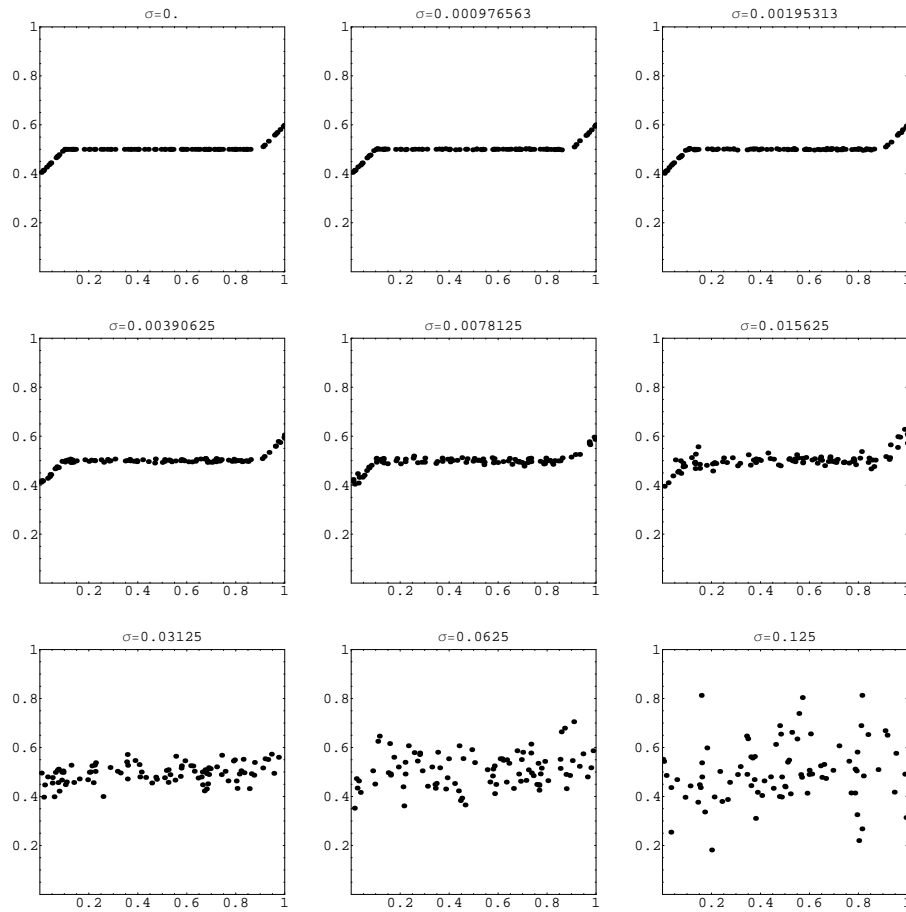


Figure 5: Different data sets obtained from contaminating a non-decreasing relationship by normally distributed noise with different standard deviations.

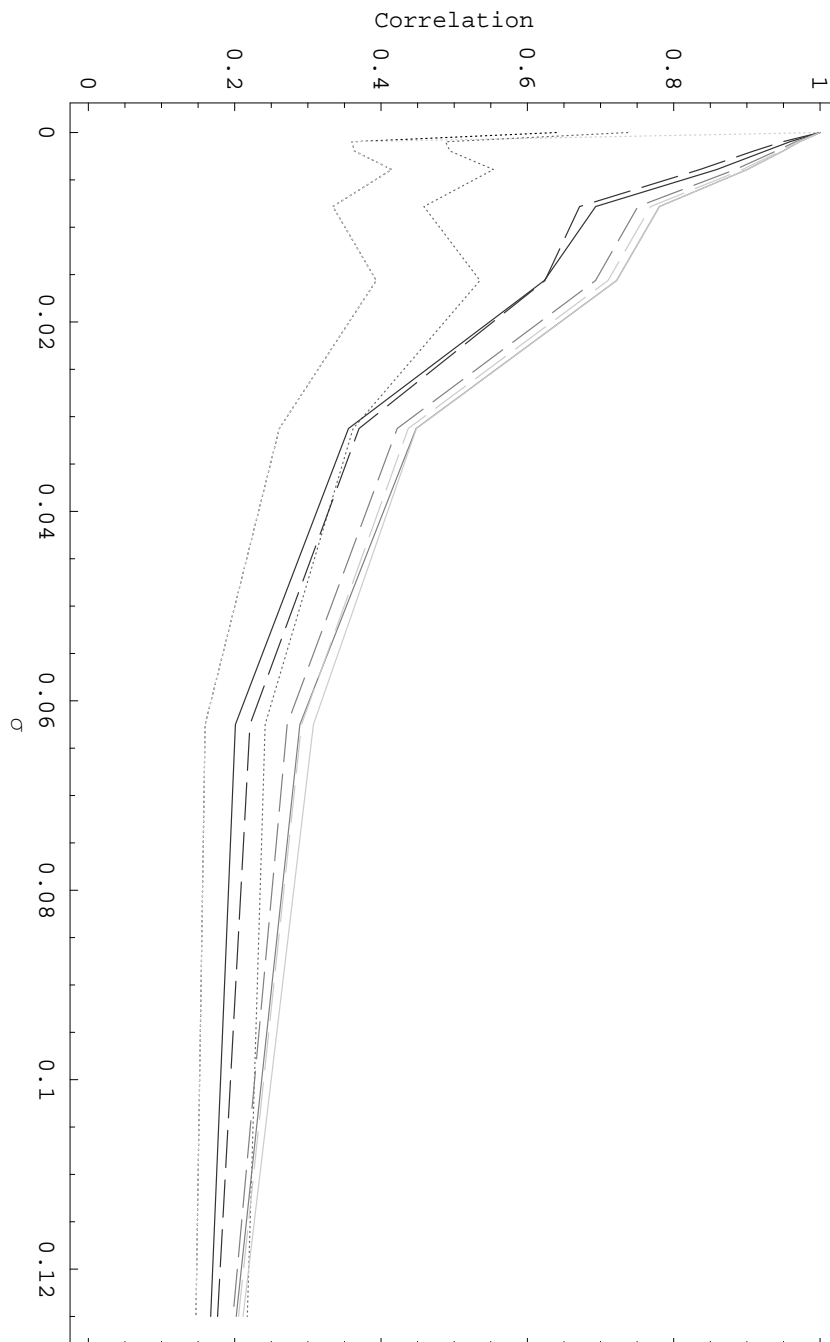


Figure 6: Results obtained by applying different rank correlation measures to the data sets shown in Figure 5.

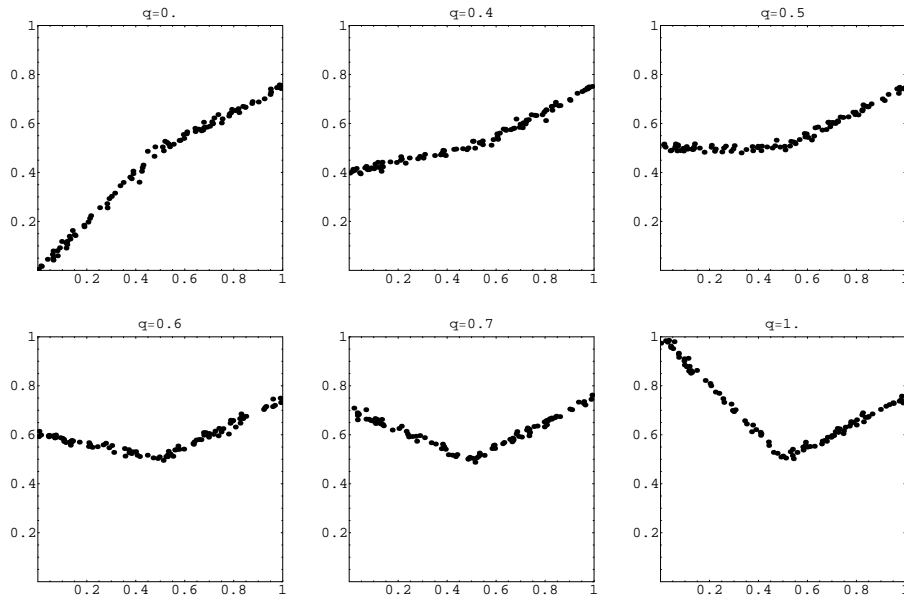


Figure 7: Noisy data sets that correspond to monotonic ($q \leq 0.5$) and non-monotonic relationships ($q > 0.5$).

and $(0.5, 0.5)$. It is clear, that this relationship is monotonic if and only if $q \leq 0.5$. The data sets are displayed in Figure 7 and the results are presented in Figure 8, where we use the same conventions to distinguish the lines as in Figure 6. We see that all variants of $\tilde{\gamma}$ show acceptable results for $q \leq 0.5$, whereas ρ , τ_b and γ again have problems to handle the noise in case of the large proportion of ties that occurs for $q = 0.5$. We also see that $\tilde{\gamma}$ already yields significantly lower values for $q = 0.6$ in the case $r = 0.05$ (no matter which of the two variants is considered). For larger r , however, we see that $\tilde{\gamma}$ cannot detect the slight non-monotonicity for $q = 0.6$ that well. These two examples demonstrate that, when choosing r , there is a trade-off between robustness (the larger r , the better) and sensitivity (the smaller r , the better).

As a third set of experiments, we have tried to figure out the variance of $\tilde{\gamma}$. For this study, we have computed all rank correlation measures used in the above experiments for different test data several times and computed the variance of the results. In all experiments, we have encountered that τ_b and γ had higher variances than all variants of $\tilde{\gamma}$. The variances we obtained for different variants of $\tilde{\gamma}$ obeyed a simple and unsurprising rule: the larger r , the smaller the variance. Interestingly, the variances we obtained for Spearman's ρ were also very low, comparable to the lower values for $\tilde{\gamma}$ with a large r .

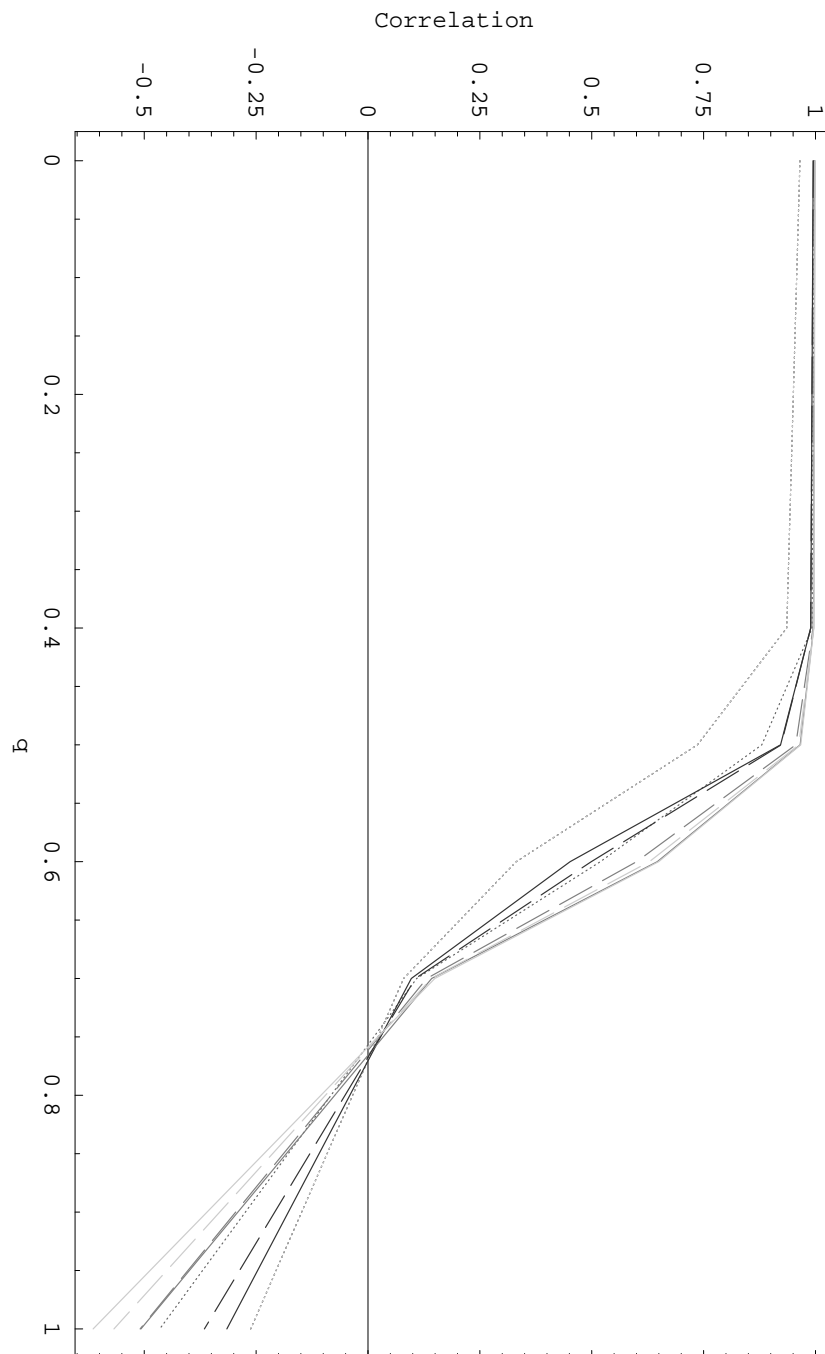


Figure 8: Results obtained by applying different rank correlation measures to the data sets shown in Figure 7.

7 Concluding Remarks

This paper, as the title suggests, attempts to present initial ideas that the authors consider promising. The examples of the previous section are intended to support this viewpoint. They are illustrative and indicative, but they cannot replace a formal investigation of the properties of $\tilde{\gamma}$. As it has been done exhaustively for Spearman's rho and Kendall's tau, a significance analysis and a variance analysis have to be carried out. Note, however, that this cannot be done analogously for $\tilde{\gamma}$. Both Spearman's rho and Kendall's tau are fully determined by the ranking of the lists (x_1, \dots, x_n) and (y_1, \dots, y_n) . Thus, combinatorial techniques can be used to study variances and significance levels [13]—not so for $\tilde{\gamma}$ that always depends on the distance relationships of the values, too, so this analysis can only be done by some distribution assumptions. These studies are left to future research.

To determine the right choice for the parameter r is another open question. As we have noted above, there is a trade-off between robustness on the one side and sensitivity/significance on the other side. So this topic goes hand in hand with a more formal statistical analysis. Profound results concerning the choice of r , again, can only be expected with specific distribution assumptions. In any case, we want to note in advance that $\tilde{\gamma}$ depends continuously on r , so at least we can be sure that $\tilde{\gamma}$ will react robustly to slightly sub-optimal choices of r .

We would like to remark that this investigation was inspired by a problem in bioinformatics: how to infer sets of co-transcribed genes in procaryotic genomes (so-called *operons*) from the gene expression levels measured by microarray experiments [3, 19, 22]. It will also be subject of future research to evaluate the rank correlation measures introduced in this paper in this domain. First experiments have been very promising.

References

- [1] H. Abdi. Coefficients of correlation, alienation and determination. In N. J. Salkind, editor, *Encyclopedia of Measurement and Statistics*. Sage, Thousand Oaks, CA, 2007.
- [2] H. Abdi. The Kendall rank correlation coefficient. In N. J. Salkind, editor, *Encyclopedia of Measurement and Statistics*. Sage, Thousand Oaks, CA, 2007.
- [3] J. Bockhorst, Y. Qiu, J. Glasner, M. Liu, F. Blattner, and M. Craven. Predicting bacterial transcription units using sequence and expression data. *Bioinformatics*, 19(Suppl. 1):i34–i43, 2003.
- [4] U. Bodenhofer. A similarity-based generalization of fuzzy orderings preserving the classical axioms. *Internat. J. Uncertain. Fuzziness Knowledge-Based Systems*, 8(5):593–610, 2000.
- [5] U. Bodenhofer and M. Demirci. Strict fuzzy orderings with a given context of similarity. *Internat. J. Uncertain. Fuzziness Knowledge-Based Systems*, 16(2):147–178, 2008.

- [6] D. Boixader, J. Jacas, and J. Recasens. Fuzzy equivalence relations: advanced material. In D. Dubois and H. Prade, editors, *Fundamentals of Fuzzy Sets*, volume 7 of *The Handbooks of Fuzzy Sets*, pages 261–290. Kluwer Academic Publishers, Boston, 2000.
- [7] B. De Baets and R. Mesiar. Pseudo-metrics and T -equivalences. *J. Fuzzy Math.*, 5(2):471–481, 1997.
- [8] B. De Baets and R. Mesiar. Metrics and T -equalities. *J. Math. Anal. Appl.*, 267:331–347, 2002.
- [9] A. DeLuca and S. Termini. A definition of a nonprobabilistic entropy in the setting of fuzzy sets theory. *Inf. Control*, 20:301–312, 1972.
- [10] J. Fodor and M. Roubens. *Fuzzy Preference Modelling and Multicriteria Decision Support*. Kluwer Academic Publishers, Dordrecht, 1994.
- [11] L. A. Goodman and W. H. Kruskal. Measures of association for cross classifications. *J. Amer. Statist. Assoc.*, 49(268):732–764, 1954.
- [12] M. G. Kendall. A new measure of rank correlation. *Biometrika*, 30:81–93, 1938.
- [13] M. G. Kendall. *Rank Correlation Methods*. Charles Griffin & Co., London, third edition, 1962.
- [14] F. Klawonn. Fuzzy sets and vague environments. *Fuzzy Sets and Systems*, 66:207–221, 1994.
- [15] E. P. Klement, R. Mesiar, and E. Pap. *Triangular Norms*, volume 8 of *Trends in Logic*. Kluwer Academic Publishers, Dordrecht, 2000.
- [16] W. H. Kruskal. Ordinal measures of association. *J. Amer. Statist. Assoc.*, 53(284):814–861, 1958.
- [17] S. V. Ovchinnikov. An introduction to fuzzy relations. In D. Dubois and H. Prade, editors, *Fundamentals of Fuzzy Sets*, volume 7 of *The Handbooks of Fuzzy Sets*, pages 233–259. Kluwer Academic Publishers, Boston, 2000.
- [18] K. Pearson. Notes on the history of correlation. *Biometrika*, 13:25–45, 1920.
- [19] C. Sabatti, L. Rohlin, M.-K. Oh, and J. C. Liao. Co-expression pattern from DNA microarray experiments as a tool for operon prediction. *Nucleic Acids Res.*, 30(13):2886–2893, 2002.
- [20] C. Spearman. The proof and measurement of association between two things. *Am. J. Psychol.*, 15(1):72–101, 1904.
- [21] C. Spearman. Demonstration of formulae for true measurement of correlation. *Am. J. Psychol.*, 18(2):161–169, 1907.

- [22] D. Steinhauser, B. H. Junker, A. Luedemann, J. Selbig, and J. Kopka. Hypothesis-driven approach to predict transcriptional units from gene expression data. *Bioinformatics*, 20(12):1928–1939, 2004.
- [23] E. Trillas and L. Valverde. An inquiry into indistinguishability operators. In H. J. Skala, S. Termini, and E. Trillas, editors, *Aspects of Vagueness*, pages 231–256. Reidel, Dordrecht, 1984.