



HUN-REN
Magyar Kutatási Hálózat

Extended Similarity Indices: Benefits of Comparing more than Two Objects Simultaneously. Theory, Speed, Consistency, and Diversity Selection

*R. A. Miranda-Quintana¹, D. Bajusz², A.
Rácz², K. Héberger²*

¹Univeristy of Florida

*² HUN-REN Research Centre for Natural
Sciences, Institute of Excellence,
Hungarian Academy of Sciences*

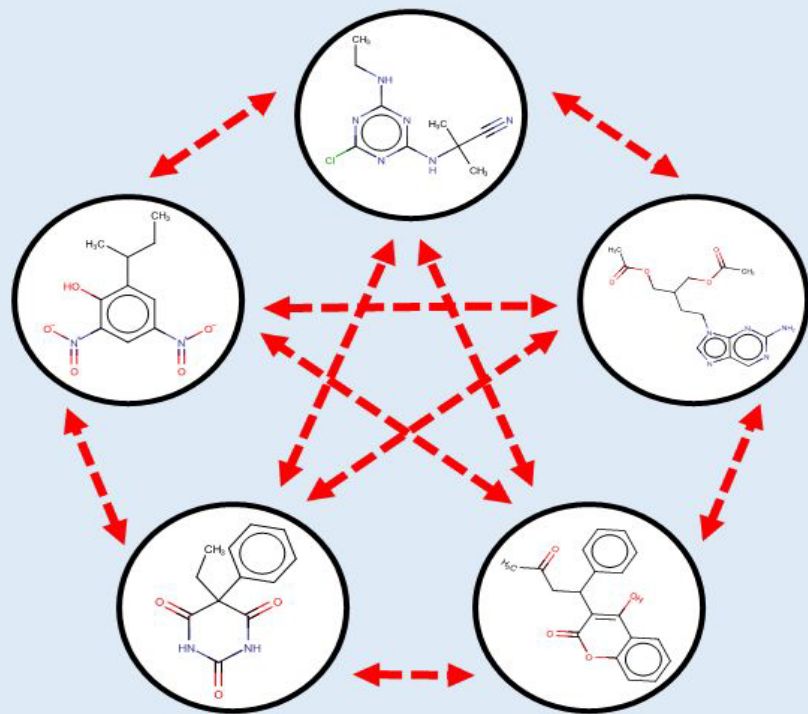
The quote

“None of the axioms employed by great generals [scientist] is difficult. Indeed, once they have been employed successfully, they reveal their innate simplicity and appear to be the obvious and sometimes only logical solution. Yet all great ideas are simple. The trick is to see them before others.”

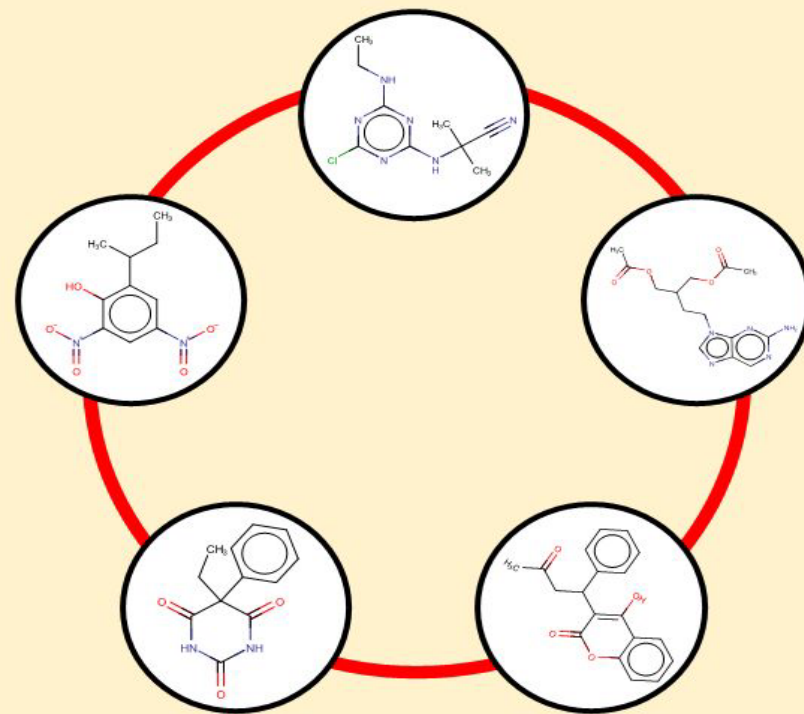
Bevin Alexander, How Great Generals Win, 1993

← COMPUTE TIME

Binary similarity



Extended (n -ary) similarity



DIVERSITY →

Definitions, terms

- **Nominal scale — binary.**
- **Fingerprint: a bit string consisting of 0-s and 1-s.**
- **Sum of rank differences (SRD): a city block (Manhattan) distance of ranks.**
- **ANOVA (a classification procedure, comparison of means)**
- **Similarity of molecules (pairwise comparison) a number:**
 $S [0,1]$ $r[-1,1]$ és $S=1-d$
- **Distance (Dissimilarity, d): $D[0,+\infty]$;**
 $S(x, y) = 1/(1 + D(x, y))$
- **Maximum likelihood principle & central limit theorem**

RESEARCH ARTICLE

Open Access

Why is Tanimoto index an appropriate choice for fingerprint-based similarity calculations?



Dávid Bajusz¹, Anita Rácz^{2,3} and Károly Héberger^{2*}

Rácz et al. *J Cheminform* (2018) 10:48
<https://doi.org/10.1186/s13321-018-0302-y>

Journal of Cheminformatics

RESEARCH ARTICLE

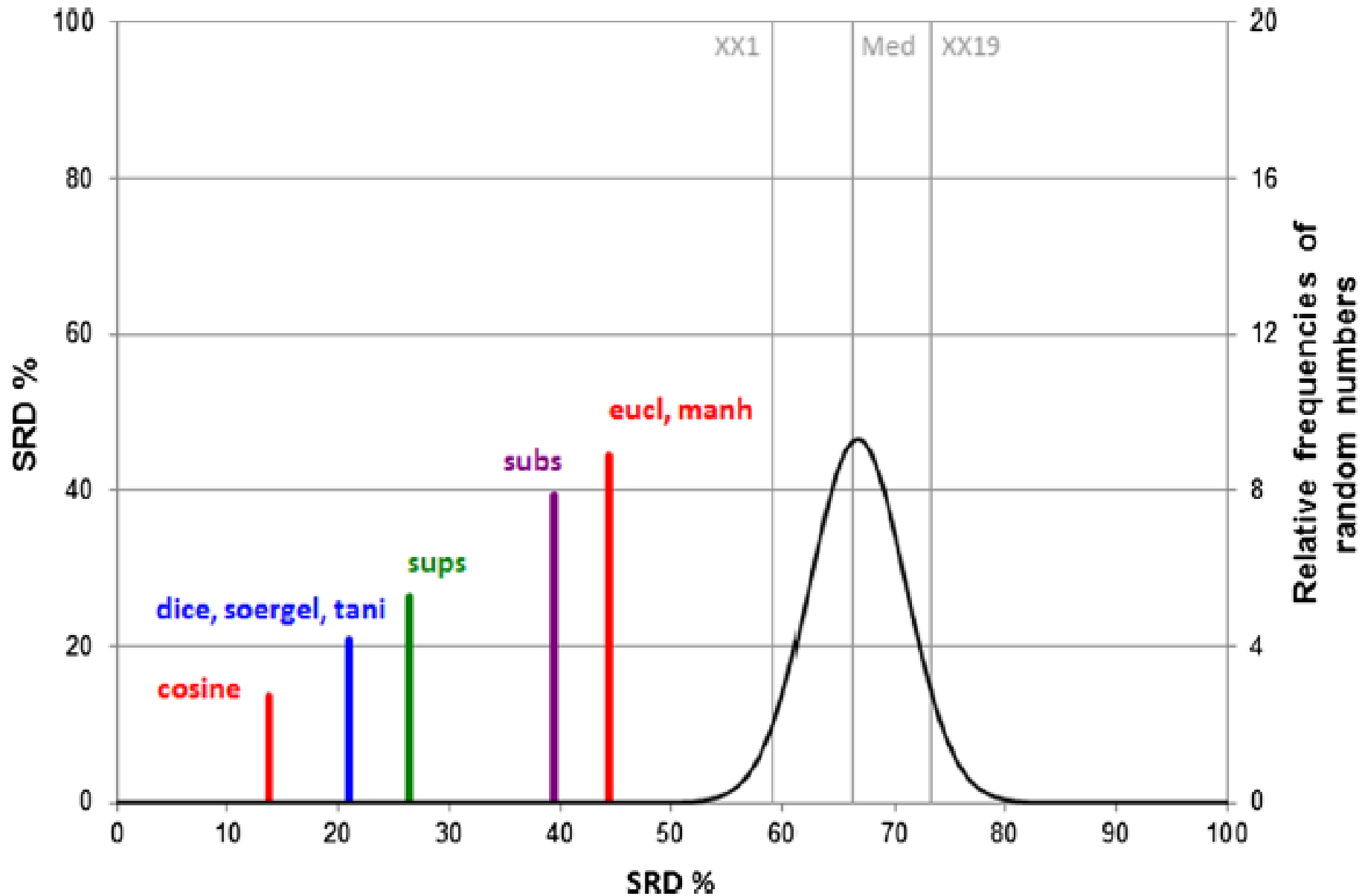
Open Access

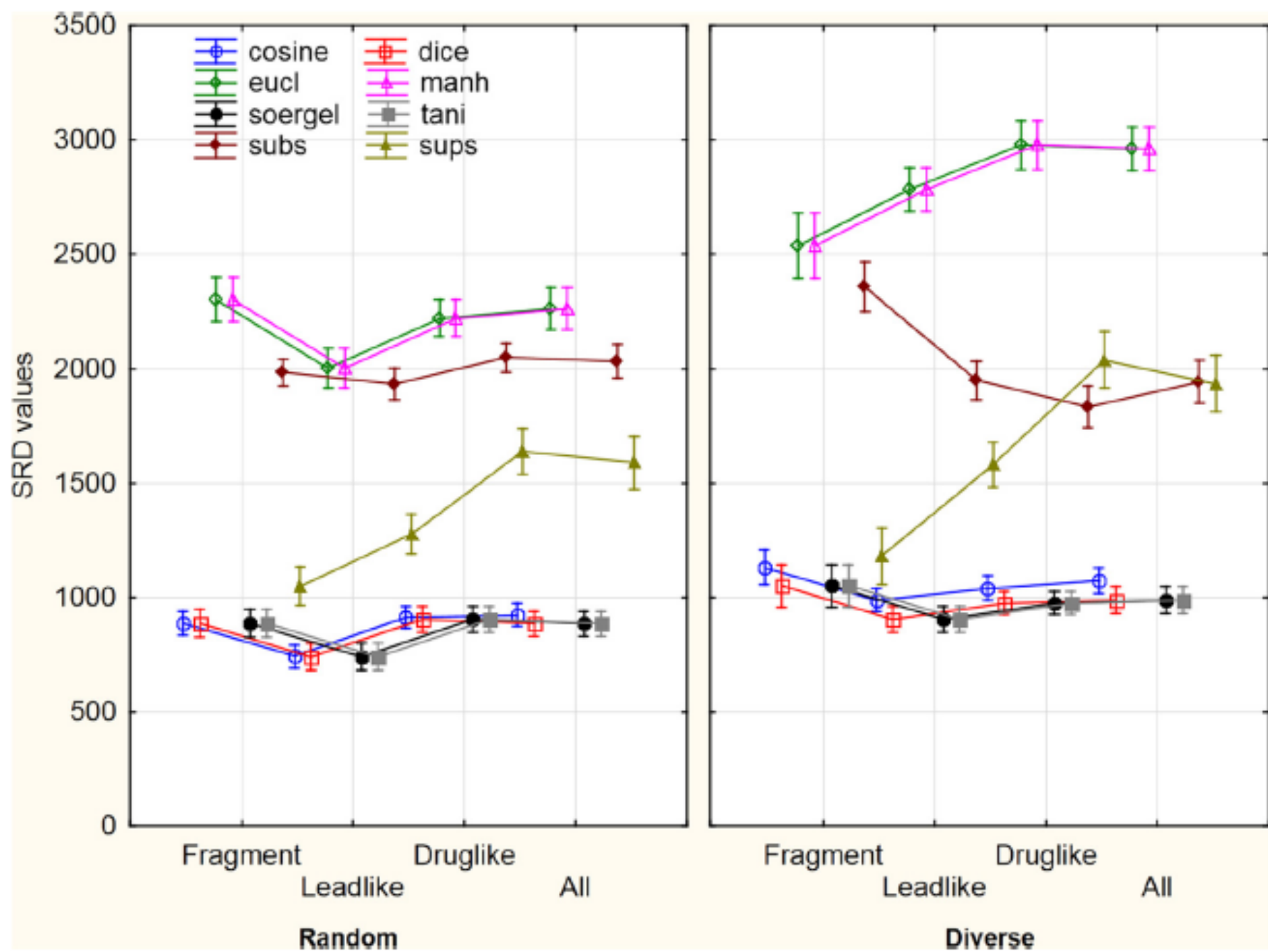
Life beyond the Tanimoto coefficient: similarity measures for interaction fingerprints



Anita Rácz¹ , Dávid Bajusz^{2*}  and Károly Héberger¹ 

Ranking of similarity measures





RESEARCH ARTICLE

Open Access

Why is Tanimoto index an appropriate choice for fingerprint-based similarity calculations?



Dávid Bajusz¹, Anita Rácz^{2,3} and Károly Héberger^{2*}

Rácz et al. *J Cheminform* (2018) 10:48
<https://doi.org/10.1186/s13321-018-0302-y>

Journal of Cheminformatics

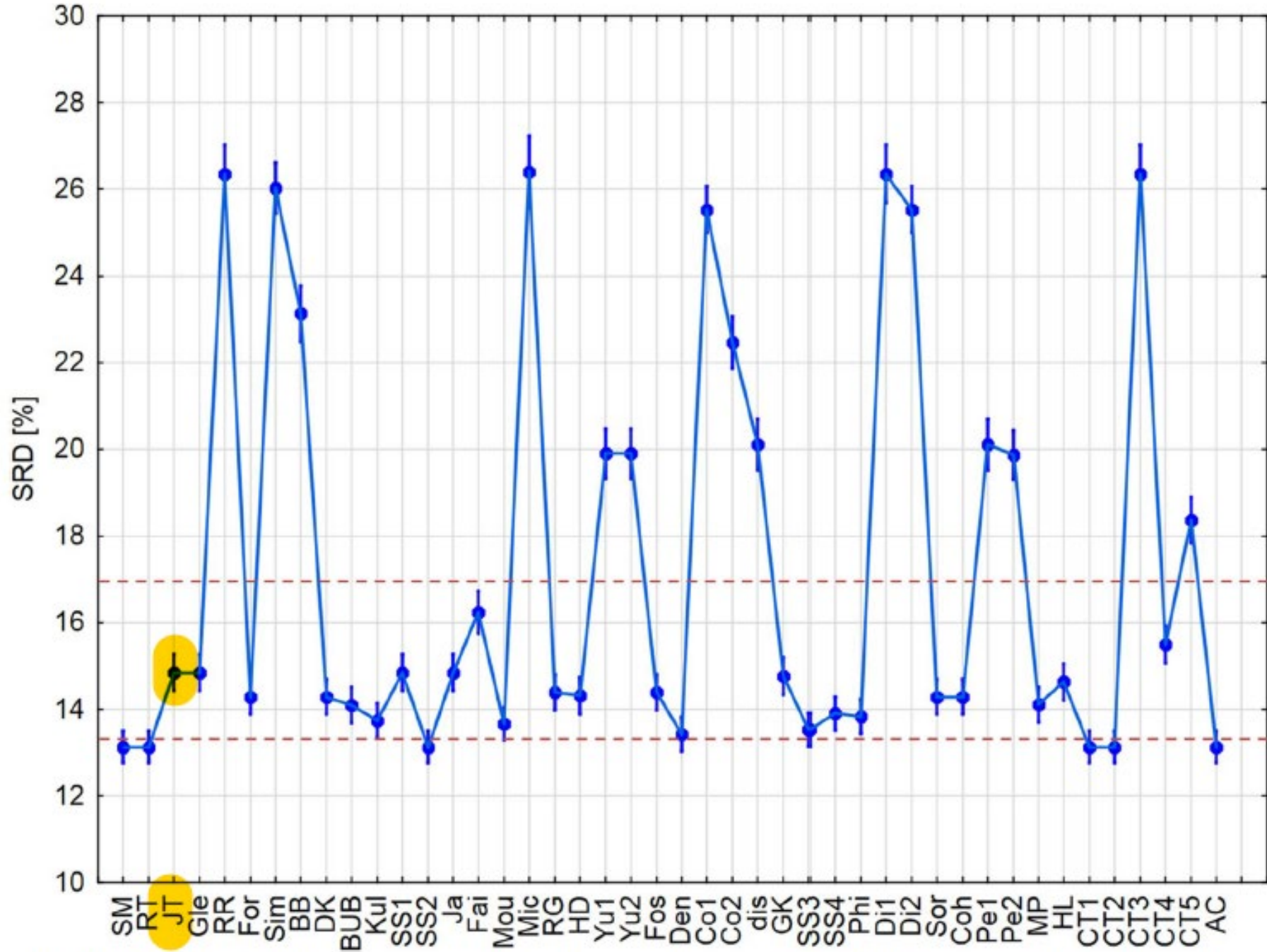
RESEARCH ARTICLE

Open Access

Life beyond the Tanimoto coefficient: similarity measures for interaction fingerprints



Anita Rácz¹ , Dávid Bajusz^{2*}  and Károly Héberger¹ 



What is the basic question of research?

- Is it possible to express the similarity of molecules better?
- To create faster algorithms?
- Well, and if so, how?

Representation of molecules

- It is not unambiguous how to express the **similarity** of two molecules.
- **Bitwise representation** (encoding of representative information) has proved **successful**.
- The bit string of a molecule is called a **fingerprint**
- The similarity of the two bands shows the similarity of the molecules, e.g. **Tanimoto coefficient**

A	<table border="1"><tr><td>1</td><td>0</td><td>1</td><td>1</td><td>0</td><td>1</td></tr></table>	1	0	1	1	0	1	$ A = 4$
1	0	1	1	0	1			
B	<table border="1"><tr><td>1</td><td>1</td><td>0</td><td>1</td><td>0</td><td>0</td></tr></table>	1	1	0	1	0	0	$ B = 3$
1	1	0	1	0	0			
$A \wedge B$	<table border="1"><tr><td>1</td><td>0</td><td>0</td><td>1</td><td>0</td><td>0</td></tr></table>	1	0	0	1	0	0	$ A \wedge B = 2$
1	0	0	1	0	0			
$A \vee B$	<table border="1"><tr><td>1</td><td>1</td><td>1</td><td>1</td><td>0</td><td>1</td></tr></table>	1	1	1	1	0	1	$ A \vee B = 5$
1	1	1	1	0	1			

$$S_T(A, B) = \frac{2}{5}$$

BV1	BV2		Freq.
1	1	<i>a</i>	3
0	1	<i>b</i>	4
1	0	<i>c</i>	2
0	0	<i>d</i>	1

Sokal – Mitchener

$$S = \frac{a + d}{a + b + c + d}$$

Jaccard – Tanimoto

$$S = \frac{a}{a + b + c}$$

1	0
0	1
0	1
1	1
...	...
0	1
1	1

Contingency table

$p = a + b + c + d$		2. molecule	
		1 (substructure present)	0 (no substructure)
1. molecule	1 (substructure present)	a	b
	0 (no substructure)	c	d

Extended Sokal-Michener index

$$F_1 = (1 \ 0 \ 1 \ 1 \ 0 \ 1 \ 0 \ 0)$$

$$F_2 = (0 \ 0 \ 1 \ 0 \ 0 \ 1 \ 0 \ 1)$$

- $F_3 = (1 \ 0 \ 1 \ 1 \ 1 \ 0 \ 0 \ 1)$

$$F_4 = (0 \ 0 \ 1 \ 1 \ 0 \ 1 \ 0 \ 0)$$

$$F_5 = (0 \ 0 \ 1 \ 1 \ 0 \ 1 \ 1 \ 0)$$

$$S = \frac{a+d}{a+b+c+d}$$

- $C_{5(5)} = 1; C_{5(4)} = 2; C_{5(3)} = 0; C_{5(2)} = 2;$

$$C_{5(1)} = 2; C_{5(0)} = 1$$

$$\Delta_{5(k)} = |2k - 5|$$

Similarity=1 (if $2k - n > \gamma$), Similarity=0 (if $n - 2k > \gamma$), and
dissimilarity (if $\Delta_{n(k)} \leq \gamma$) counter

Extended Sokal-Michener index

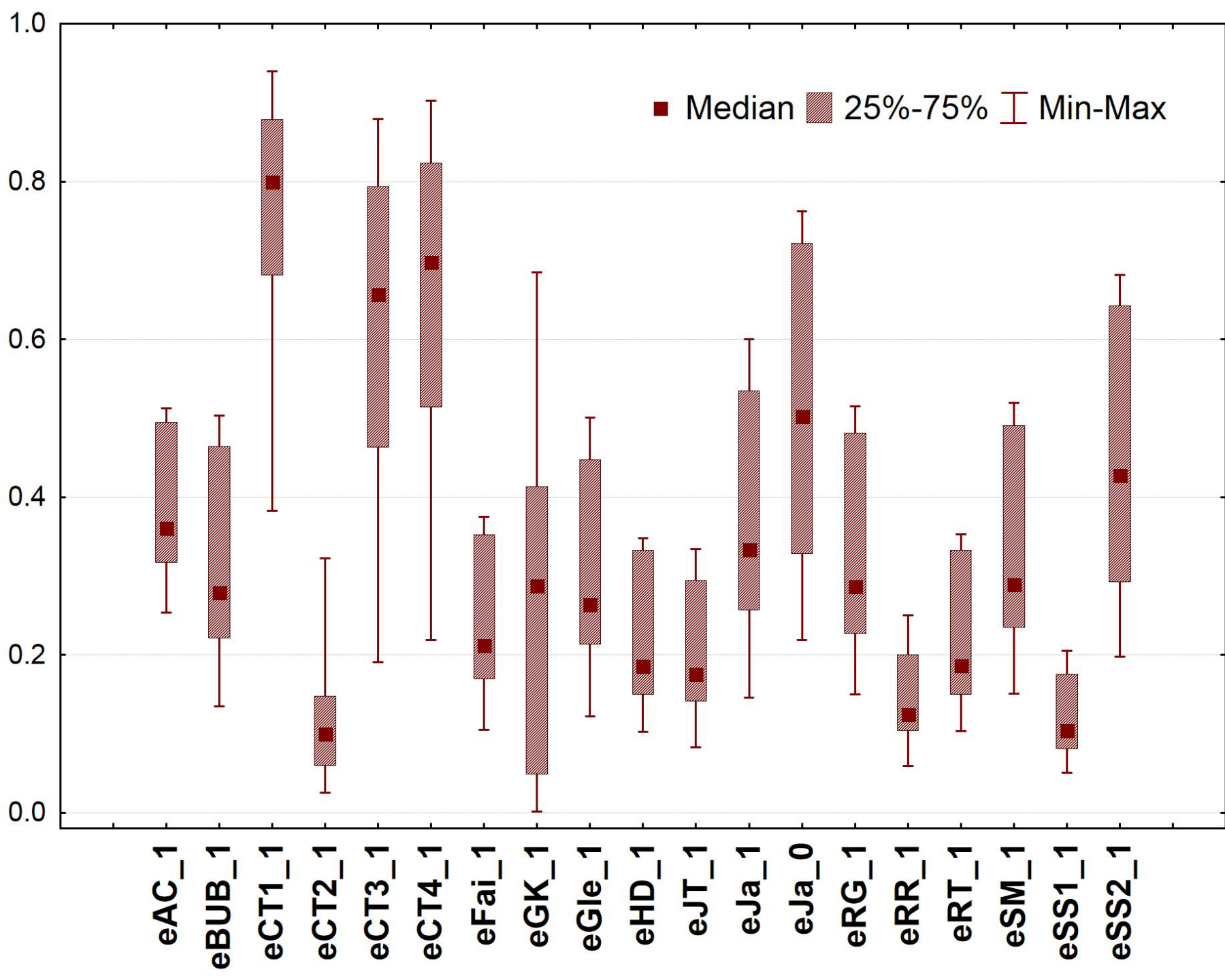
- $S_{eSM(1s_wd)} = \frac{\sum_s f_s(\Delta_{n(k)}) C_{n(k)}}{\sum_s f_s(\Delta_{n(k)}) C_{n(k)} + \sum_d f_d(\Delta_{n(k)}) C_{n(k)}}$
- $S_{eSM(1s_d)} = \frac{\sum_s f_s(\Delta_{n(k)}) C_{n(k)}}{\sum_s C_{n(k)} + \sum_d C_{n(k)}}$
- $f_s(\Delta_{n(k)}) = \frac{\Delta_{n(k)}}{n}$ és $f_d(\Delta_{n(k)}) = 1 - \frac{\Delta_{n(k)} - n \bmod 2}{n}$
- $\gamma = n \bmod 2$

Factorial ANOVA

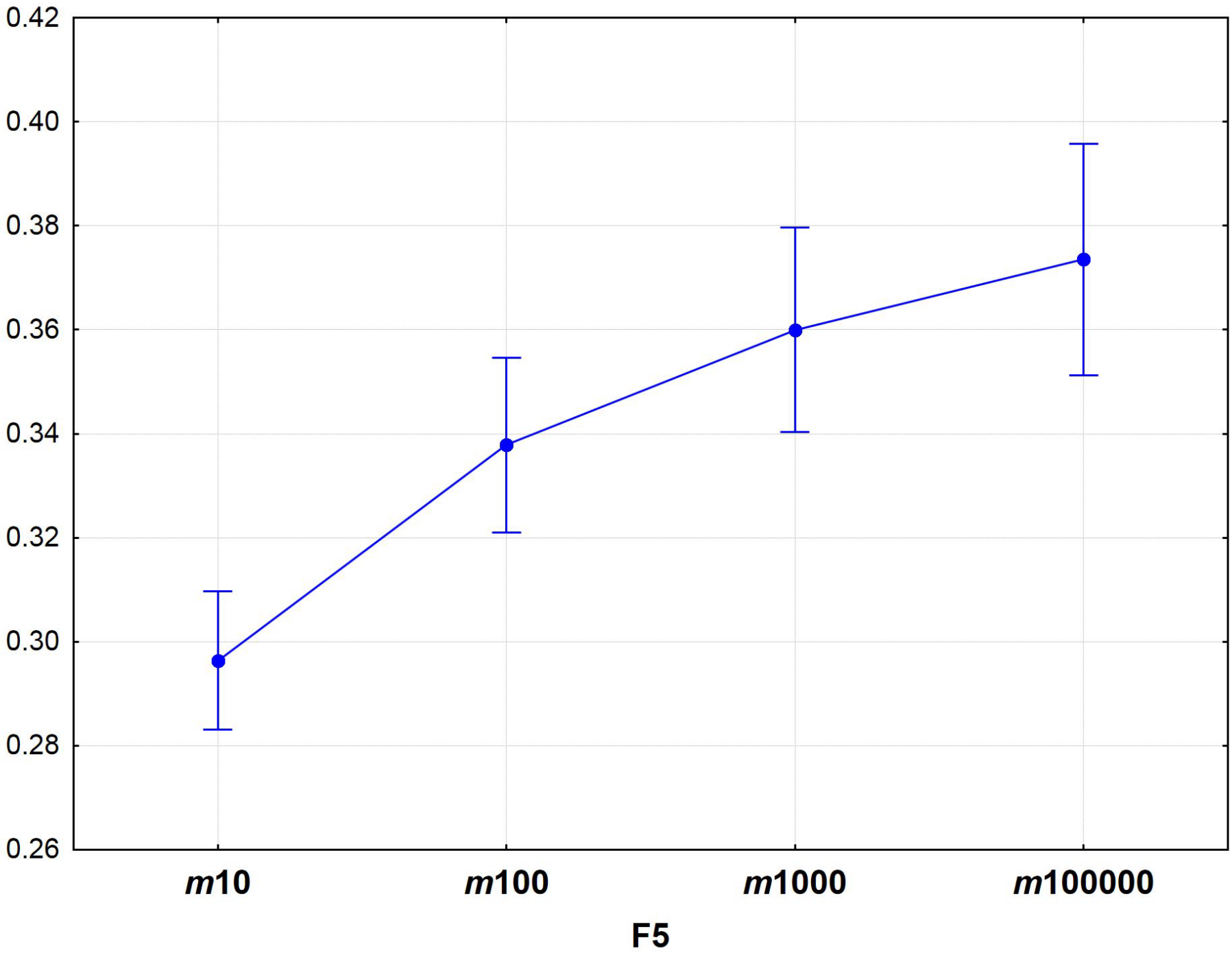
Factors taken into account:

- **F2**: Number of comparisons. Number (n) of bitstreams (e.g. fingerprints) compared, **14 levels**:
 $n = 2, 3, \dots, 15$ n -ary;
- **F3**: role of **weighting**, **two levels**: weighted and unweighted versions of the new similarity coefficients and
- **F4**: the **similarity coefficients** themselves, **19 levels**
- **F5**: m - length of fingerprints, **four levels**:
 $m = 10, 100, 1000, 100\ 000$ (**fingerprints** are random dichotomous vectors of **length m**);

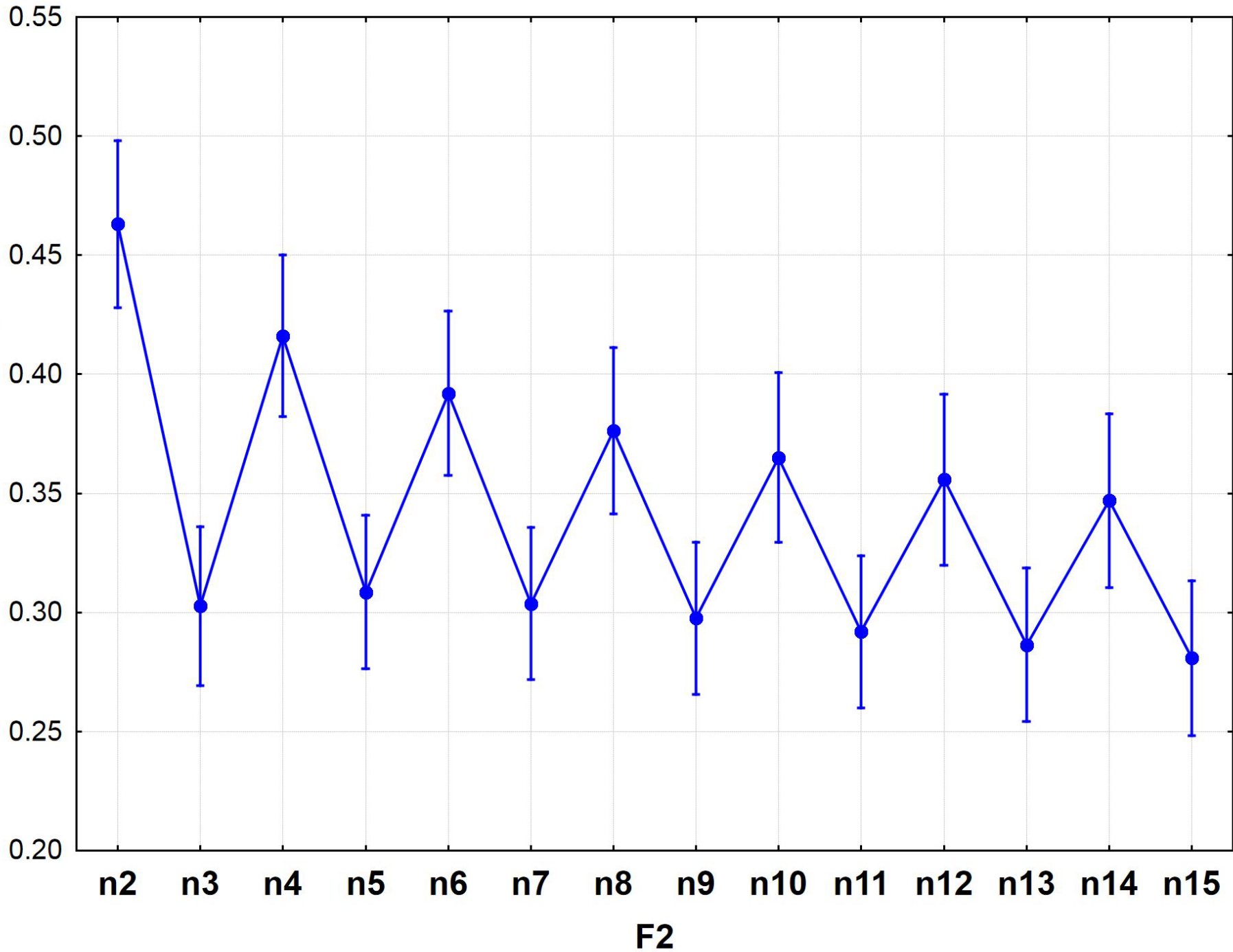
Mean of extended similarity coefficients

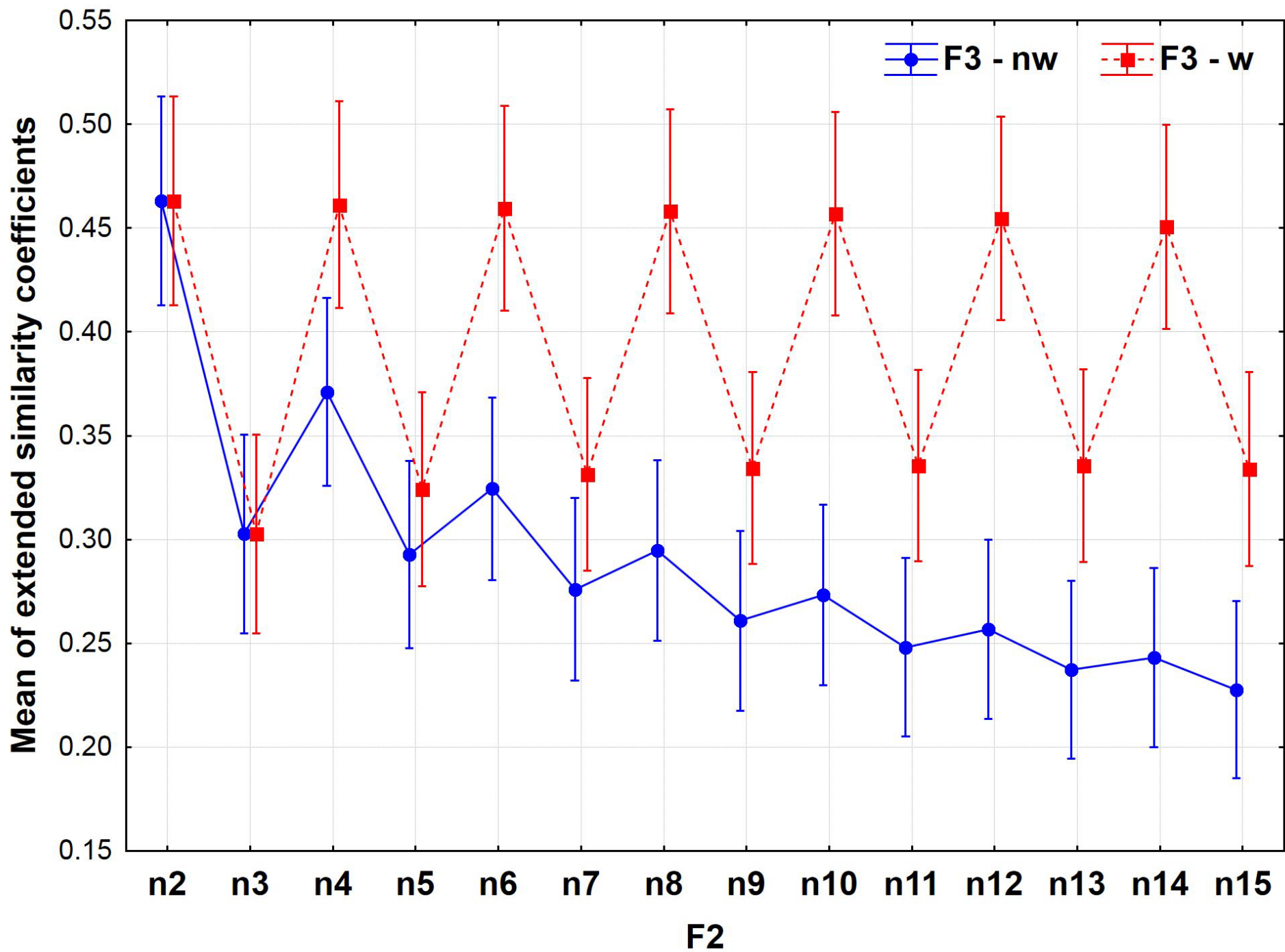


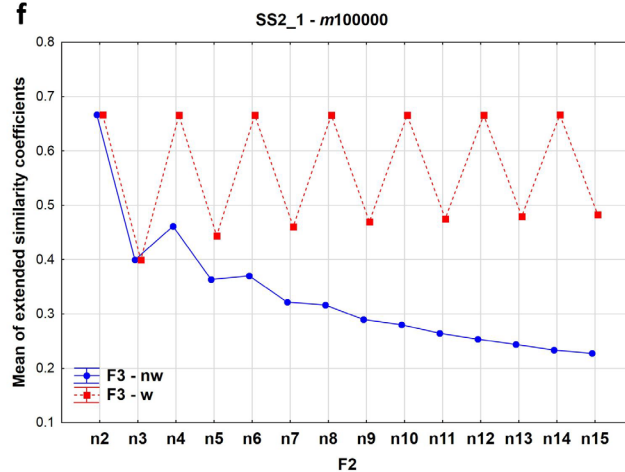
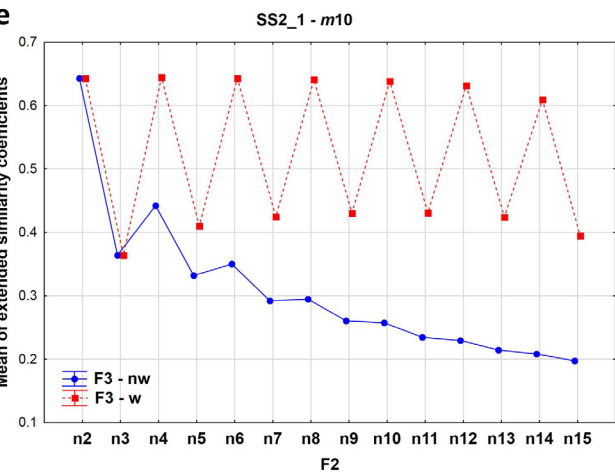
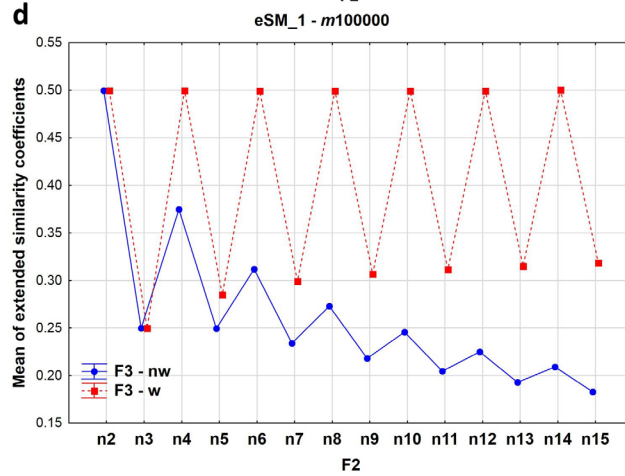
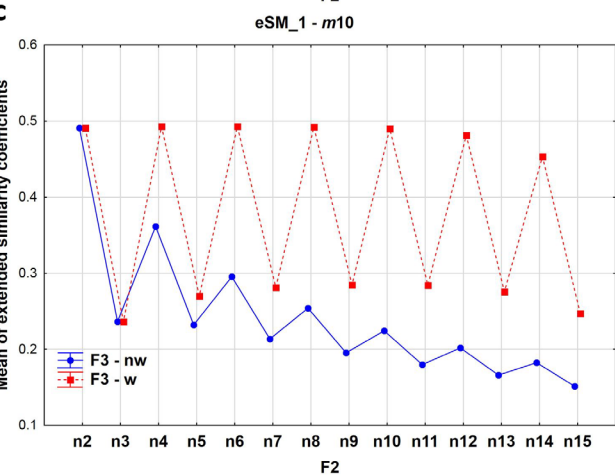
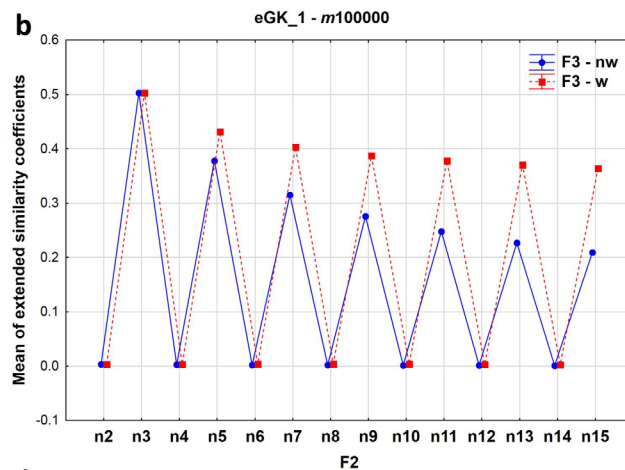
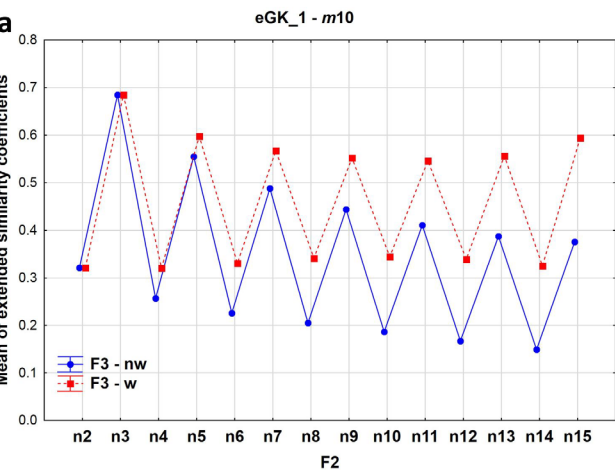
Mean of extended similarity coefficients



Mean of extended similarity coefficients







Lefthand side $m=10$

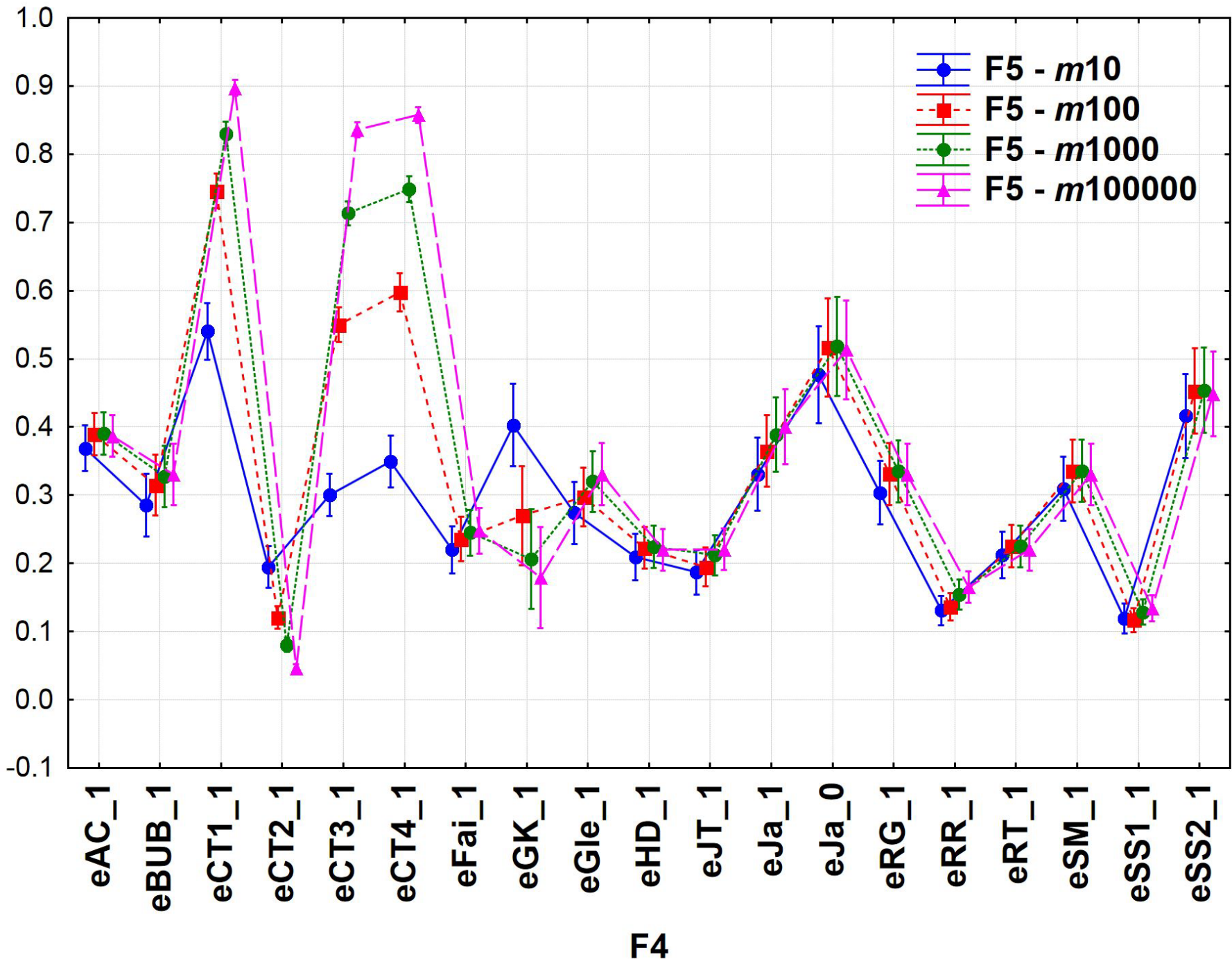
Righthand side $m=100\ 000$

eGK_1:
extended Goodman-Kruskal,

eSM_1:
extended Sokal-Michener

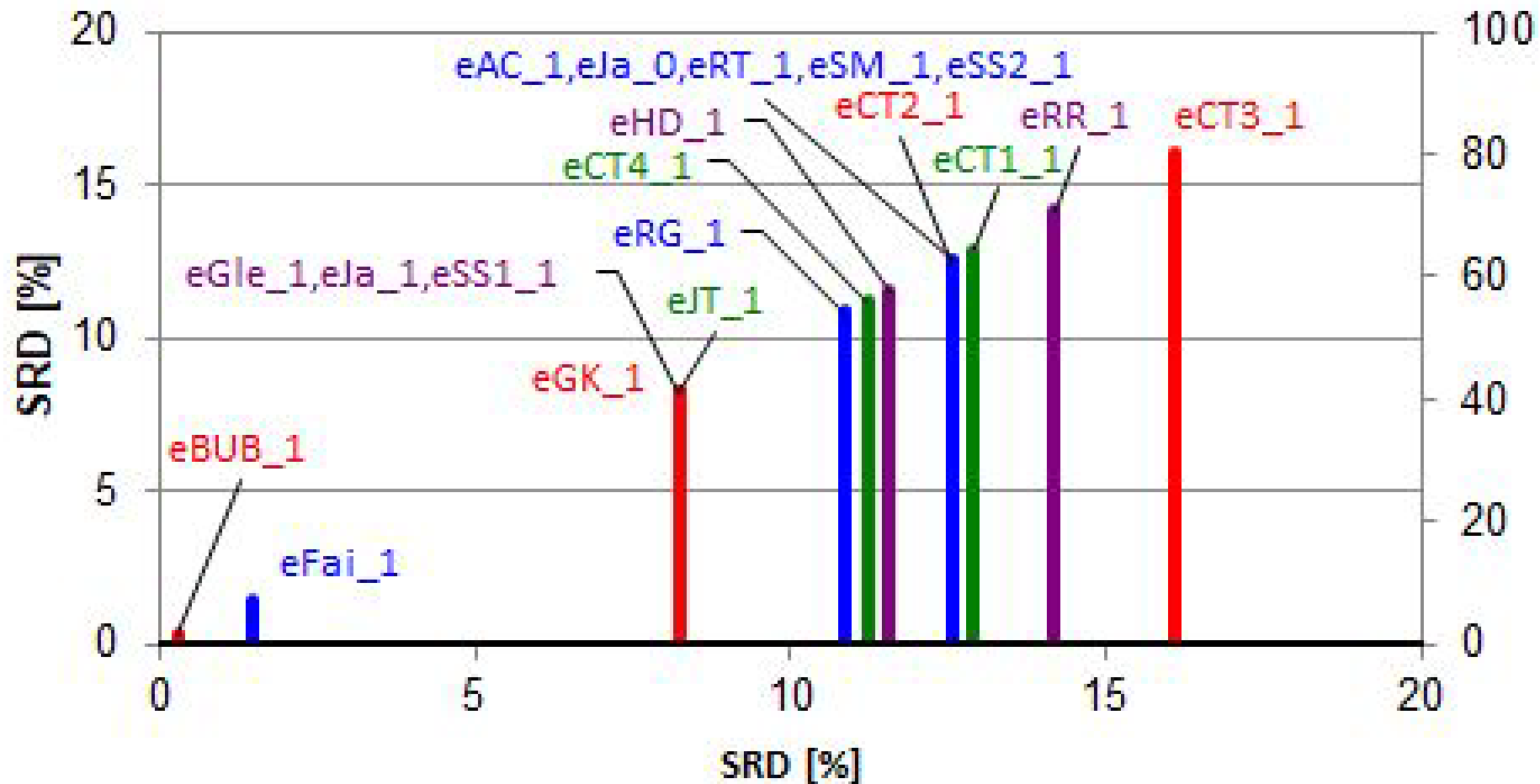
eSS_2:
extended Sokal-Sneath 2

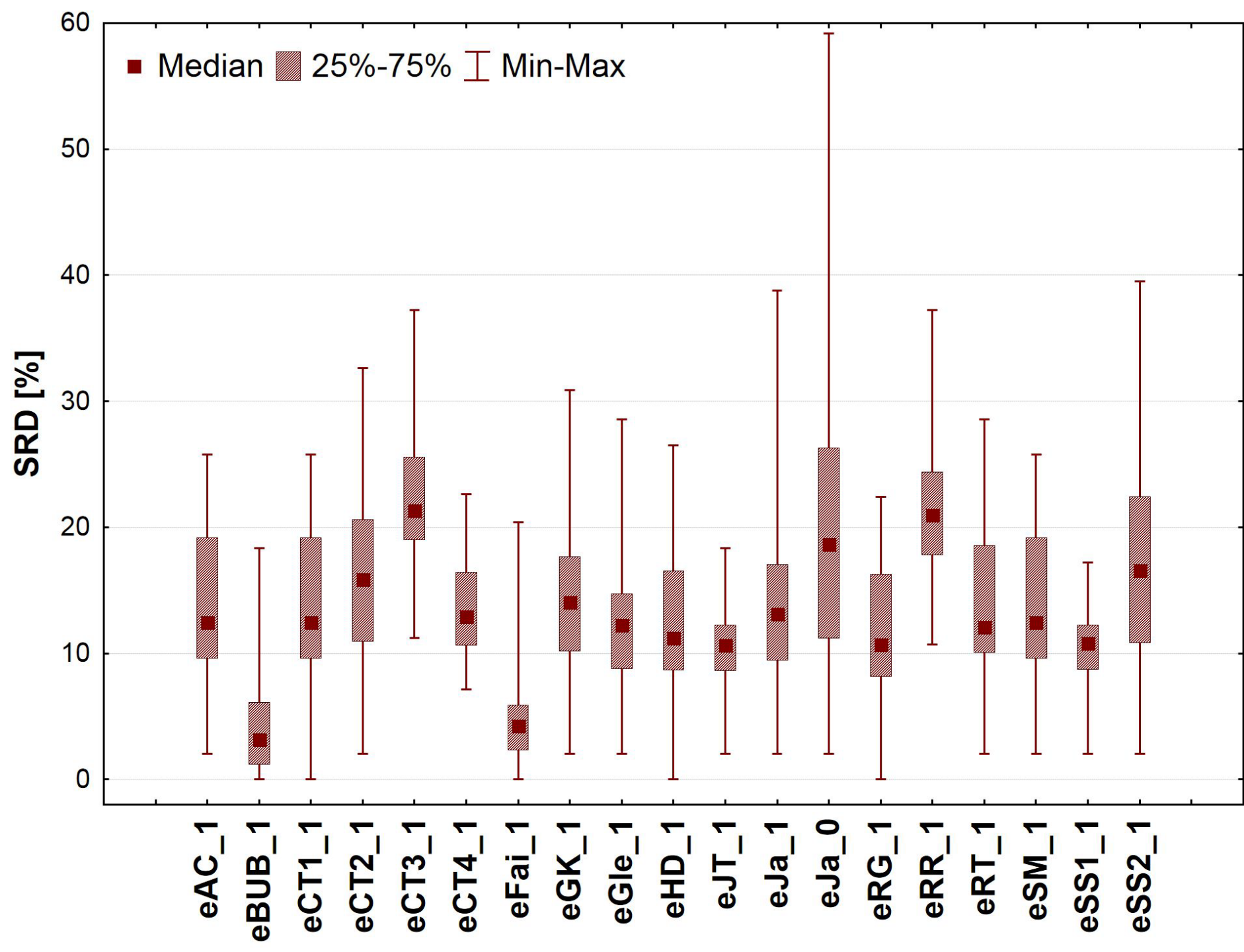
Mean of extended similarity coefficients

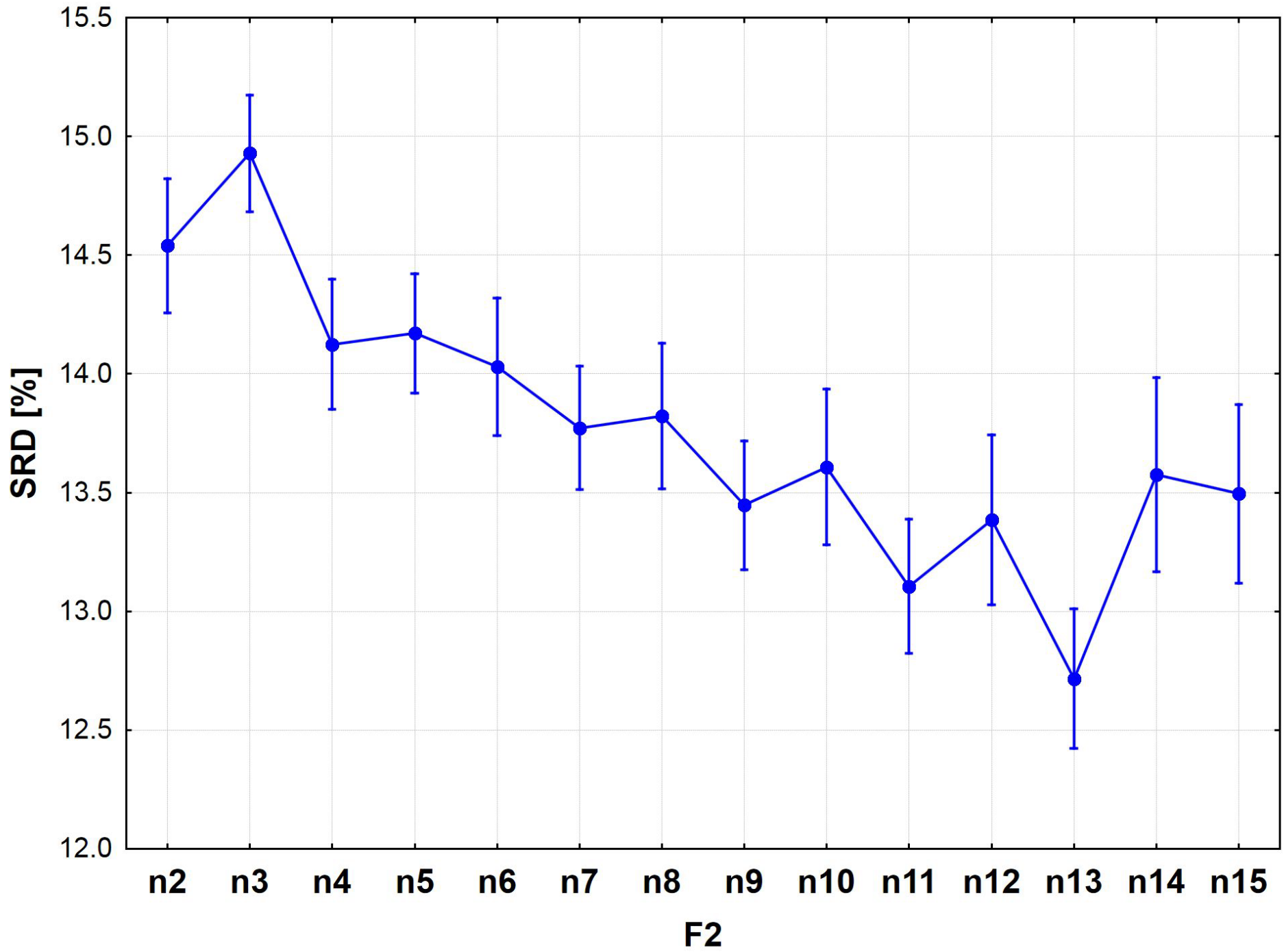


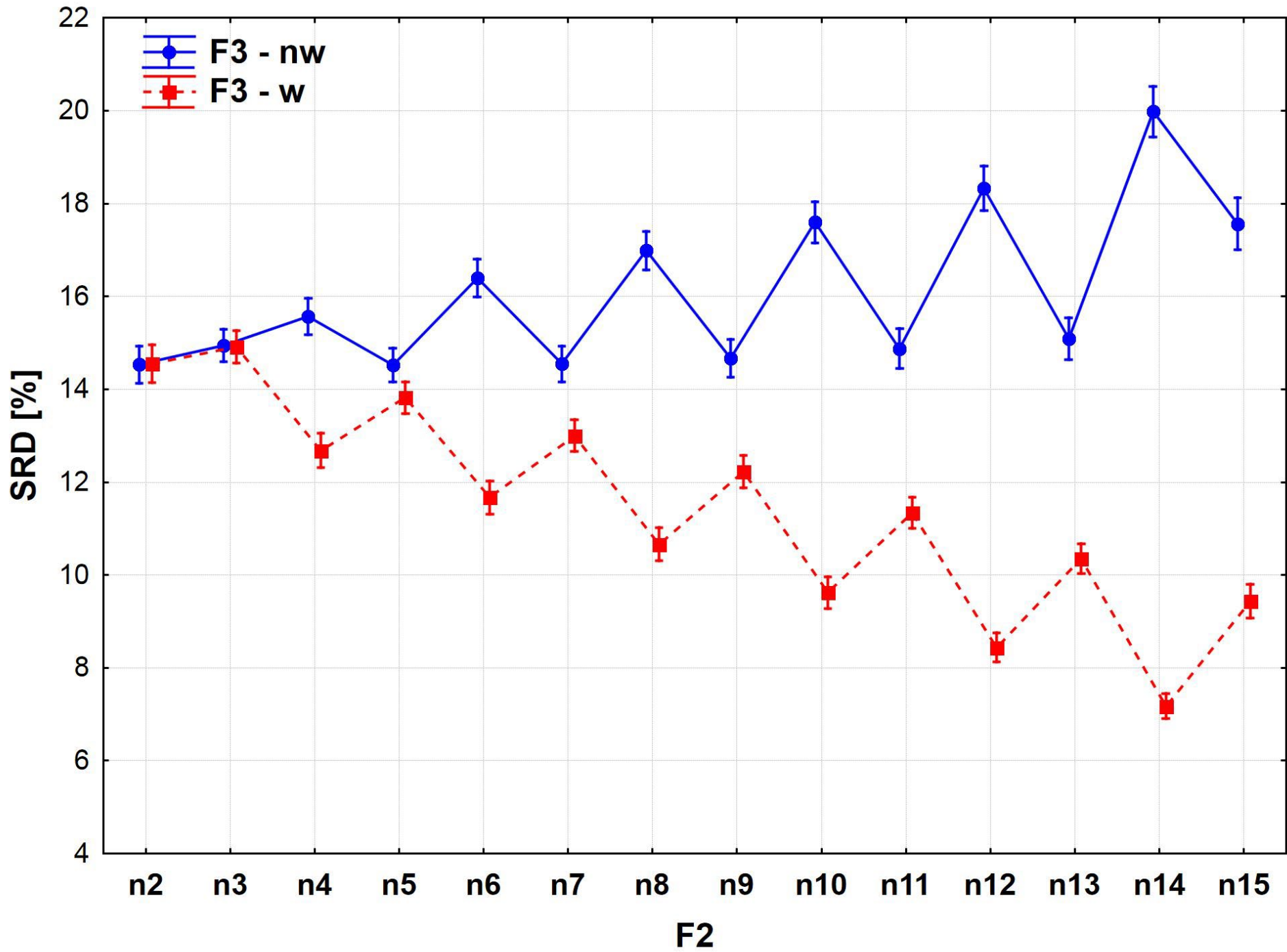
Ranking of extended indices

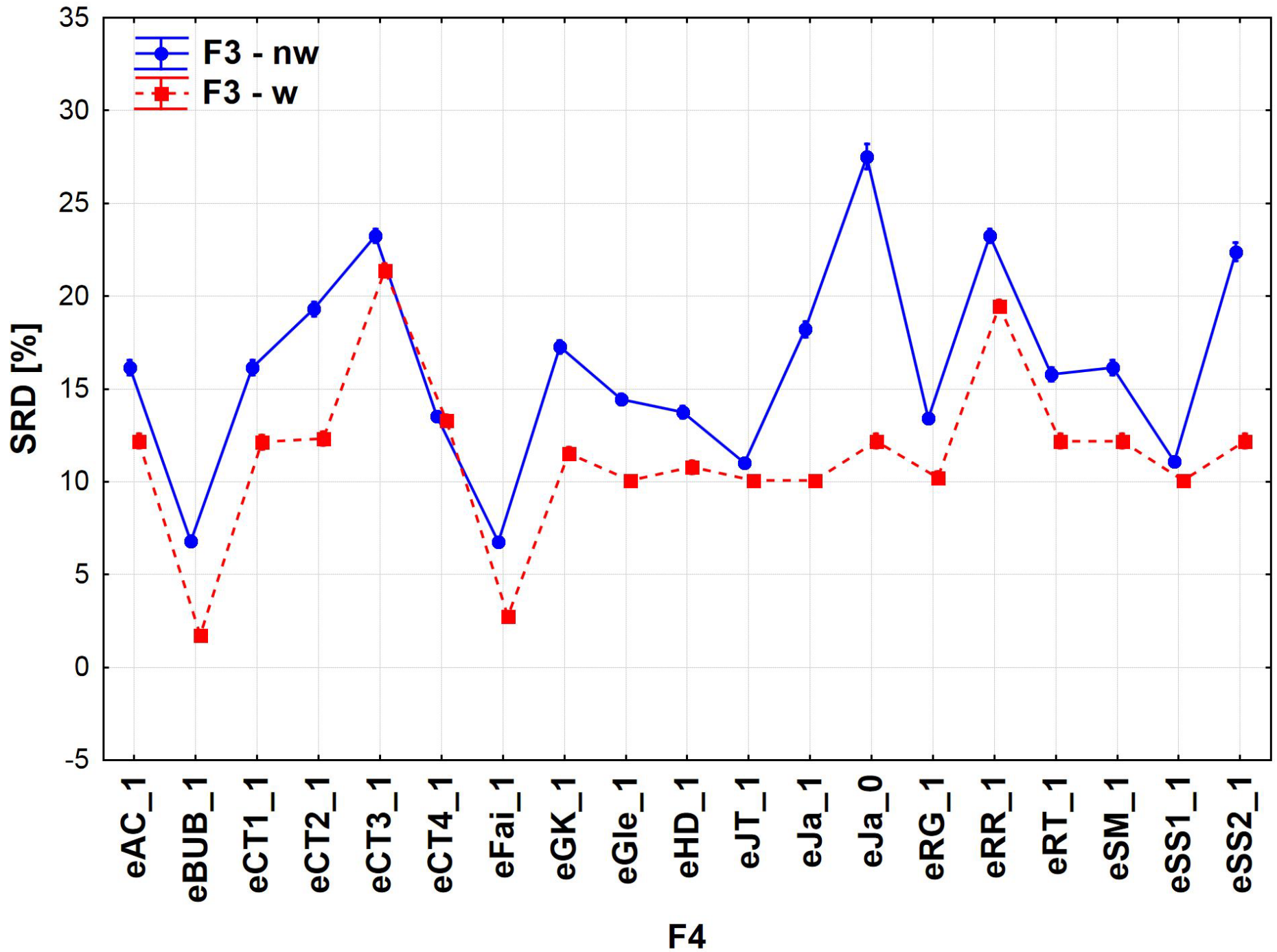
CRRN (NormApp, n=560 ; SRD%rep=0 ; d=0 ; ClassTH=1)

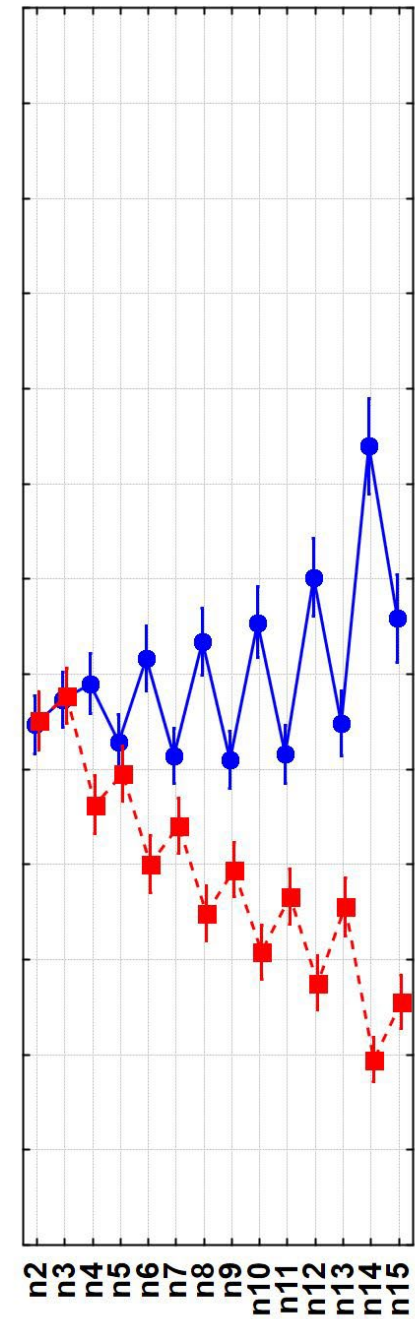
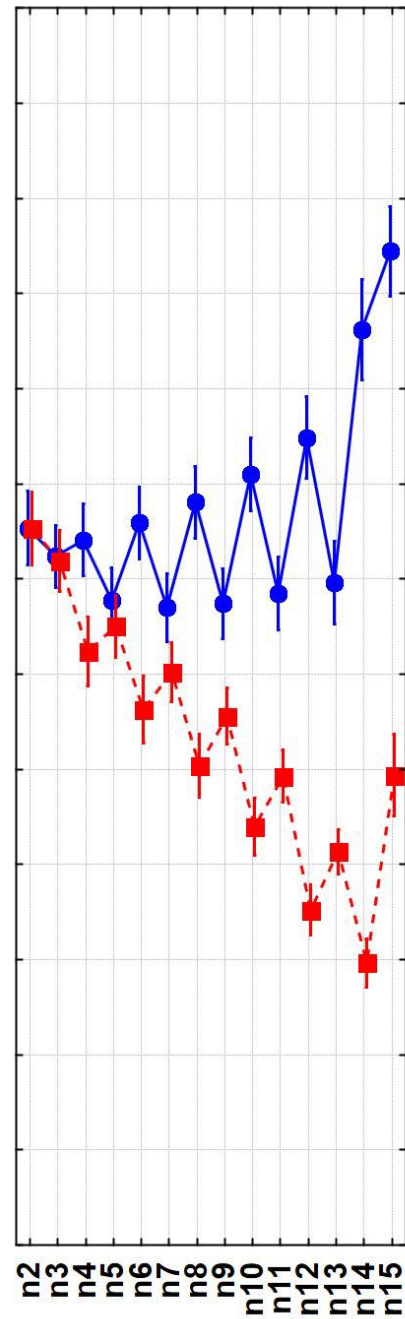
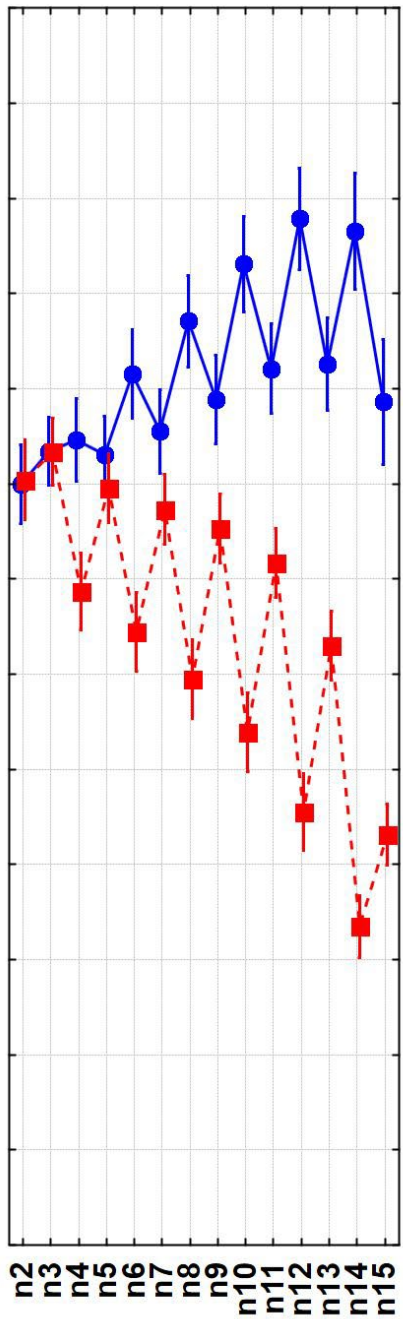
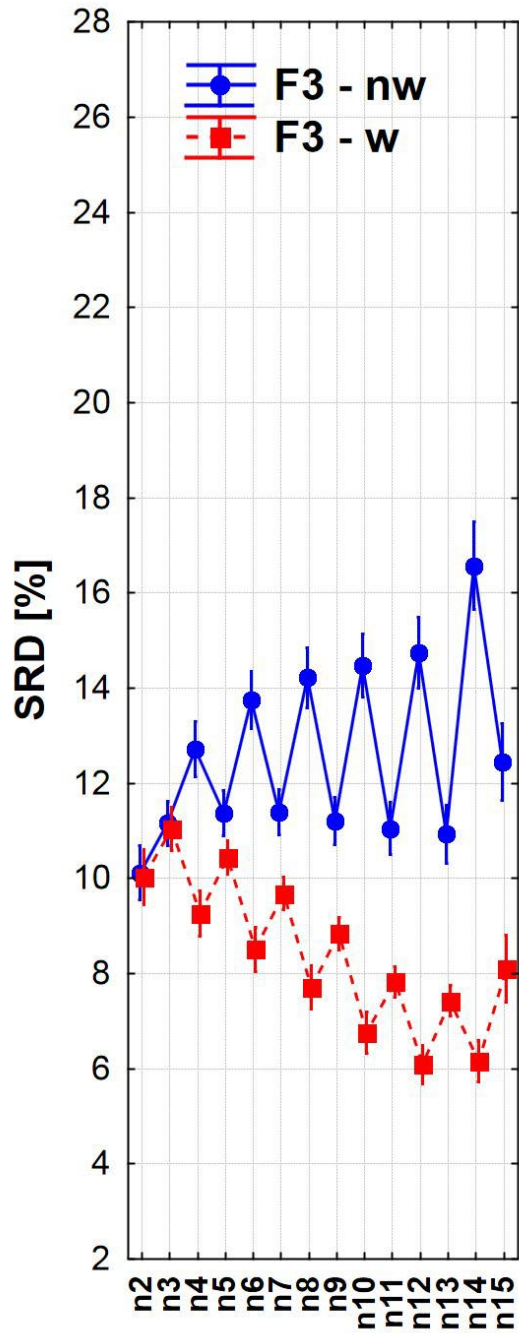










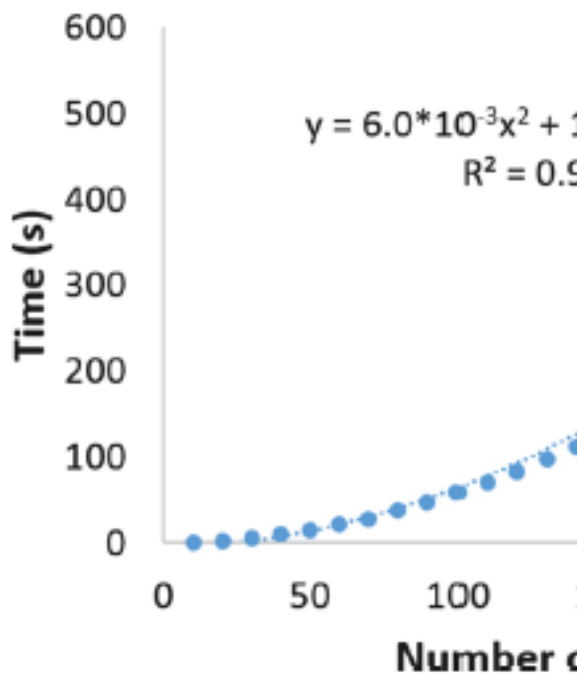


Summary

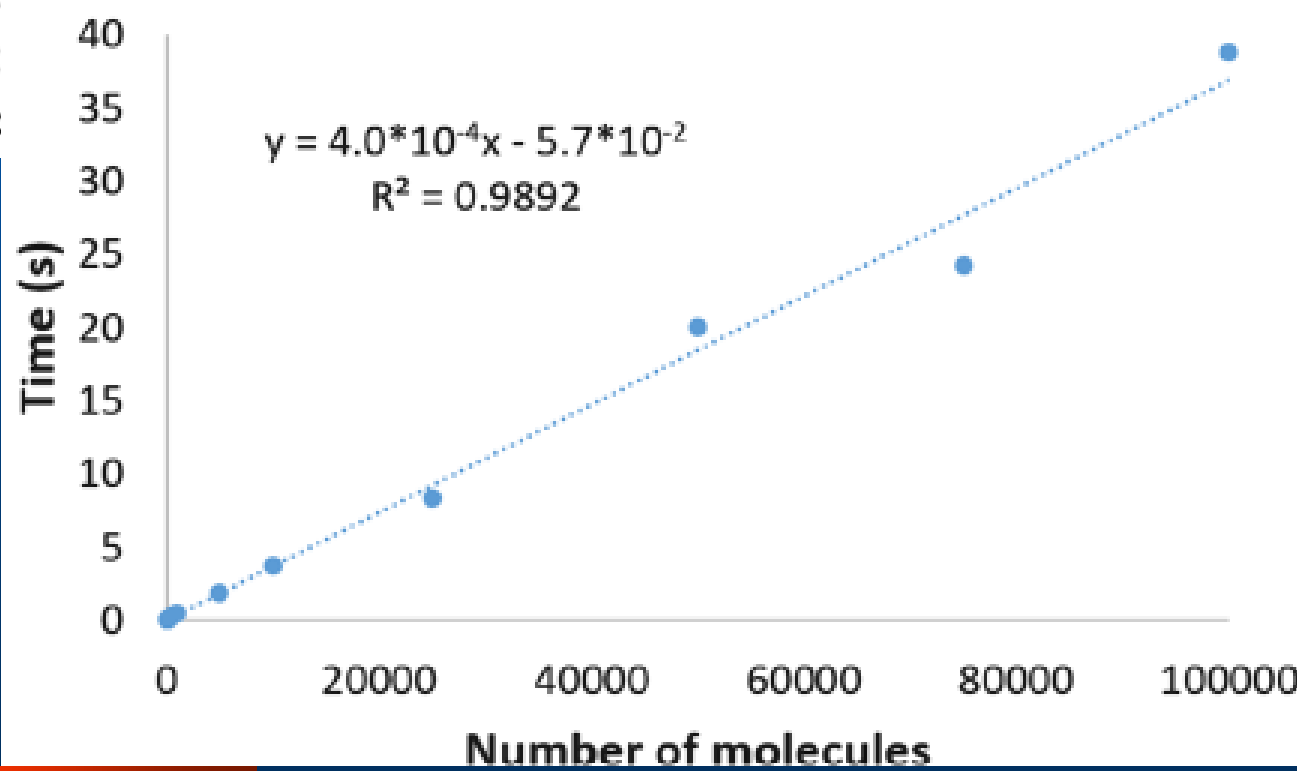
- By analogy (?) with the previous similarity indices, we have created **19 new**, symmetric **indices** allowing **multiple comparisons**.
- Their properties were extracted:
- Using ANOVA, we decomposed the effect of the factors: n , m , w , S , $n*w$, $w*S$, ...
- Using **SRD**, we found **optimal** factors, consistent solutions.

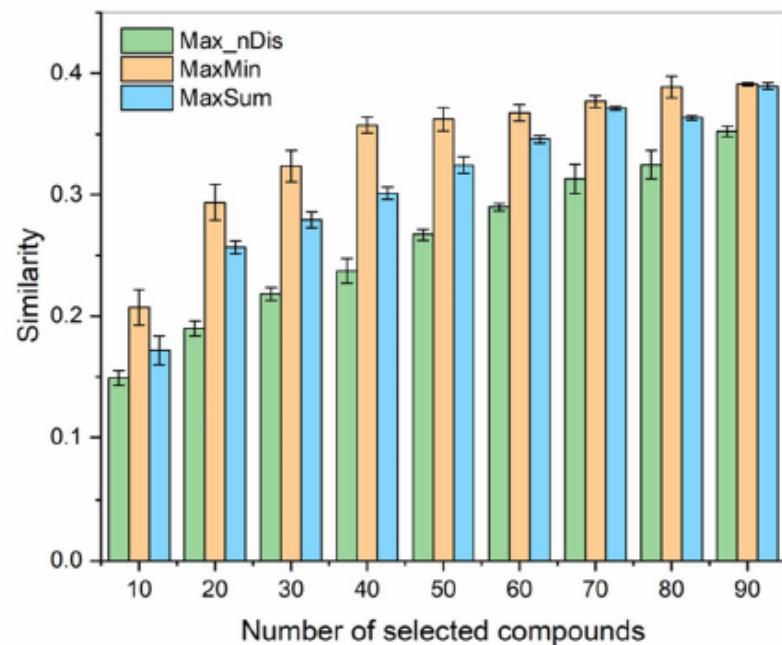
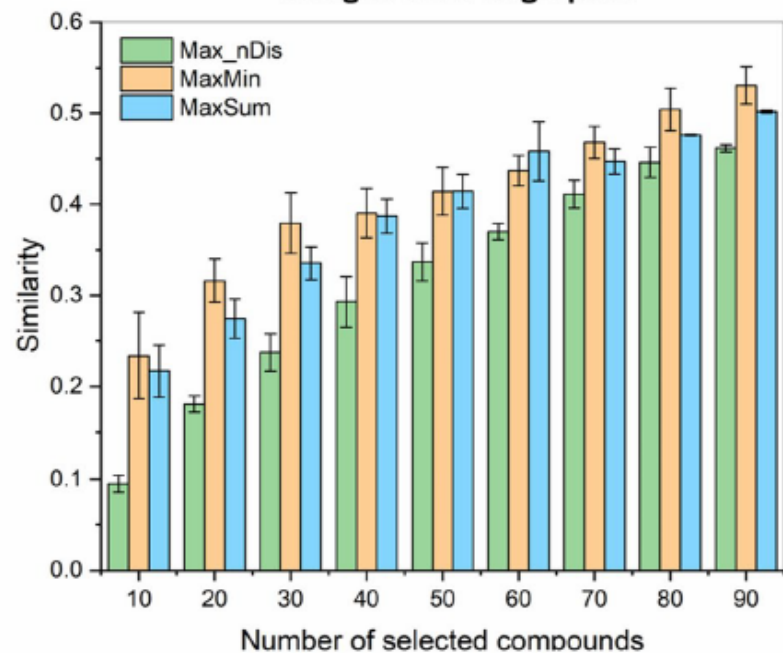
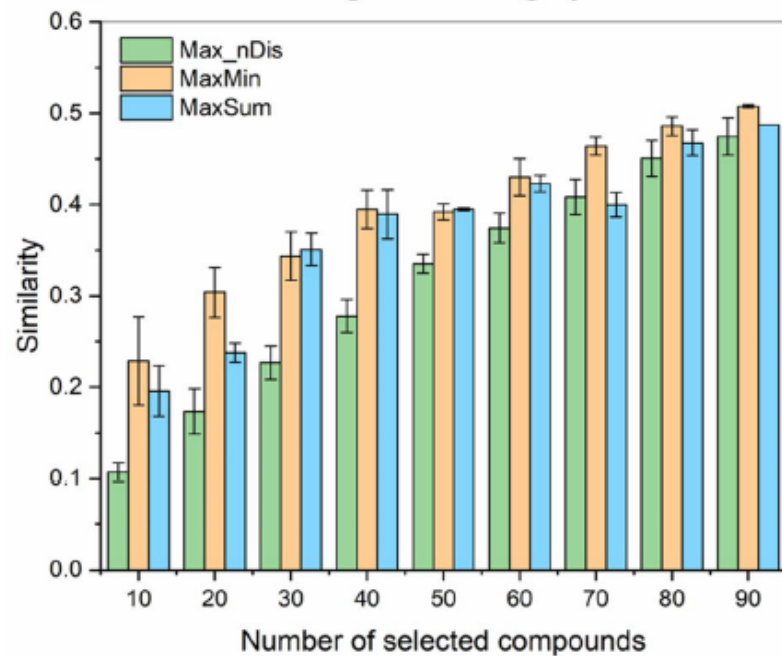
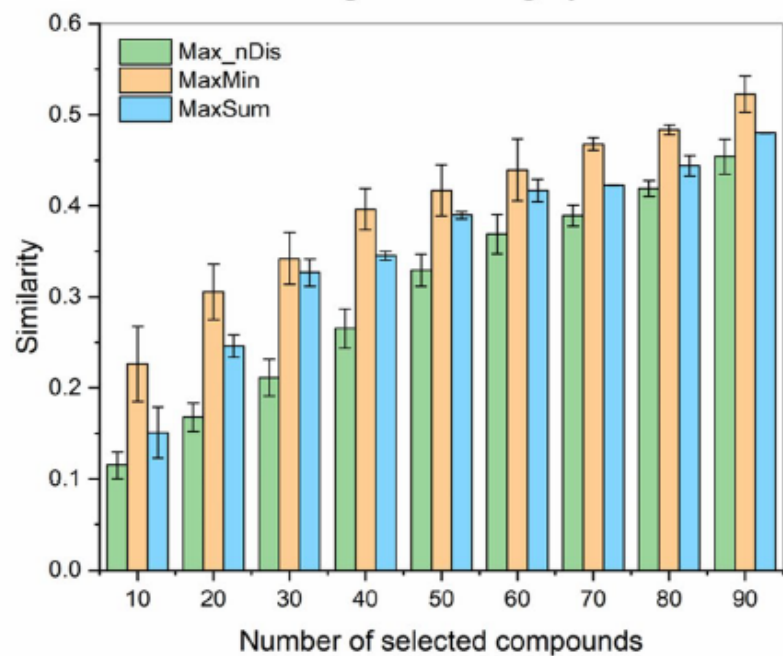
Conclusions

- The **optimal comparison number**? ($n=13$)
- Do you need weighting? **Yes** (**optimum: $n=14$**)
- **Extensions** of the **Baroni-Urbani-Buser** (eBUB) and **Faith** (eFai) coefficients (counting similarities of 1) are most recommended.
- The **fingerprint length** dependence of some indices (**eCTi** and **eGK**) is significant.
- Optimal fingerprint length? (10 or 100000, but an **optimal index** can be searched for a realistic 1000)

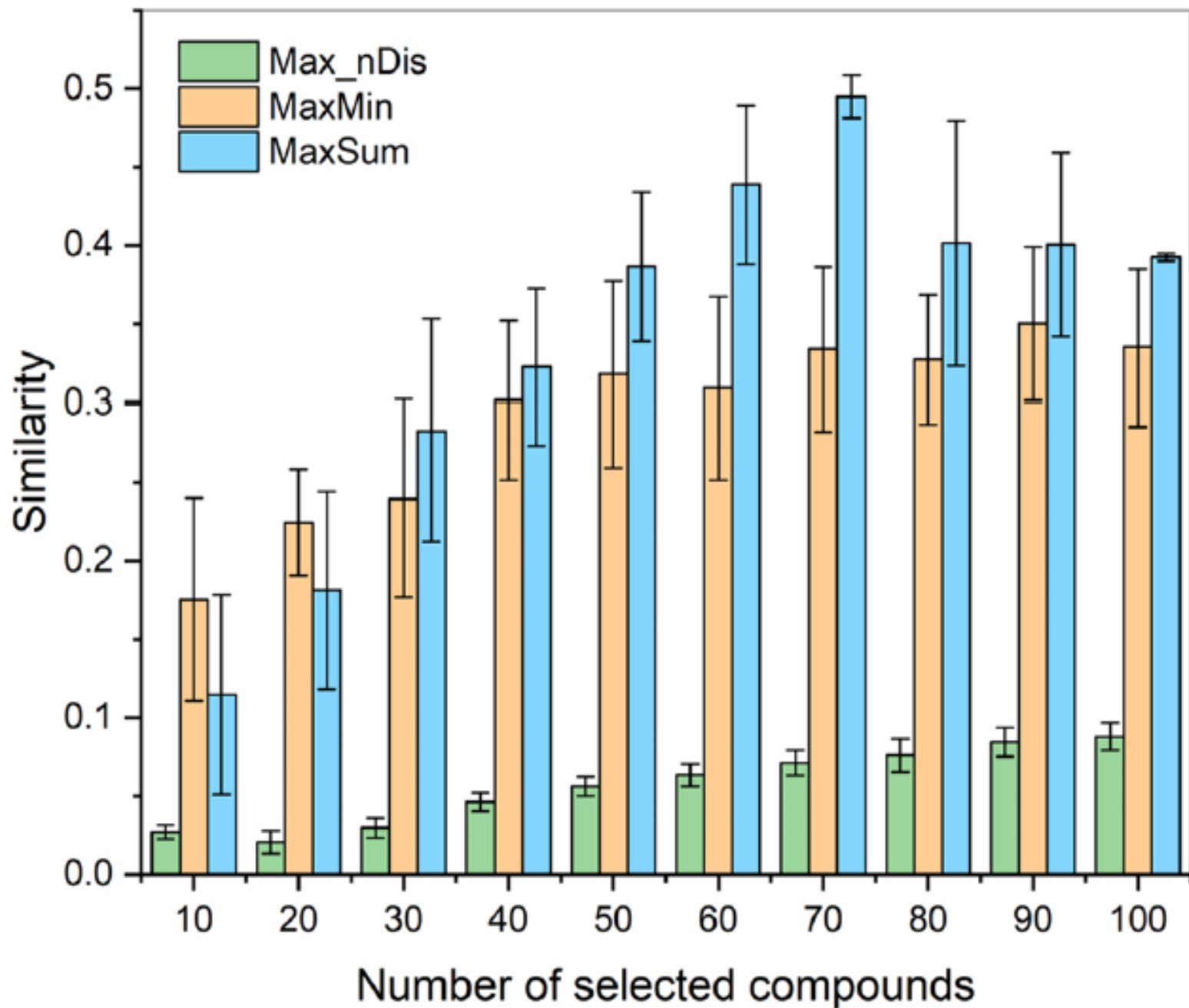
a**MACCS**

ion time

b**MACCS**

a MACCS fingerprint**b** Morgan-1024 fingerprint**c** Morgan-2048 fingerprint**d** Morgan-4096 fingerprint

CYP 2C9 dataset



Advantages of *n*-ary

- **Less computing power**, faster algorithm, larger datasets,
- **Diversity detection** significantly better,
- Possibility of **clustering** (HCA),
- Definition of **internal and external consistency indices** & their optimization is possible.

Acknowledgement

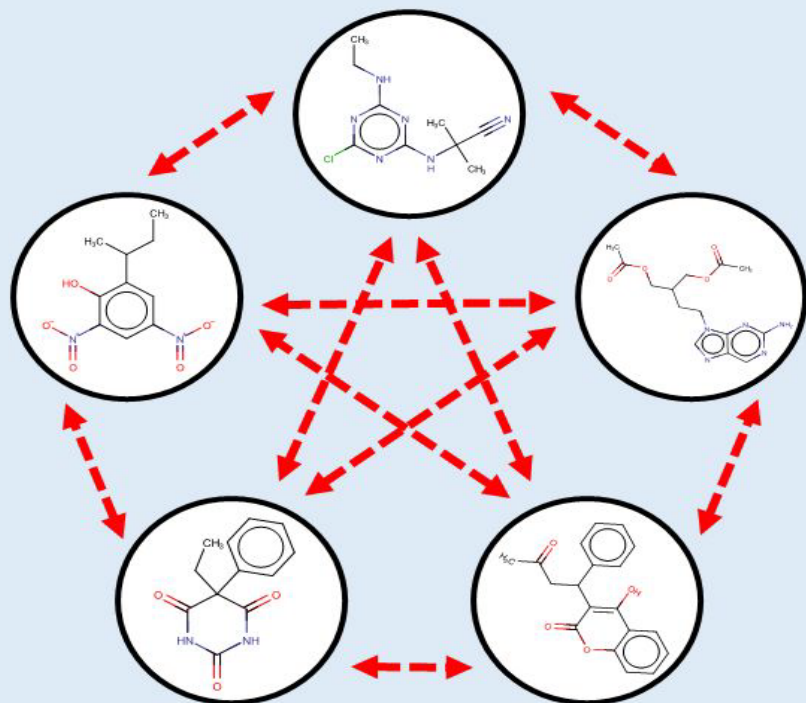
- National Research Development and Innovation Fund (Hungary): OTKA **K134260**
- University of Florida startup grant (RAMQ)
- **Covid pandemic**

Releted papers

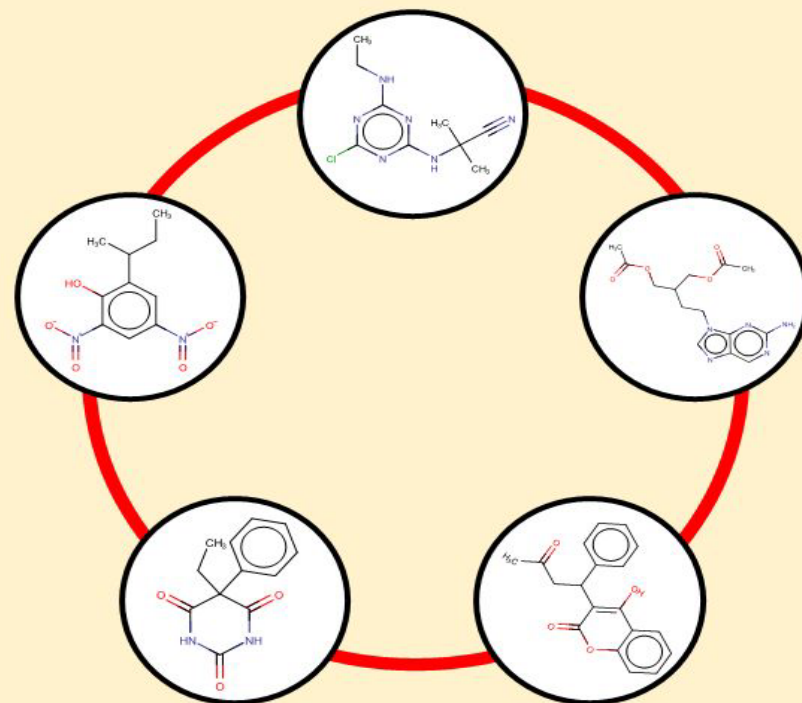
- Ramón Alain Miranda-Quintana*, Dávid Bajusz, Anita Rácz, Károly Héberger*, Extended similarity indices: the benefits of comparing more than two objects simultaneously. Part 1: Theory and characteristics *Journal of Cheminformatics*, 13, Article No: 32 (2021) if(2020)=5.514 (D1)
- Ramón Alain Miranda-Quintana*, Dávid Bajusz, Anita Rácz, Károly Héberger*, Extended similarity indices: the benefits of comparing more than two objects simultaneously. Part 2: speed, consistency, diversity selection *Journal of Cheminformatics* 13, Article No: 33 (2021) if(2020)=5.514 (D1)
- Ramon Miranda-Quintana*, Dávid Bajusz, Anita Rácz, Károly Héberger*, Differential Consistency Analysis: Which similarity measures can be applied in drug discovery? *Molecular Informatics*, 40, Article No.: 2060017 (2021) if(2020)=3.353 (Q2)
- Dávid Bajusz, Ramón Alain Miranda-Quintana*, Anita Rácz, Károly Héberger*, Extended many-item similarity indices for sets of nucleotide and protein sequences *Computational and Structural Biotechnology Journal* 19, 3628-3639 (2021) if(2020)=7.271 (Q1)

← COMPUTE TIME

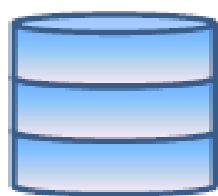
Binary similarity



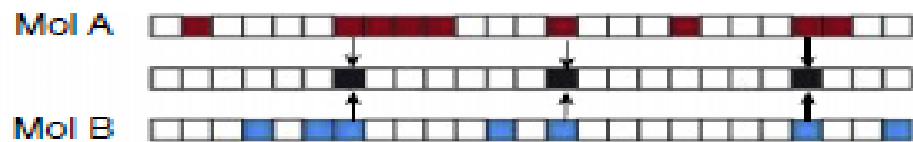
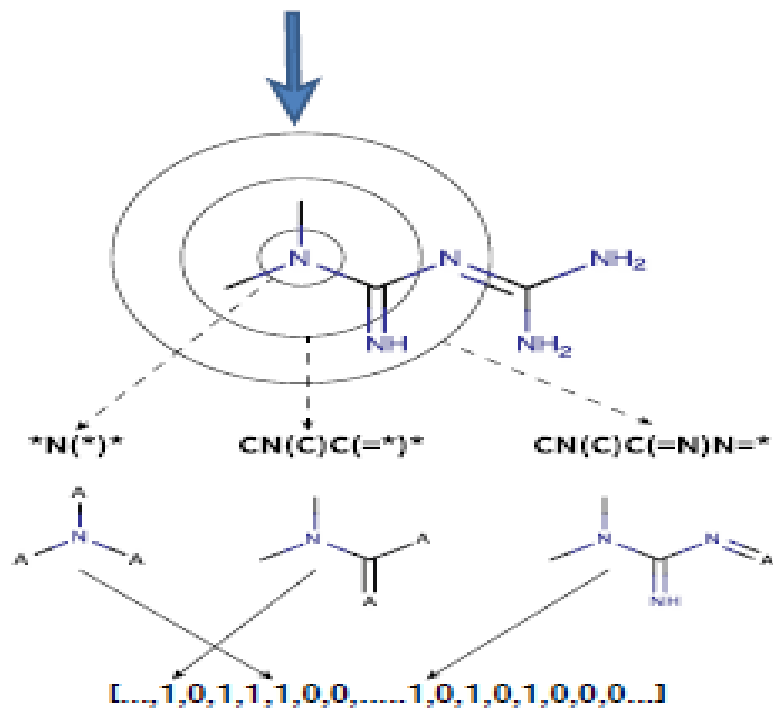
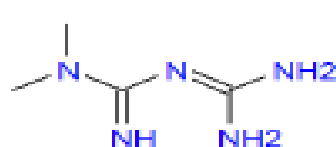
Extended (n -ary) similarity



DIVERSITY →



Compound database

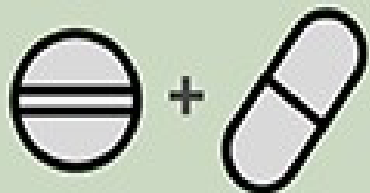


$$Tc(\text{Mol A } \blacksquare, \text{Mol B } \blacksquare) = \frac{\blacksquare}{\blacksquare + \blacksquare - \blacksquare} = \frac{3}{9 + 7 - 3} = 0.23$$

Yu-Chen Lo, Stefano E. Rensi, Wen Torng and Russ B. Altman, Machine learning in chemoinformatics and drug discovery *Drug Discovery Today* **23**, 1538-1546 (2018)

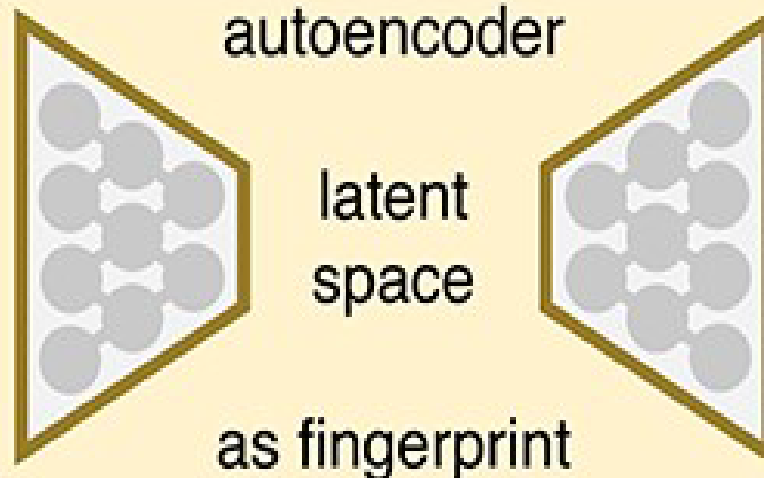
Input

drug pairs
tested
in cell lines



Data-driven

autoencoder



Rule-based

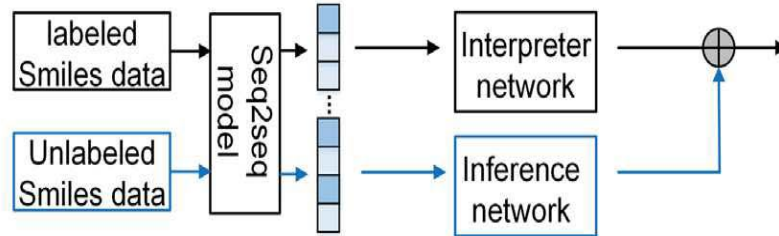
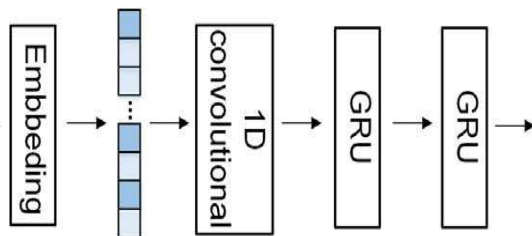
2D/3D
molecular
fingerprints

Compare via

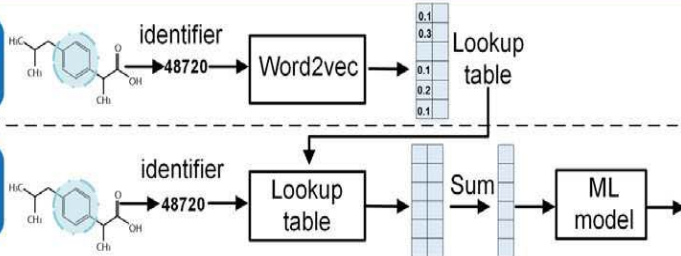
- synergy prediction
- clustering
- representation similarity

B. Zagidullin, Z. Wang, Y. Guan, E. Pitkänen and J. Tang: Comparative analysis of molecular fingerprints in prediction of drug combination effects *Briefings in Bioinformatics*, 00(00), 2021, 1–15 <https://doi.org/10.1093/bib/bbab291>

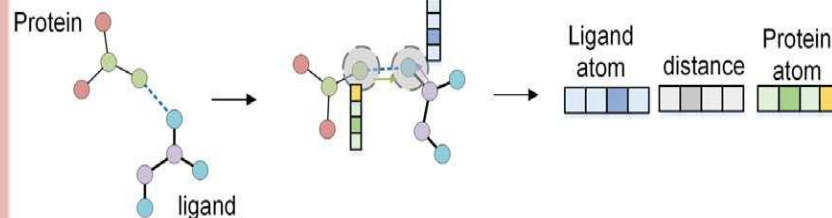
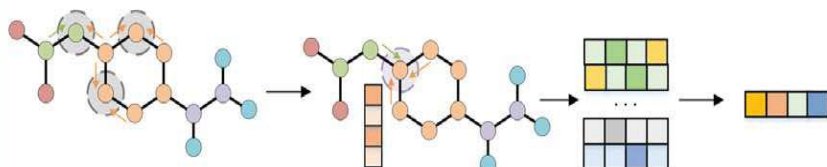
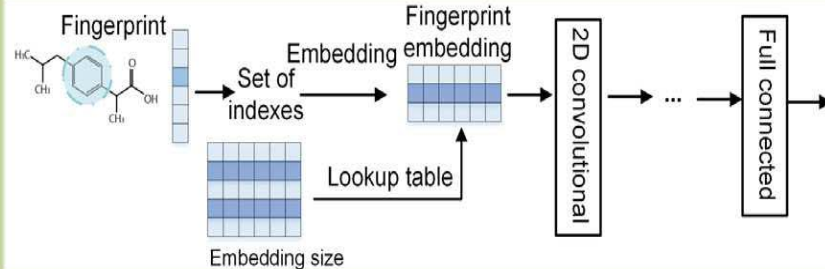
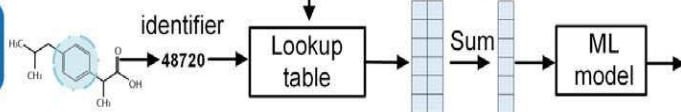
Smiles Data:
s1c2ncnc(O)c2cc1CC
COCCOCCO
C1=CC=C(C=C1)C=O



Unsupervised pre-training



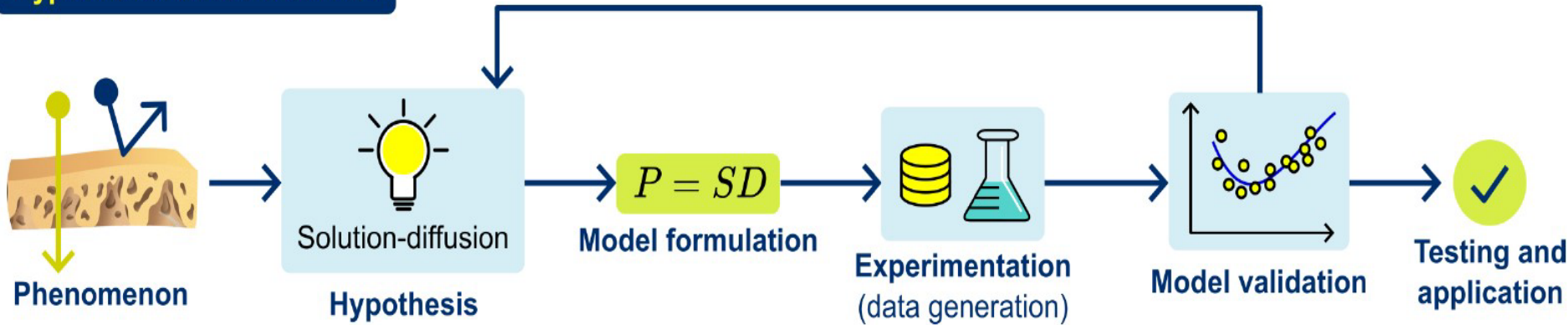
Supervised learning



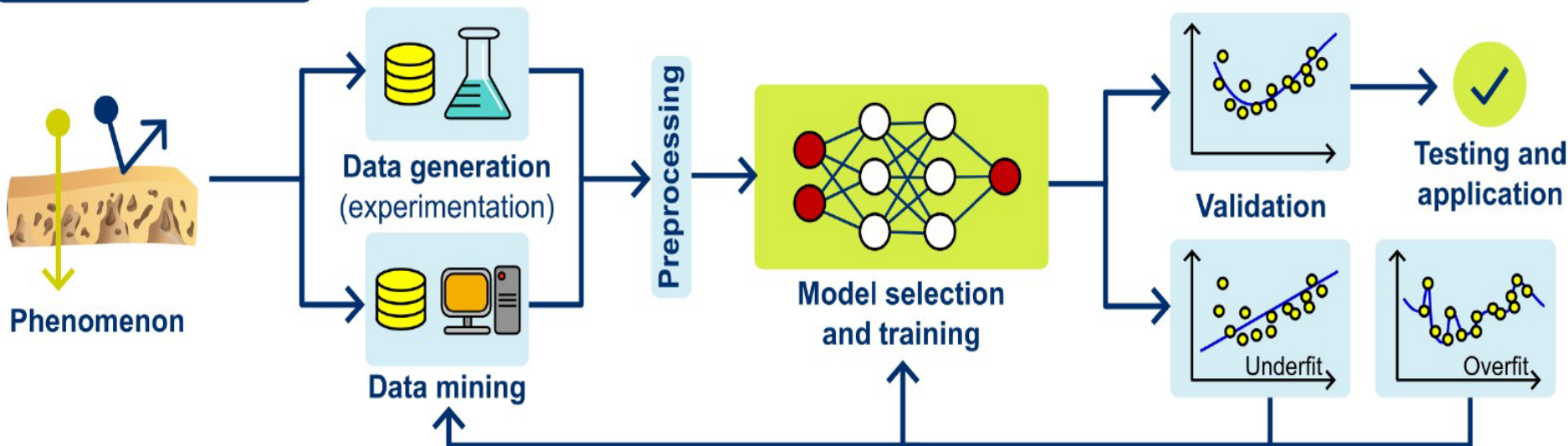
Youjun Xu, Chenjing Cai, Shiwei Wang, Luhua Lai, Jianfeng Pei: Efficient molecular encoders for virtual screening. *Drug Discov Today Technol.* 2019 Dec; **32-33:19-27. doi: 10.1016/j.ddtec.2020.08.004 Epub 2020 Oct 4.**

Machine learning vs. classic

Hypothesis-driven research

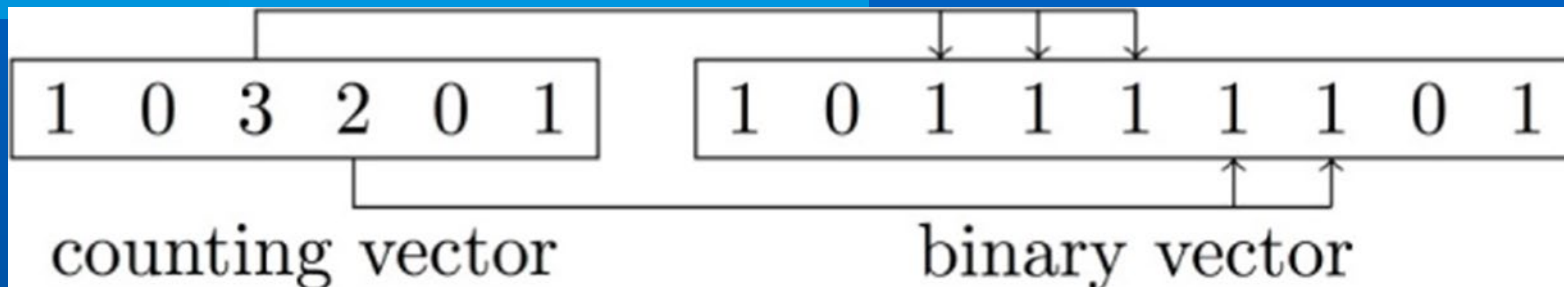


Data-driven research

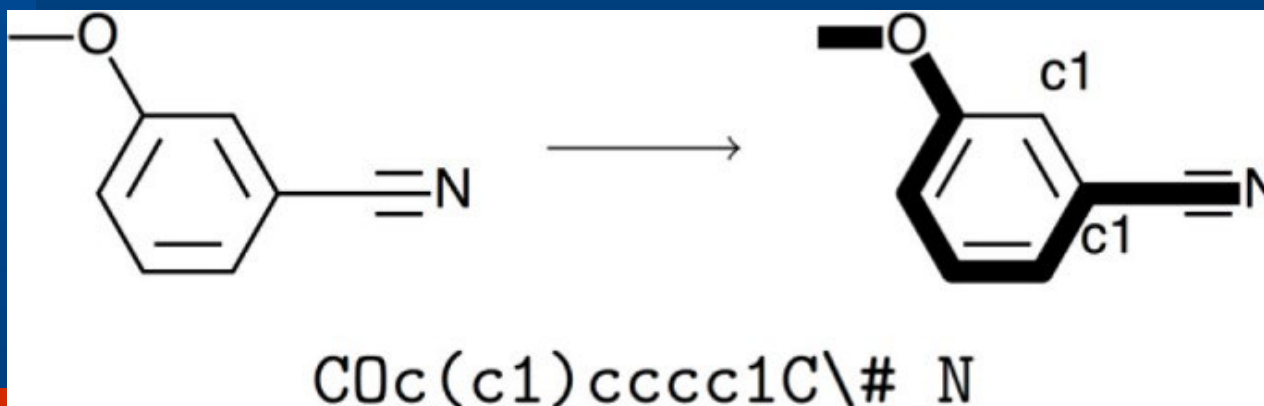


Molecular representations

- A **binary fingerprint** describes the molecule as a set of features:



- A **counting vector** allows for a more detailed description of the molecule as a multi-set of features
- Linear representation: SMILES:



Main steps of drug design

