

# Selection of Relevant and Non-Redundant Multivariate Ordinal Patterns for Time Series Classification - Supplementary Material

Arvind Kumar Shekar

Marcus Pappik

Patricia Iglesias Sanchez

Emmanuel Mueller

## 1 Experiments

**1.1 Comparison of other information theoretic relevance measures** In the work of *ordex*, we introduced a novel scoring method for ordinal patterns. We compare bivariate mutual information and multivariate KL-divergence based relevance measure [1] against our proposed scoring. We evaluate the run times and the test data accuracy on 18 UCR and 2 UCI datasets.

We observe that our relevance function, based on the Chebyshev-Inequality, performs better in comparison to the KL-divergence in the context of our algorithm. For 70% of the data sets, using our relevance measure yields a better accuracy. On average, the achieved accuracy was 2.77% better than with KLD. In all cases, our relevance measure outperforms the KLD in terms of run time. Even if both run times seem to grow quite similar with respect to the data set properties, the KLD needed 4.77% more time on average.

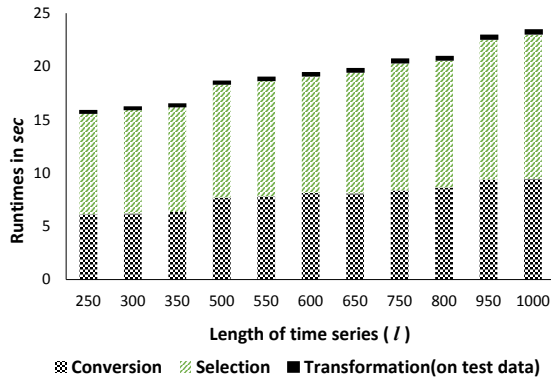


Figure 1: Scalability w.r.t. increasing  $l$ , where  $m=5$ ,  $n=100$ ,  $d=5$ ,  $o=10$  and  $I=200$

By comparing our Chebyshev-based relevance function with Mutual Information, our relevance measure outperformed MI on 75% of the data sets in terms of accuracy and on 80% of the data sets in terms of standard deviation. Also our scoring outperformed on all data sets w.r.t. the run time. Using MI, the average run time was 18.09% higher for the selection phase.

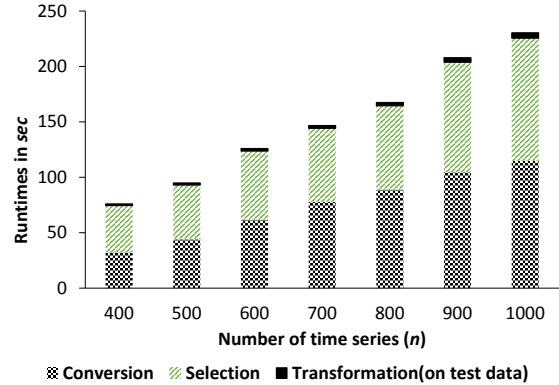


Figure 2: Scalability w.r.t increasing  $n$ , where  $m=5$ ,  $l=600$ ,  $d=5$ ,  $o=10$  and  $I=200$

**1.2 Scalability Experiments** Figure 1 and Figure 2 show the scalability of *ordex* with increasing length of time series and number of samples in a dataset.

## 2 Upper bound of mis-classification

Consider  $f$  as the feature extracted using  $s$  to classify classes  $c_a$  and  $c_b$ .  $f|_{c_a}$  is the distribution of the extracted feature for class  $c_a$  and  $E[f|_{c_a}]$  and  $Var[f|_{c_a}]$  are the expected value and the variance of the distribution. Similarly, for class  $c_b$ , we define the distribution  $f|_{c_b}$ , expected value  $E[f|_{c_b}]$  and variance  $Var[f|_{c_b}]$ . Without loss of generality, we assume  $E[f|_{c_a}] < E[f|_{c_b}]$ .  $P(M_k^{c_a, c_b})$  denotes the probability that  $c_a$  is mis-classified as  $c_b$  or  $c_b$  is misclassified as  $c_a$ .

Following the rule of Chebychev's inequality [2], a sample is classified as  $c_a$  or  $c_b$  based on the arbitrary threshold value  $k$  |  $0 < k < |E[f|_{c_b}] - E[f|_{c_a}]|$ . Under the assumption that  $f|_{c_a}$  and  $f|_{c_b}$  are symmetrically distributed around their expected values, a sample is classified as  $c_b$  when its expected value is greater than  $E[f|_{c_a}] + k$ . Hence, to estimate  $P(M_k^{c_a, c_b})$  we need to quantify the maximum number of  $c_a$  that exceed the threshold and likewise for  $c_b$ . The upper bound of mis-

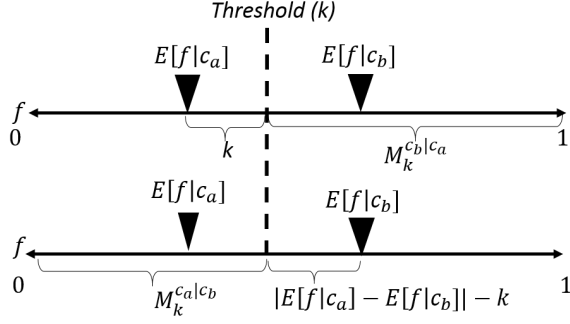


Figure 3: Example: A number line with limits  $[0,1]$

classification is represented as,

$$P(M_k^{c_a, c_b}) \leq \frac{Var[f|c_a]}{2k^2} + \frac{Var[f|c_b]}{2(|E[f|c_b] - E[f|c_a]| - k)^2}.$$

**2.1 Proof** The upper bound of mis-classification is strongly founded by the principles of Chebyshev-inequality. For a random variable  $X$  with arbitrary distribution and  $a > 0$ , the inequality states,

$$(2.1) \quad P(|X - E[X]| \geq a) \leq \frac{Var[X]}{a^2}.$$

Given that the expected value  $E[X]$  and variance  $Var[X]$  of the random variable, the inequality represents the probability of a sample that is greater than  $a$ . The approach is commonly used for finding outliers, i.e., instances with a high probability of being greater than  $E[X] + a$  are outliers.

For a classification task, using our threshold  $E[f|c_a] + k$ , we aim to estimate:

$M_k^{c_b|c_a}$ : a sample with class  $c_a$  is misclassified as  $c_b$  based on this threshold.

$M_k^{c_a|c_b}$ : a sample with class  $c_b$  is misclassified as  $c_a$  based on this threshold.

Figure 3 visualizes both cases on a simple number-line, by applying the Chebychev inequality,

$$\begin{aligned} P(M_k^{c_b|c_a}) &= P((f|c_a) \geq E[f|c_a] + k) \\ &= P((f|c_a) - E[f|c_a] \geq k). \end{aligned}$$

We assumed that  $f|c_a$  is distributed symmetrically around its expected value. Thus  $P(M_k^{c_b|c_a})$  is half the probability that the value of  $f|c_a$  has at least a distance of  $k$  to its expected value,

$$P(M_k^{c_b|c_a}) = \frac{1}{2}P(|(f|c_a) - E[f|c_a]| \geq k).$$

Using the Chebyshev-Inequality and setting  $a = k$  in Eq 2.1, we can estimate an upper bound of  $P(M_k^{c_b|c_a})$

as,

$$(2.2) \quad P(M_k^{c_b|c_a}) \leq \frac{Var[f|c_a]}{2k^2}.$$

On applying the same symmetric assumption on  $(f|c_b)$ ,  $M_k^{c_a|c_b}$  is half of the probability, that values of  $(f|c_b)$  are at least  $|E[f|c_b] - E[f|c_a]| - k$  away from their expected value

$$\begin{aligned} P(M_k^{c_a|c_b}) &= P((f|c_b) \leq E[f|c_a] + k) \\ &= P((f|c_b) - E[f|c_b] \leq E[f|c_a] - E[f|c_b] + k) \\ &= P(E[f|c_b] - (f|c_b) \geq E[f|c_b] - E[f|c_a] - k) \\ &= P(E[f|c_b] - (f|c_b) \geq |E[f|c_b] - E[f|c_a]| - k) \\ &= \frac{1}{2}P(|(f|c_b) - E[f|c_b]| \geq |E[f|c_b] \\ &\quad - E[f|c_a]| - k) \end{aligned}$$

Comparing the above result with Eq 2.1, we derive,  $a = |E[f|c_b] - E[f|c_a]| - k$ . To estimate an upper bound

$$(2.3) \quad P(M_k^{c_a|c_b}) \leq \frac{Var[f|c_b]}{2(|E[f|c_b] - E[f|c_a]| - k)^2}$$

From Eq 2.2 and 2.3 this, we derive an upper bound for the total mis-classification probability for  $c_a$  and  $c_b$  as,

$$\begin{aligned} P(M_k^{c_a, c_b}) &= P(M_k^{c_b|c_a} \cup M_k^{c_a|c_b}) \\ &\leq P(M_k^{c_b|c_a}) + P(M_k^{c_a|c_b}) \\ &\leq \frac{Var[f|c_a]}{2k^2} + \frac{Var[f|c_b]}{2(|E[f|c_b] - E[f|c_a]| - k)^2}. \end{aligned}$$

This means, given an optimal value of the threshold  $E[f|c_a] + k$ , we can calculate an upper bound for the minimal misclassification of each pair of classes. However, we cannot assume that all classifiers find such an optimal  $k$ , based on the data. More over, finding this bound costs additional computation time. In order to be independent of the classifier and have a better efficiency, we use the fact, that  $0 < k < |E[f|c_b] - E[f|c_a]|$ . Due to that, the upper bound of the mis-classification grows approximately as fast as

$$\frac{Var[f|c_a] + Var[f|c_b]}{2(|E[f|c_b] - E[f|c_a]|)^2}.$$

### 3 Parameters for real world dataset

Table 1: Real world data experiment parameter settings

Dataset	$d$	$I$	$m'$	$\alpha$	$o$
EMG limb sen	5	200	3	0.1	10
EMG limb pie	5	200	3	0.1	10
EMG limb mar	5	200	3	0.1	10
Character	5	300	3	0.1	50
Activity recognition	5	100	2	0.3	20
User Movement	5	300	3	0.8	30
Occupancy	5	100	3	0.5	10
Bosch	5	200	3	0.1	10

### References

- [1] A. K. Shekar, T. Bocklisch, C. N. Straehle, P. I. Sánchez, and E. Mller, “Including multi-feature interactions and redundancy for feature ranking in mixed datasets,” in *Machine Learning and Knowledge Discovery in Databases - European Conference, ECML PKDD 2017, Macedonia, Skopje, September 18-22, 2017, Proceedings*, ser. Lecture Notes in Computer Science. Springer, 2017.
- [2] S. Karlin and W. J. Studden, *Tchebycheff systems: With applications in analysis and statistics*. Interscience New York, 1966.