

Dwa nowe algorytmy klasyfikacji farmakokinetycznych szeregów czasowych

Two novel classification algorithms of the pharmacokinetic time series

Piotr Wilczek¹

Streszczenie: Przedkładana praca zawiera propozycję dwóch nowych miar odległości pomiędzy szeregami czasowymi. Mianowicie, wprowadzamy tzw. rekurencyjno-kanoniczną miarę odległości oraz jej krzyżową formę. Nasze podejście opieramy na bezprogowych (krzyżowych) macierzach rekurencyjnych, których koncepcje zostały rozwinięte w ramach (krzyżowej) analizy rekurencyjnej szeregów czasowych oraz na kanonicznej mierze odległości pomiędzy dwoma wielowymiarowymi zbiorami danych, której koncepcja została rozwinięta w obrębie chemometrii. Wyczerpujące symulacje komputerowe przeprowadzone na 27 farmakokinetycznych szeregach czasowych pokazały, iż algorytmy bazujące na nowo opracowanych funkcjach odległości są bardziej efektywne niż algorytmy bazujące na klasycznych miarach $L2$ i DTW oraz na ich modyfikacjach.

Abstract: In the present contribution, we proposed two novel distance measures between time series. Namely, we introduced the so-called recurrence-canonical measure of distance as well as its cross form. Our approach is based on the notion of the unthresholded (cross-)recurrence matrices developed in the field of the (Cross-)Recurrence Quantification Analysis and the notion of the canonical measure of distance between multidimensional datasets developed in the field of chemometrics. The extensive computer simulations carried out on 27 pharmacokinetic time series showcased that the algorithms based on the newly designed distance functions outperform the protocols based on the classical $L2$ and DTW functions and on their modifications.

Słowa kluczowe: szereg czasowy, (krzyżowa) analiza rekurencyjna, kanoniczna miara odległości, farmakokinetyka

Keywords: time series, (cross-)recurrence quantification analysis, canonical measure of distance, pharmacokinetics

¹ Computer Laboratory, 61-163 Poznań, e-mail: piotr.wilczek.net@onet.pl, ORCID 0000-0003-2758-343X.

1. Wprowadzenie

W ostatnich latach można zaobserwować wzrost popularności analiz opartych na tzw. „danych temporalnych”. Mianowicie, biorąc pod uwagę powszechne wykorzystanie nowoczesnych technologii cyfrowych, możliwa jest rejestracja dużej liczby pomiarów czasowych w obrębie nauk podstawowych (np. w fizyce, chemii, biologii czy też geologii), jak również nauk stosowanych (np. w farmakologii, medycynie, robotyce czy też automatyce). W związku z tym nastąpił dramatyczny wzrost zainteresowania gromadzeniem oraz eksploracją danych temporalnych, co z kolei zaowocowało dużą liczbą artykułów naukowych proponujących nowatorskie techniki indeksacji, klasyfikacji, grupowania, aproksymacji oraz prognozowania szeregów czasowych. W szczególności zaproponowano wiele nowych miar odległości pomiędzy danymi temporalnymi. Przedkładana praca ma na celu wprowadzenie dwóch nowych funkcji odległości pomiędzy szeregami czasowymi. Efektywność nowo zaproponowanych miar temporalnych zostanie przetestowana na 27 farmakokinetycznych przykładowych zbiorach danych.

1. Podstawy teoretyczne i terminologia

Z matematycznego punktu widzenia sygnałem nazywamy zmianę jednej (lub kilku) wielkości fizycznie mierzalnych w zależności od zmian innej wielkości. W powszechnym użyciu pojęcie sygnału stosuje się do zmian pewnych fizycznie mierzalnych wielkości (np. chemicznych, biologicznych czy też biochemicznych) w funkcji czasu. Taki sygnał nazywamy szeregiem czasowym. Przypomnijmy, iż szereg czasowy to ciąg pomiarów (obserwacji) uporządkowany w czasie. W prezentowanej pracy przyjmujemy, iż zmienna niezależna (tj. zmienna czasowa) jest dyskretna. A więc, z formalnego punktu widzenia, szereg czasowy T to ciąg par uporządkowanych o postaci: $T = [(t_1, x_1), (t_2, x_2), \dots, (t_i, x_i), \dots, (t_n, x_n)]$ dla $t_1 < t_2 < \dots < t_i < \dots < t_n$ oraz dla zbioru indeksów $I = \{1, 2, \dots, i, \dots, n\}$, gdzie każdy wyraz x_i to wynik pomiaru (obserwacji) zmiennej zależnej w d -wymiarowej przestrzeni cech (ang. *feature space*) oraz każdy element t_i to punkt (krok) czasowy (ang. *timestamp*), w którym dany pomiar (dana obserwacja) został zarejestrowany. W prezentowanej pracy zakładamy, iż wszystkie analizowane dane temporalne są jednowymiarowe (ang. *univariate*), a więc $d = 1$ dla każdego z rozpatrywanych szeregów czasowych. Ponadto przyjmujemy, iż każdy badany ciąg temporalny jest regularny, a więc punkty czasowe t_i wszystkich rozpatrywanych szeregów czasowych są równomiernie rozmieszczone. Liczba zarejestrowanych pomiarów (obserwacji), (n), to długość szeregu czasowego T . Z kolei, i -ty wyraz szeregu czasowego T oznaczamy

przez $T(i)$. Etykietowany temporalny zbiór danych \mathbf{TD} to kolekcja $\mathbf{T} = \{T_g\}_{g=1}^m$ szeregów czasowych z których każdy ma długość n razem z określonym (ang. *predefined*) dyskretnym wektorem² etykiet \mathbf{C} o długości m . Konkludując, etykietowany temporalny zbiór danych to struktura o postaci $\mathbf{TD} = (\mathbf{T}, \mathbf{C}) = (T_g, C^{T_g})$, gdzie $T_g \in \mathbf{T}$ to g -ty szereg czasowy, natomiast $C^{T_g} \in \mathbf{C}$ to jego etykieta. W przedkładanej pracy przyjmujemy, iż dwie etykiety $C^{T_g}, C^{T_h} \in \mathbf{C}$ są równoważne (\cong) wtedy i tylko wtedy gdy są identyczne, tj. $C^{T_g} \cong C^{T_h} \leftrightarrow C^{T_g} = C^{T_h}$, gdzie \leftrightarrow to logiczny spójnik równoważności międzyzdaniowej. Dla szeregu czasowego T o długości n , jego pierwsza dyskretna pochodna T' to nowy szereg czasowy o długości $n - 1$, którego wyrazy dane są zależnością³: $T'(i) := T(i + 1) - T(i)$, gdzie $i = 1, 2, \dots, n - 1$.

3. Nowe algorytmy obliczania odległości pomiędzy szeregami czasowymi

Proponowane w tej części pracy nowe techniki obliczania odległości pomiędzy danymi temporalnymi oparte są na analizie rekurencyjnej (ang. *recurrence quantification analysis*) szeregów czasowych RQA^4 , jej krzyżowej formie (ang. *cross-recurrence quantification analysis*) $CRQA$ oraz na metodzie obliczania odległości pomiędzy wielowymiarowymi zbiorami danych opracowanej przez Roberta Todeschiniego oraz jego współpracowników⁵. Techniki te składają się z następujących kroków:

1. Dla każdego szeregu czasowego w analizowanej bazie danych, tj. $T_g \in \mathbf{T}$ obliczamy jego bezprogową macierz rekurencyjną (ang. *unthresholded recurrence matrix*) UR^g lub jej krzyżową formę (ang. *unthresholded cross-recurrence matrix*) UCR^g zgodnie z następującymi formułami⁶:

$$UR_{ij}^g := \|T_g(i) - T_g(j)\|_2$$

oraz

² Z matematycznego punktu widzenia wektor \mathbf{C} to multizbiór.

³ T. Górecki, M. Łuczak, *Multivariate time series classification with parametric derivative dynamic time warping*, „Expert Systems with Applications” 2015, nr 42, s. 2307.

⁴ J.-P. Eckmann, S. O. Kamphorst, D. Ruelle, *Recurrence plots of dynamical systems*, „Europhysics Letters” 1987, nr 5, s. 973–977; N. Marwan, M. C. Romano, M. Thiel, J. Kurths, *Recurrence plots for the analysis of complex systems*, „Physics Reports” 2007, nr 438, s. 237–329.

⁵ R. Todeschini, D. Ballabio, V. Consonni, A. Manganaro, A. Mauri, *Canonical measure of correlation (CMC) and canonical measure of distance (CMD) between sets of data. Part 1. Theory and simple chemometric applications*, „Analytica Chimica Acta” 2009, nr 648, s. 45–51.

⁶ N. Marwan *et al.*, *op. cit.*

$$UCR_{ij}^g := \|T_g(i) - T_g(j)\|_2,$$

gdzie $\|\cdot\|_2$ to norma *euklidesowa*, a $T_g(j)$ to j -ty wyraz pierwszej dyskretnej pochodnej T_g' szeregu temporalnego T_g . Macierz UR^g to symetryczna względem głównej przekątnej macierz kwadratowa o wymiarach $n \times n$, natomiast macierz UCR^g to prostokątna macierz o wymiarach $n \times (n - 1)$ ⁷.

2. W drugim kroku wyrazy otrzymanych macierzy, tj. UR_{ij}^g oraz UCR_{ij}^g , transformujemy według zależności⁸:

$$TUR_{ij}^g := e^{-UR_{ij}^g}$$

oraz

$$TUCR_{ij}^g := e^{-UCR_{ij}^g}$$

otrzymując transformowaną bezprogową macierz rekurencyjną (ang. *transformed unthresholded recurrence matrix*) TUR^g oraz jej krzyżową wersję (ang. *transformed unthresholded cross-recurrence matrix*) $TUCR^g$.

3. W trzecim kroku zapisujemy macierze TUR^g oraz $TUCR^g$ w formie kolumnowej⁹:

$$TUR^g = \begin{bmatrix} TUR_{11}^g & TUR_{12}^g & \dots & TUR_{1n}^g \\ TUR_{21}^g & TUR_{22}^g & \dots & TUR_{2n}^g \\ \dots & \dots & \ddots & \dots \\ \dots & \dots & \ddots & \dots \\ TUR_{n1}^g & TUR_{n2}^g & \dots & TUR_{nn}^g \end{bmatrix} = [TUR_{\cdot,1}^g, TUR_{\cdot,2}^g, \dots, TUR_{\cdot,n}^g]$$

oraz

⁷ Przypomnijmy, iż macierze rekurencyjne (progowe lub bezprogowe) zawsze są symetrycznymi macierzami kwadratowymi, mającymi na celu porównywanie odległości pomiędzy wyrazami tego samego szeregu czasowego. Natomiast, krzyżowe macierze rekurencyjne progowe lub bezprogowe zawsze służą do porównywania odległości pomiędzy wyrazami dwóch różnych szeregów czasowych i tylko w wyjątkowych przypadkach, gdy długości obu szeregów są równe, są one kwadratowe. Por. *Ibidem*. Bezprogowe krzyżowe macierze rekurencyjne prównujące odległości pomiędzy wyrazami oryginalnego (wejściowego) szeregu czasowego i jego pierwszej dyskretnej pochodnej rozpatrywane były w: P. Wilczek, *An application of the local binary pattern algorithm and its uniform variant to improve the recurrence and cross-recurrence quantification analyses of the pharmacologically important time series* [w:] *Recent Advances in Computational Oncology and Personalized Medicine. Vol. 2. The Challenges of the Future !*, red. K. Krukiewicz, M. Marczyk, M. Bugdol, S. Bajkacz, Z. Ostrowski, Gliwice 2022, s. 128–152.

⁸ D. Eroglu, T. K. DM. Peron, N. Marwan, F. A. Rodrigues, L. da F. Costa, M. Sebek, I. Z. Kiss, J. Kurths, *Entropy of weighted recurrence plots*, „Physical Review E” 2014, nr 90, s. 042919.

⁹ Por. formuły 3 i 4 w R. Todeschini *et al.*, *op. cit.*

$$TUR^g = \begin{bmatrix} TUCR_{11}^g & TUCR_{12}^g & \dots & TUCR_{1(n-1)}^g \\ TUCR_{21}^g & TUCR_{22}^g & \dots & TUCR_{2(n-1)}^g \\ \dots & \dots & \ddots & \dots \\ \dots & \dots & \ddots & \dots \\ TUCR_{n1}^g & TUCR_{n2}^g & \dots & TUCR_{n(n-1)}^g \end{bmatrix} = [TUCR_{\cdot 1}^g, TUCR_{\cdot 2}^g, \dots, TUCR_{\cdot (n-1)}^g].$$

W powyższym zapisie symbole TUR_i^g oraz $TUCR_i^g$ oznaczają i -tą kolumnę, odpowiednio macierzy TUR^g oraz macierzy $TUCR^g$.

4. W czwartym kroku obliczamy macierz współczynników korelacji Pearsona r pomiędzy dwiema macierzami TUR^g oraz TUR^h , odpowiadającymi dwóm szeregom czasowym $T_g, T_h \in \mathbf{T}$, pomiędzy którymi chcemy obliczyć odległość rekurencyjno-kanoniczną, lub pomiędzy dwiema macierzami $TUCR^g$ oraz $TUCR^h$, odpowiadającymi dwóm szeregom czasowym $T_g, T_h \in \mathbf{T}$, pomiędzy którymi chcemy obliczyć krzyżową odległość rekurencyjno-kanoniczną. W pierwszym przypadku macierz korelacji oznaczona jest przez r_{TUR^g, TUR^h} , a w drugim przez $r_{TUCR^g, TUCR^h}$. Mają one postać:

$$r_{TUR^g, TUR^h} = \begin{bmatrix} r(TUR_{\cdot 1}^g, TUR_{\cdot 1}^h) & r(TUR_{\cdot 1}^g, TUR_{\cdot 2}^h) & \dots & r(TUR_{\cdot 1}^g, TUR_{\cdot n}^h) \\ r(TUR_{\cdot 2}^g, TUR_{\cdot 1}^h) & r(TUR_{\cdot 2}^g, TUR_{\cdot 2}^h) & \dots & r(TUR_{\cdot 2}^g, TUR_{\cdot n}^h) \\ \dots & \dots & \ddots & \dots \\ \dots & \dots & \ddots & \dots \\ r(TUR_{\cdot n}^g, TUR_{\cdot 1}^h) & r(TUR_{\cdot n}^g, TUR_{\cdot 2}^h) & \dots & r(TUR_{\cdot n}^g, TUR_{\cdot n}^h) \end{bmatrix}$$

oraz

$$r_{TUCR^g, TUCR^h} = \begin{bmatrix} r(TUCR_{\cdot 1}^g, TUCR_{\cdot 1}^h) & r(TUCR_{\cdot 1}^g, TUCR_{\cdot 2}^h) & \dots & r(TUCR_{\cdot 1}^g, TUCR_{\cdot (n-1)}^h) \\ r(TUCR_{\cdot 2}^g, TUCR_{\cdot 1}^h) & r(TUCR_{\cdot 2}^g, TUCR_{\cdot 2}^h) & \dots & r(TUCR_{\cdot 2}^g, TUCR_{\cdot (n-1)}^h) \\ \dots & \dots & \ddots & \dots \\ \dots & \dots & \ddots & \dots \\ r(TUCR_{\cdot (n-1)}^g, TUCR_{\cdot 1}^h) & r(TUCR_{\cdot (n-1)}^g, TUCR_{\cdot 2}^h) & \dots & r(TUCR_{\cdot (n-1)}^g, TUCR_{\cdot (n-1)}^h) \end{bmatrix}.$$

Wyraz na pozycji (i, j) macierzy r_{TUR^g, TUR^h} to współczynnik korelacji liniowej Pearsona pomiędzy i -tą kolumną macierzy TUR^g a j -tą kolumną macierzy TUR^h . A więc w ogólnym przypadku ma on postać $r(TUR_{\cdot i}^g, TUR_{\cdot j}^h)$ ¹⁰. Ponieważ obie macierze mają n kolumn, macierz r_{TUR^g, TUR^h} ma wymiary $n \times n$. Podobnie, wyraz na pozycji (i, j) macierzy $r_{TUCR^g, TUCR^h}$

¹⁰ Por. formuła 5 w *Ibidem*.

to współczynnik korelacji liniowej Pearsona pomiędzy i -tą kolumną macierzy $TUCR^g$ a j -tą kolumną macierzy $TUCR^h$. Wtedy ma on postać $r(TUCR^g_i, TUCR^h_j)$. Ponieważ obie macierze mają $n - 1$ kolumn, macierz ich współczynników korelacji $r_{TUCR^g, TUCR^h}$ ma wymiary $(n - 1) \times (n - 1)$. Zauważmy, iż w ogólnym przypadku $r_{A,B} \neq r_{B,A}$, gdzie A oraz B to dowolne macierze mające równą ilość rzędów.

5. W piątym kroku obie macierze współczynników korelacji liniowych są symetryzowane względem głównej przekątnej za pomocą mnożenia macierzowego¹¹:

$$Q_{RCMD} = r_{TUR^g, TUR^h} \times r_{TUR^h, TUR^g}$$

lub

$$Q_{RCMD}^* = r_{TUR^h, TUR^g} \times r_{TUR^g, TUR^h},$$

gdzie $(r_{TUR^g, TUR^h})^T = r_{TUR^h, TUR^g}$

oraz

$$Q_{CRCMD} = r_{TUCR^g, TUCR^h} \times r_{TUCR^h, TUCR^g}$$

lub

$$Q_{CRCMD}^* = r_{TUCR^h, TUCR^g} \times r_{TUCR^g, TUCR^h},$$

gdzie $(r_{TUCR^g, TUCR^h})^T = r_{TUCR^h, TUCR^g}$. Działanie $(\cdot)^T$ to transpozycja macierzowa. Macierze produktowe Q_{RCMD} oraz Q_{RCMD}^* mają wymiary $n \times n$, natomiast macierze Q_{CRCMD} oraz Q_{CRCMD}^* mają wymiary $(n - 1) \times (n - 1)$.

6. W szóstym kroku obliczymy wartości własne macierzy produktowych Q_{RCMD} (lub Q_{RCMD}^*) oraz Q_{CRCMD} (lub Q_{CRCMD}^*). Macierz produktowa Q_{RCMD} ma takie same niezerowe wartości własne jak macierz Q_{RCMD}^* . Podobnie, macierz produktowa Q_{CRCMD} ma takie same niezerowe wartości własne jak macierz Q_{CRCMD}^* .

7. W siódmym kroku obliczamy wartości odległości rekurencyjno-kanonicznej pomiędzy szeregami czasowymi $T_g, T_h \in \mathbf{T}$ według wzoru¹²:

¹¹ Por. Formuła 6 w *Ibidem*.

¹² W oryginalnej pracy Todeschiniego i współautorów kanoniczna miara odległości $CMD(A, B)$ pomiędzy dwoma wielowymiarowymi zbiorami danych A i B wyrażona jest formułą: $CMD(A, B) := p_A + p_B - 2 \sum_j^k \sqrt{\lambda_j}$, gdzie p_A oraz p_B to, odpowiednio liczba zmiennych (kolumn) w zbiorze A i w zbiorze B , λ_i to wartości własne

$$RCMD(T_g, T_h) := 2n - 2 \sum_{j=1}^k \sqrt{\lambda_j^{RCMD}},$$

gdzie n to długość obu szeregów, λ_j^{RCMD} to wartości własne macierzy Q_{RCMD} (lub Q_{RCMD}^*), a indeks k równy jest ilości niezerowych wartości własnych λ_j^{RCMD} . Przy tych założeniach zachodzi, iż dla dowolnych $T_g, T_h \in \mathbf{T}$, $0 \leq RCMD(T_g, T_h) \leq 2n$. Z kolei krzyżowa odległość rekurencyjno-kanoniczna pomiędzy szeregami $T_g, T_h \in \mathbf{T}$ dana jest zależnością:

$$CRCMD(T_g, T_h) := 2(n - 1) - 2 \sum_{j=1}^k \sqrt{\lambda_j^{CRCMD}},$$

gdzie n to długość obu szeregów, λ_j^{CRCMD} to wartości własne macierzy Q_{CRCMD} (lub Q_{CRCMD}^*), a indeks k równy jest ilości niezerowych wartości własnych λ_j^{CRCMD} . W tym przypadku zachodzi, iż dla dowolnych $T_g, T_h \in \mathbf{T}$, $0 \leq CRCMD(T_g, T_h) \leq 2(n - 1)$.

Aby zilustrować powyższy algorytm, rozpatrzmy dwa czteroelementowe szeregi czasowe T_1 oraz T_2 , tj. $T_1 = \langle 1, 3, 8, 5 \rangle$ oraz $T_2 = \langle 2, 0, 9, 7 \rangle$. Ich bezprogowe macierze rekurencyjne oraz ich transponowane formy mają postać:

$$UR^1 = \begin{bmatrix} 0 & 2 & 7 & 4 \\ 2 & 0 & 5 & 2 \\ 7 & 5 & 0 & 3 \\ 4 & 2 & 3 & 0 \end{bmatrix}, UR^2 = \begin{bmatrix} 0 & 2 & 7 & 5 \\ 2 & 0 & 9 & 7 \\ 7 & 9 & 0 & 2 \\ 5 & 7 & 2 & 0 \end{bmatrix}$$

oraz

$$TUR^1 = \begin{bmatrix} 1 & 0,1353 & 0,0009 & 0,0183 \\ 0,1353 & 1 & 0,0067 & 0,1353 \\ 0,0009 & 0,0067 & 1 & 0,0498 \\ 0,0183 & 0,1353 & 0,0498 & 1 \end{bmatrix},$$

$$TUR^2 = \begin{bmatrix} 1 & 0,1353 & 0,0009 & 0,0067 \\ 0,1353 & 1 & 0,0001 & 0,0009 \\ 0,0009 & 0,0001 & 1 & 0,1353 \\ 0,0067 & 0,0009 & 0,1353 & 1 \end{bmatrix}.$$

Z kolei macierze korelacji liniowych Pearsona mają postać:

odpowiednich macierzy produktowych, a indeks k równy jest ilości niezerowych wartości własnych λ_i . Por. formuła 7 w *Ibidem*.

$$r_{TUR^1, TUR^2} = \begin{bmatrix} 0,9999 & -0,083 & -0,4688 & -0,4424 \\ -0,1417 & 0,9913 & -0,5105 & -0,3441 \\ -0,4229 & -0,4137 & 0,9961 & -0,1655 \\ -0,445 & -0,3001 & -0,2315 & 0,9821 \end{bmatrix}$$

oraz

$$(r_{TUR^1, TUR^2})^T = r_{TUR^2, TUR^1} = \begin{bmatrix} 0,9999 & -0,1417 & -0,4229 & -0,445 \\ -0,083 & 0,9913 & -0,4137 & -0,3001 \\ -0,4688 & -0,5105 & 0,9961 & -0,2315 \\ -0,4424 & -0,3441 & -0,1655 & 0,9821 \end{bmatrix}$$

Obie macierze produktowe dane są poniżej:

$$Q_{RCMD} = r_{TUR^1, TUR^2} \times r_{TUR^2, TUR^1} = \begin{bmatrix} 1,4223 & 0,1676 & -0,7823 & -0,746 \\ 0,1676 & 1,3819 & -0,8018 & -0,4542 \\ -0,7823 & -0,8018 & 1,3696 & -0,0808 \\ -0,746 & -0,4542 & -0,0808 & 1,3061 \end{bmatrix}$$

oraz

$$Q_{RCMD}^* = r_{TUR^2, TUR^1} \times r_{TUR^1, TUR^2} = \begin{bmatrix} 1,3968 & 0,0851 & -0,7147 & -0,7606 \\ 0,0851 & 1,2509 & -0,8099 & -0,5306 \\ -0,7147 & -0,8099 & 1,5262 & -0,0091 \\ -0,7606 & -0,5306 & -0,0091 & 1,306 \end{bmatrix}$$

Wartości własne obu macierzy produktowych Q_{RCMD} oraz Q_{RCMD}^* są równe 2,838209; 1,513631; 1,128044; 0. Trzy pierwsze z nich służą do obliczenia wartości odległości $RCMD(T_1, T_2)$. Wynosi ona 0,04582047. Odległość $CRCMD(T_1, T_2)$ obliczana jest w sposób podobny. Wynosi ona 2,256611.

Nowo zaproponowane protokoły obliczania odległości pomiędzy danymi temporalnymi będą testowane w przykładowych zadaniach klasyfikacyjnych w części czwartej poniższej pracy. Przypomnijmy, iż dla kolekcji $\mathbf{T} = \{T_g\}_{g=1}^m$ szeregów czasowych i dyskretnego wektora ich etykiet \mathbf{C} (w naszym przypadku jest to wektor binarny) zagadnienie klasyfikacji polega na aproksymacji funkcji o postaci: $c: \mathbf{T} \rightarrow \mathbf{C}$ (gdzie $c(T_g) = C^{T_g}$ dla $\forall T_g \in \mathbf{T}$) funkcją $\hat{c}: \mathbf{T} \rightarrow \mathbf{C}$, tak aby spełniony był warunek: $\hat{c}(T_g) = C^{T_g}$ dla $\forall T_g \in \mathbf{T}$. Funkcja c to klasyfikator aprioryczny (ang. *predefined*), natomiast funkcja \hat{c} to klasyfikator aposterioryczny¹³. A więc, jeżeli:

¹³ Oczywiście pojęcia klasyfikatora apriorycznego oraz aposteriorycznego mają sens tylko w odniesieniu do klasyfikacji formalnej (matematycznej). Oznacza to, iż chociaż elementy dyskretnego wektora etykiet \mathbf{C} mogły zostać przyporządkowane poszczególnym danym temporalnym na podstawie procedur empirycznych, to podczas

$c(T_g) = \hat{c}(T_g)$ dla $\forall T_g \in \mathbf{T}$, to mówimy, że klasyfikator aposterioryczny w pełni (perfekcyjnie, doskonale) aproksymuje klasyfikator aposterioryczny. Natomiast jeżeli tylko dla większości $T_g \in \mathbf{T}$ zachodzi warunek $c(T_g) = \hat{c}(T_g)$, to mówimy, iż klasyfikator aposterioryczny tylko częściowo przybliży klasyfikator aprioryczny. Problem klasyfikacji polega na znalezieniu jak najdokładniejszego klasyfikatora aposteriorycznego. Tym samym wydajność klasyfikatora aposteriorycznego może być mierzona w procentach poprawnie sklasyfikowanych szeregów czasowych w analizowanym temporalnym zbiorze danych. Postępując zgodnie z sugestią Eamonna Keogha oraz ShrutiKasetty¹⁴, w części czwartej poniższej pracy nowo zaproponowane funkcje odległości pomiędzy danymi temporalnymi będą testowane przez pryzmat klasycznego klasyfikatora jednego najbliższego sąsiada (ang. *1-nearest neighbor classifier* (1NN)) na 27 przykładowych farmakokinetycznych zbiorach danych.

4. Porównawcze algorytmy obliczania odległości pomiędzy szeregami czasowymi

Zgodnie z ogólnie przyjętą metodologią walidacji nowo zaproponowanych miar odległości pomiędzy danymi temporalnymi¹⁵ w części czwartej prezentowanej pracy porównamy efektywność nowych technik klasyfikacji szeregów czasowych z efektywnością technik bazujących na klasycznych miarach odległości oraz na niedawno wprowadzonych funkcjach odległości. Mianowicie, wydajność w przykładowych zadaniach klasyfikacyjnych nowych algorytmów będzie porównana z wydajnością algorytmu euklidesowego ($L2$), jego trzech modyfikacji ($DL2$, $PDL2$ oraz $CIL2$), algorytmu DTW oraz jego trzech ulepszeń ($DDTW$, $PDDTW$ oraz $CIDTW$). Przypomnijmy, iż dla dwóch szeregów czasowych T_g oraz T_h , których wyrazy indeksowane są tym samym zbiorem I , algorytm euklidesowy ma postać $L2(T_g, T_h) := \sqrt{\sum_{i=1}^n [T_g(i) - T_h(i)]^2}$. Z kolei miara DTW obliczana jest według procedury zaimplementowanej w pakietach obliczeniowych języka R , tj. w pakietach dtw ¹⁶ oraz $TSdist$ ¹⁷. Dla dwóch szeregów czasowych T_g oraz T_h , których wyrazy indeksowane są tym samym zbiorem I , niech $M(T_g, T_h)$ będzie kwadratową $n \times n$ dodatnią macierzą odległości (tzw.

matematycznego procesu klasyfikacji owo przyporządkowanie jest czymś z góry ustalonym (ang. *predefined*), a więc z wewnętrznego punktu widzenia klasyfikacji formalnej jest to przyporządkowanie aprioryczne.

¹⁴ E. Keogh, S. Kasetty, *On the need for time series data mining benchmarks: a survey and empirical demonstration* [w:] *Proceedings of the eight ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, New York 2002, s. 102–111.

¹⁵ *Ibidem*, a także P. Wilczek, *op. cit.*

¹⁶ T. Giorgino, *Computing and visualizing dynamic time warping alignments in R: The dtw package*, „Journal of Statistical Software” 2009, nr (7), s. 1–24.

¹⁷ U. Mori, A. Mendiburu, J. A. Lozano, *Distance measures for time series in R: The TSdist Package*, „The R Journal” 2016, nr 2, s. 451–459.

macierzą kosztów lokalnych), której wyraz M_{ij} dany jest zależnością: $M_{ij} := |T_g(i) - T_g(j)|$. Wtedy tak zwana krzywa (ang. *warping path*) P o postaci: $P := [(e_1, f_1), (e_2, f_2), \dots, (e_i, f_i), \dots, (e_s, f_s)]$ to ciąg punktów (tj. par indeksów) definiujących trawersalę macierzy $M(T_g, T_h)$. W definicji odległości *DTW* zakładamy, iż krzywa P musi spełniać następujące warunki: $(e_1, f_1) = M_{11}$ oraz $(e_s, f_s) = M_{nn}$, jak również $0 \leq e_{i+1} - e_i \leq 1$ oraz $0 \leq f_{i+1} - f_i \leq 1$ dla $\forall i < n$. A więc krzywa P to ścieżka biegnąca od elementu M_{11} macierzy $M(T_g, T_h)$ do elementu M_{nn} tejże macierzy, przy czym przy przejściu od wyrazu M_{11} do wyrazu M_{nn} dozwolone są tylko kroki: $(0,1)$, $(1,0)$ oraz $(1,1)$. Niech $p_i = \gamma M_{e_i f_i}$ będzie odległością globalną (tj. kumulacyjną) pomiędzy elementem w pozycji e_i szeregu T_g a elementem w pozycji f_i szeregu T_h dla i -tej pary punktów (tj. indeksów) na krzywej P . Wtedy odległość pomiędzy szeregami czasowymi T_g i T_h wzdłuż krzywej P , D_P , dana jest zależnością $D_P(T_g, T_h) := \sum_{i=1}^s p_i$. Niech \mathcal{P} będzie przestrzenią wszystkich możliwych krzywych P . Wtedy ścieżka *DTW*, P^* , to krzywa minimalizująca odległość D_P , tj. $P^* := \min_{P \in \mathcal{P}} D_P(T_g, T_h)$. Odległość *DTW* pomiędzy dwoma ciągami temporalnymi T_g oraz T_h to odległość wzdłuż ścieżki P^* ¹⁸. Można ją znaleźć za pomocą procedur programowania dynamicznego¹⁹ na podstawie następującego warunku rekurencyjnego, umożliwiającego obliczenie globalnej odległości $\gamma M_{e_i f_i}$ pomiędzy elementami szeregów T_g oraz T_h na podstawie znajomości wartości kosztów lokalnych: $\gamma M_{e_1 f_1} = M_{11}$ oraz dla $i, j \neq 1$ $\gamma M_{e_i f_i} = M_{ij} + \min\{\gamma M_{e_{i-1} f_{i-1}}, \gamma M_{e_{i-1} f_i}, \gamma M_{e_i f_{i-1}}\}$. A więc obliczenie kumulacyjne odległości $\gamma M_{e_i f_i}$ pomiędzy elementem w pozycji e_i szeregu T_g a elementem w pozycji f_i szeregu T_h dla i -tej pary indeksów na krzywej P redukuje się do znajomości odległości lokalnej bieżącego elementu M_{ij} oraz kumulacyjnych odległości elementów przyległych (ang. *adjacent*), tj. odległości $\gamma M_{e_{i-1} f_{i-1}}$, $\gamma M_{e_{i-1} f_i}$ oraz $\gamma M_{e_i f_{i-1}}$. Dla dwóch szeregów czasowych T_g oraz T_h , których wyrazy indeksowane są tym samym zbiorem I , oraz funkcji $DF \in \{L2, DTW\}$, tzw. algorytm odległości pochodnej *DDF* ma postać $DDF(T_g, T_h) := DF(T'_g, T'_h)$, gdzie T'_g oraz T'_h

¹⁸ D. J. Berndt, J. Clifford, *Using dynamic time warping to find patterns in time series* [w:] *Proceedings of the 3rd International Conference on Knowledge Discovery and Data Mining*, 1994, s. 359–370.

¹⁹ Wyczerpujące informacje na temat procedur programowania dynamicznego czytelnik znajdzie w następujących pozycjach: T. Giorgino, *op. cit.*, E. Keogh, C. A. Ratanamahatana, *Exact indexing of dynamic time warping*, „Knowledge and Information Systems” 2005, nr 7, s. 358–386, L. Rabiner, B.–H. Juang, *Fundamentals of Speech Recognition*, New Jersey 1993, D. Sankoff, J. Kruskal, *Time Warps, String Edits, and Macromolecules. The Theory and Practice of Sequence Comparison*, Stanford 1999.

to pierwsze dyskretne pochodne szeregów T_g i T_h ²⁰. Z kolei dla dwóch szeregów czasowych T_g oraz T_h , których wyrazy indeksowane są tym samym zbiorem I , oraz funkcji $DF \in \{L2, DTW\}$, tzw. algorytm pochodnej odległości parametrycznej $PDDF$ to ważona wypukła kombinacja o postaci $PDDF(T_g, T_h) := aDF(T_g, T_h) + bDDF(T_g, T_h)$, gdzie DDF to algorytm odległości pochodnej oraz $b = 1 - a$, gdzie $a \in [0,1]$ to parametry rzeczywiste²¹. W części czwartej poniższej pracy algorytm $PDDF$ obliczany jest względem parametrów równych: $a = b = \cos\alpha = \sin\alpha = 0,7071068$, gdzie $\alpha = \frac{\pi}{2}$. Natomiast dla dwóch szeregów czasowych T_g oraz T_h , których wyrazy indeksowane są tym samym zbiorem I , oraz funkcji $DF \in \{L2, DTW\}$, tzw. algorytm niezmienniczy ze względu na złożoność szeregów czasowych (ang. *complexity-invariant*) ma postać $CIDF(T_g, T_h) := DF(T_g, T_h) \times CF(T_g, T_h)$, gdzie $CF(T_g, T_h)$ to tzw. czynnik korygujący złożoność (ang. *complexity correction factor*), obliczany dla T_g oraz T_h według zależności: $CF(T_g, T_h) := \frac{\max\{CE(T_g), CE(T_h)\}}{\min\{CE(T_g), CE(T_h)\}}$. W powyższej zależności $CE(T)$ (dla $T \in \{T_g, T_h\}$) to tzw. wyznacznik złożoności (ang. *complexity estimate*) szeregu T . Ma on postać:

$$CE(T) := \sqrt{\sum_{i=1}^{n-1} [T(i) - T(i-1)]^2} \text{ dla } i = 1, 2, \dots, n - 1^{22}.$$

5. Materiały i metody

Nowo zaproponowane algorytmy obliczania funkcji odległości pomiędzy szeregami czasowymi zostaną przetestowane na części publicznie dostępnego zbioru danych HTS007²³. Zbiór ten został wygenerowany w laboratorium *High Throughput Screening* Uniwersytetu Vanderbilta²⁴. Obejmuje on pięciodniowe pomiary proliferacji komórek przeprowadzone na 8 liniach komórkowych (BT20, HCC1143, MCF10A-HMS, MCF10A-VU, MDAMB231, MDAMB453, MDAMB468 oraz SUM149) traktowanych 27 lekami przeciwnowotworowymi o różnym stężeniu (próba testowa) oraz nietraktowanych żadnym lekiem (próba kontrolna). Do naszych symulacji komputerowych użyliśmy tylko szeregów czasowych, będących wynikami

²⁰ E. J. Keogh, M. J. Pazzani, *Dynamic time warping with higher order features* [w:] *Proceedings of the 2001 SIAM International Conference on Data Mining*, red. V. Kumar, R. Grossman, 2001, s. 1–11.

²¹ T. Górecki, M. Łuczak, *Using derivatives in time series classification*, „Data Mining and Knowledge Discovery” 2013, nr 2, s. 310–331.

²² G. E. A. P. A. Batista, X. Wang, E. J. Keogh, *A complexity-invariant distance measure for time series* [w:] *Proceedings of the 2011 SIAM International Conference on Data Mining*, red. B. Liu, H. Liu, C. Clifton, T. Washio, C. Kamath, 2011, s. 699–710.

²³ Thunor, <https://www.thunor.net>, (on-line 16.02.2024).

²⁴ A. L. R. Lubbock, L. A. Harris, V. Quaranta, D. R. Tyson, C. F. Lopez, *Thunor: visualization and analysis of high-throughput dose-response datasets*, „Nucleic Acids Research” 2021, nr 46, s. w633–w640.

pomiarów przeprowadzonych na linii BT20. Nazwy poszczególnych zbiorów danych pomiarowych odpowiadają nazwą leków testowanych na komórkach, których te pomiary dotyczą. Są one wyszczególnione w Tabeli 1. Kolumny Test oraz Kontrola odpowiadają kolejno liczbie testowych oraz kontrolnych szeregów czasowych w każdym ze zbiorów danych.

Lp.	<i>TD</i>	Test	Kontrola	Lp.	<i>TD</i>	Test	Kontrola
1	Abemaciclib	20	31	15	Neratinib	15	31
2	Alpelisib	15	31	16	Osimertinib	20	31
3	Azd7762	20	31	17	Paclitaxel	20	31
4	Bleomycin	20	31	18	Palbociclib	17	31
5	Buparlisib	19	31	19	Panobinostat	19	31
6	Cediranib	17	31	20	Pictilisib	17	31
7	Certitinib	19	31	21	Saracatinib	17	31
8	Dasatinib	17	31	22	Taselisib	20	31
9	Doxorubicin	20	31	23	Tivantinib	19	31
10	Etoposide	19	31	24	Torin2	20	31
11	Everolimus	19	31	25	Trametinib	16	31
12	Ink128	15	31	26	Volasertib	20	31
13	Ipatasertib	18	31	27	Vorinostat	19	31
14	Luminespib	20	31				

Tab. 1 Przykładowe farmakokinetyczne zbiory danych ze zbioru HTS007

Każdy z 27 pojedynczych zbiorów danych składa się z 31 dwudziestoczerolementowych kontrolnych szeregów czasowych oraz z od 15 do 20 (por. Tabela 1) dwudziestoczerolementowych testowych szeregów czasowych. Tym samym każdy zbiór danych ma dwie etykiety: 1/ kontrola oraz 2/ test. Wszystkie obliczenia zostały wykonane w języku programowania R oraz w jego pakietach²⁵.

²⁵ T. Giorgino *op. cit.*, U. Mori *et al.*, *op. cit.*, R Core Team, *R: A language and environment for statistical computing*. R Foundation for Statistical Computing, Vienna 2022, <https://www.R-project.org/> (on-line 16.02.2024), W. N. Venables, B. D. Ripley, *Modern Applied Statistics with S*, New York 2002.

6. Wyniki

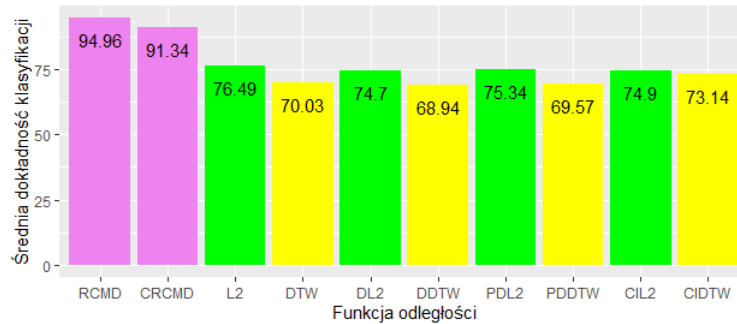
Tabela 2 zawiera wyniki analizy porównawczej wydajności algorytmów klasyfikacyjnych opartych na nowo zaproponowanych funkcjach odległości, tj. na funkcjach *RCMD* oraz *CRCMD*, jak również na referencyjnych miarach typu *L2* (tj. miarach *L2*, *DL2*, *PDL2* oraz *CIL2*) oraz typu *DTW* (tj. miarach *DTW*, *DDTW*, *PDDTW* oraz *CIDTW*). Wykres na Rycinie 1 podsumowuje wyniki z Tabeli 2. Dane te jednoznacznie wskazują, iż w 96,3% pomiarów efektywność techniki *RCMD* jest równa lub wyższa niż 90% poprawnie sklasyfikowanych szeregów czasowych w każdym zbiorze danych. Natomiast, skuteczność techniki *CRCMD* jest równa lub wyższa niż 90% poprawnie sklasyfikowanych szeregów w 66,67% testowanych przypadków. Z drugiej strony, na podstawie danych z Tabeli 2 można wywnioskować, iż efektywność technik typu *L2* jest równa lub wyższa niż 90% tylko w 13,89% dokonanych pomiarów. Z kolei wydajność technik typu *DTW* jest równa lub wyższa niż 90% poprawnie sklasyfikowanych szeregów tylko w 5,56 % przeprowadzonych symulacji.

<i>TD</i>	<i>RCMD</i>	<i>CRCMD</i>	<i>L2</i>	<i>DTW</i>	<i>DL2</i>	<i>DDTW</i>	<i>PDL2</i>	<i>PDDTW</i>	<i>CIL2</i>	<i>CIDTW</i>
1	90,2	<u>92,16</u>	68,6 3	64,7 1	70,5 9	58,82	74,5 1	62,75	68,6 3	74,51
2	<u>97,83</u>	91,3	69,5 7	58,7	63,0 4	58,7	63,0 4	56,52	69,5 7	58,7
3	<u>100</u>	96,08	68,6 3	62,7 5	66,6 7	50,98	68,6 3	64,71	58,8 2	64,71
4	<u>96,08</u>	<u>96,08</u>	72,5 5	56,8 6	49,0 2	60,78	60,7 8	52,94	68,6 3	66,67
5	<u>92</u>	88	68	58	80	56	64	50	62	74
6	<u>91,67</u>	87,5	83,3 3	70,8 3	66,6 7	70,83	75	68,75	70,8 3	66,67
7	<u>98</u>	96	66	58	66	64	62	64	64	70
8	<u>100</u>	85,42	85,4 2	81,2 5	83,3 3	62,5	85,4 2	79,17	83,3 3	87,5
9	<u>92,16</u>	88,24	72,5 5	78,4 3	72,5 5	68,63	72,5 5	74,51	68,6 3	76,47
10	<u>96</u>	92	68	68	64	62	68	62	72	54
11	98	<u>100</u>	<u>100</u>	<u>100</u>	98	<u>100</u>	<u>100</u>	<u>100</u>	<u>100</u>	<u>100</u>

12	91,3	<u>93,48</u>	91,3	89,1	91,3	89,13	<u>93,4</u>	89,13	89,1	93,48
				3			<u>8</u>		3	
13	<u>95,92</u>	93,88	79,5	63,2	75,5	53,06	69,3	63,27	73,4	63,27
			9	7	1		9		7	
14	<u>92,16</u>	84,31	68,6	74,5	84,3	76,47	72,5	72,55	70,5	62,75
			3	1	1		5		9	
15	<u>100</u>	93,48	73,9	50	80,4	65,22	76,0	60,87	71,7	56,52
			1		3		9		4	
16	<u>96,08</u>	92,16	72,5	68,6	64,7	66,67	74,5	74,51	78,4	74,51
			5	3	1		1		3	
17	80,39	<u>96,08</u>	90,2	84,3	82,3	86,27	92,1	80,39	90,2	82,35
				1	5		6			
18	<u>100</u>	95,83	83,3	68,7	66,6	56,25	77,0	70,83	77,0	72,92
			3	5	7		8		8	
19	<u>90</u>	78	72	68	82	72	74	66	70	74
20	<u>97,92</u>	95,83	68,7	62,5	68,7	60,42	66,6	66,67	66,6	72,92
			5		5		7		7	
21	<u>95,83</u>	89,58	66,6	58,3	75	68,75	68,7	60,42	66,6	58,33
			7	3			5		7	
22	<u>94,12</u>	88,24	72,5	52,9	68,6	70,59	70,5	50,98	66,6	50,98
			5	4	3		9		7	
23	<u>96</u>	92	72	74	70	78	74	72	80	88
24	<u>90,2</u>	78,43	80,3	78,4	80,3	68,63	80,3	68,63	78,4	74,51
			9	3	9		9		3	
25	<u>100</u>	97,87	76,6	72,3	80,8	74,47	74,4	76,6	76,6	80,85
				4	5		7			
26	92,16	94,12	92,1	88,2	92,1	<u>96,08</u>	90,2	88,24	90,2	<u>96,08</u>
			6	4	6					
27	<u>100</u>	90	82	80	74	66	86	82	90	80

Tab. 2 Procentowe wyniki klasyfikacji 27 farmakokinetycznych zbiorów danych

Ponadto dane zawarte w Tabeli 2 wskazują, iż algorytmy *RCMD*, *CRCMD*, typu *L2* oraz typu *DTW* osiągają najlepsze wyniki spośród wszystkich analizowanych przypadków, odpowiednio w 81,48%, 18,52%, 3,7% oraz w 5,56% wykonanych pomiarów. Tak więc można stwierdzić, iż skuteczność nowo wprowadzonych algorytmów *RCMD* oraz *CRCMD* jest znacznie wyższa niż skuteczność metod bazujących na klasycznych miarach odległości *L2*, *DTW* oraz ich udoskonaleniach.



Ryc. 1. Wykres ilustrujący średnią procentową efektywność porównywanych funkcji odległości

Ściśle rzecz biorąc, zestawiając średnie wartości efektywności porównywanych protokołów klasyfikacyjnych, można stwierdzić, iż schemat *RCMD* średnio przewyższa najlepszy schemat typu *L2* (tj. algorytm *L2*) oraz typu *DTW* (tj. algorytm *CIDTW*) o odpowiednio, 19,45% oraz o 22,98% poprawnie sklasyfikowanych szeregów temporalnych (por. Ryc. 1). Natomiast drugi z nowo zaproponowanych protokołów, tj. protokół *CRCMD*, przewyższa średnio schematy *L2* oraz *CIDTW* o odpowiednio, 16,26% oraz o 19,93% poprawnie sklasyfikowanych ciągów czasowych (por. Ryc. 1).

7. Dyskusja

W swojej obszernej i ważnej pracy *The Elements of Statistical Learning. Data Mining, Inference, and Prediction*, Trevor Hastie, Robert Tibshirani oraz Jerome Friedman zauważają, iż „wyszczególnienie adekwatnej miary niepodobieństwa [pomiędzy analizowanymi obiektami – przyp. mój – P. W.] jest o wiele bardziej istotne w osiągnięciu sukcesu w procesie klasteryzacji niż wybór samego algorytmu analizy skupień. Ten aspekt problemu jest mniej podkreślany w literaturze poświęconej klasteryzacji, ponieważ aspekt ten zależy od specyficzności dziedziny przedmiotowej [poddawanej analizie skupień – przyp. mój – P.W.]

oraz jest mniej podatny na ogólne analizy”²⁶. W naszej opinii te deklaracje, choć odnoszą się do zagadnienia „klasyfikacji bez nadzoru” (ang. *unsupervised learning*), mogą być uogólnione względem każdego rodzaju klasyfikacji, z włączeniem zagadnienia klasyfikacji szeregów czasowych. Dlatego też, aby poprawić efektywność znanych algorytmów klasyfikacji szeregów czasowych (np. algorytmu $L2$ lub DTW połączonych z klasyfikatorem $1NN$), zaproponowaliśmy nowe miary odległości pomiędzy danymi temporalnymi. Z przeprowadzonych symulacji komputerowych wynika, iż wydajność algorytmów klasyfikacji opartych na nowych bezparametrycznych funkcjach odległości jest znacznie wyższa niż wydajność protokołów opartych na funkcjach referencyjnych (np. na klasycznych miarach $L2$, DTW oraz na ich parametrycznych modyfikacjach). Rezultat ten wydaje się być bardzo istotny, ponieważ jak zauważają Tomasz Górecki oraz Maciej Łuczak „[...] prosta metoda łącząca klasyfikator jednego najbliższego sąsiada ($1NN$) oraz pewną formę miary odległości DTW okazała się być jedną z najwydajniejszych technik klasyfikacji szeregów czasowych. [...]. Euklidesowa miara odległości ma kilka zalet. Mianowicie, złożoność wyznaczania tej metryki jest liniowa, jest ona łatwa do zaimplementowania, może być łączona z dowolnymi innymi metodami oraz jest bezparametryczna. [...]. Zostało empirycznie dowiedzione, iż prosta Euklidesowa metryka odległości jest konkurencyjna lub lepsza względem wielu złożonych miar odległości oraz spełnia ważną nierówność trójkąta”²⁷. Ci sami autorzy w dalszych częściach swojego tekstu twierdzą, iż „metryka odległości Euklidesowej jest najbardziej oczywistą miarą podobieństwa dla szeregów czasowych, natomiast miara DTW jest jedną z najbardziej wydajnych funkcji odległości dla danych temporalnych”²⁸. Powyższe fragmenty zaczerpnięte z bieżącej literatury dotyczącej analizy i klasyfikacji szeregów czasowych mogą być bezpośrednio skonfrontowane z wynikami zamieszczonymi w Tabeli 2 oraz na Rycinie 1. Z postulowanego zestawienia można jednoznacznie wywnioskować, iż rezultaty uzyskane przez proponowany w naszej pracy rekurencyjno-kanoniczny schemat klasyfikacji oraz jego krzyżowa forma są kontrprzykładami względem stwierdzenia Góreckiego i Łuczaka, a także trudno jest znaleźć metodologię klasyfikacji danych temporalnych przewyższającą swoją skutecznością metodologię opartą na mierze DTW . Można zaznaczyć, że proponowany w tej pracy algorytm jest na tyle elastyczny, iż dzięki jego modyfikacjom można określić nowe rodziny miar odległości opartych na dyskretnych

²⁶ T. Hastie, R. Tibshirani, J. Friedman, *The Elements of Statistical Learning. Data Mining, Inference, and Prediction*, New York 2009, s. 506.

²⁷ T. Górecki, M. Łuczak, *Using...*, *op. cit.*, s. 311.

²⁸ *Ibidem.*, s. 320.

pochodnych analizowanych szeregów czasowych. Np. dla szeregu czasowego T_g o długości n z rozpatrywanej kolekcji szeregów \mathbf{T} można zdefiniować następujące oparte na pochodnych macierze rekurencyjne: $UR_{ij}^{(1)g} := \|T'_g(i) - T'_g(j)\|_2$, $UR_{ij}^{(2)g} := \|T''_g(i) - T''_g(j)\|_2$, gdzie T''_g to druga dyskretna pochodna szeregu czasowego T_g ²⁹ oraz następujące oparte na pochodnych krzyżowe macierze rekurencyjne: $UCR_{ij}^{(0,2)g} := \|T_g(i) - T''_g(j)\|_2$ oraz $UCR_{ij}^{(1,2)g} := \|T'_g(i) - T''_g(j)\|_2$ ³⁰. Analizie tego typu 2D struktur poświęcona będzie oddzielna praca.

7. Uwagi końcowe

Podsumowując, można stwierdzić, iż cel przedkładanej pracy został osiągnięty, a nowe rekurencyjno-kanoniczne funkcje obliczania odległości pomiędzy danymi temporalnymi znajdują praktyczne zastosowanie w analizie i klasyfikacji danych farmakokinetycznych i tym samym przyczynią się do zwiększenia stopnia automatyzacji w badaniach podstawowych stosowanych.

Podziękowania: Autor powyższej pracy pragnie wyrazić podziękowania anonimowej recenzentce tekstu, której uwagi przyczyniły się do znacznego ulepszenia przedkładanej pracy.

Bibliografia

1. Batista G. E. A. P. A., Wang X., Keogh E. J., *A complexity-invariant distance measure for time series* [w:] *Proceedings of the 2011 SIAM International Conference on Data Mining*, red. B. Liu, H. Liu, C. Clifton, T. Washio, C. Kamath, 2011, s. 699–710.
2. Berndt D. J., Clifford J., *Using dynamic time warping to find patterns in time series* [w:] *Proceedings of the 3rd International Conference on Knowledge Discovery and Data Mining*, 1994, s. 359–370.
3. Eckmann J.–P., Kamphorst S. O., Ruelle D., *Recurrence plots of dynamical systems*, „Europhysics Letters” 1987, nr 5, s. 973–977.
4. Eroglu D., T. K. DM. Peron, N. Marwan, F. A. Rodrigues, L. da F. Costa, M. Sebek, I. Z. Kiss, J. Kurths, *Entropy of weighted recurrence plots*, „Physical Review E” 2014, nr 90, 042919-1–024919-7.

²⁹ Druga dyskretna pochodna T'' szeregu czasowego T o długości n to nowy szereg czasowy o długości $n - 2$, będący pierwszą dyskretną pochodną pierwszej dyskretnej pochodnej T' wejściowego (oryginalnego) szeregu T .

³⁰ W powyższej notacji macierze z kroku pierwszego prezentowanego algorytmu to odpowiednio macierze $UR_{ij}^{(0)g}$ oraz $UCR_{ij}^{(0,1)g}$, gdyż wejściowy (oryginalny) szereg czasowy T_g może być traktowany jako swoja dyskretna zero-rzędowa pochodna.

5. Giorgino T., *Computing and visualizing dynamic time warping alignments in R: The dtw package*, „Journal of Statistical Software” 2009 (7), s. 1–24.
6. Górecki T., Łuczak M., *Using derivatives in time series classification*, „Data Mining and Knowledge Discovery” 2013, nr 2, s. 310–331.
7. Górecki T., Łuczak M., *Multivariate time series classification with parametric derivative dynamic time warping*, „Expert Systems with Applications” 2015, nr 42, s. 2305–2312.
8. Hastie T., Tibshirani R., Friedman J., *The Elements of Statistical Learning. Data Mining, Inference, and Prediction*, New York 2009.
9. Keogh E., Kasetty S., *On the need for time series data mining benchmarks: a survey and empirical demonstration* [w:] *Proceedings of the eight ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, New York 2002, s. 102–111.
10. Keogh E., Pazzani M. J., *Dynamic time warping with higher order features* [w:] *Proceedings of the 2001 SIAM International Conference on Data Mining*, red. V. Kumar, R. Grossman, 2001, s. 1–11.
11. Keogh E., Ratanamahatana C. A., *Exact indexing of dynamic time warping*, „Knowledge and Information Systems” 2005 (7), s. 358–386.
12. Lubbock A. L. R., Harris L. A., Quaranta V., Tyson D. R., Lopez C. F., *Thunor: visualization and analysis of high-throughput dose-response datasets*, „Nucleic Acids Research” 2021, nr 46, s. w633–w640.
13. Marwan N., Romano M. C., Thiel M., Kurths J., *Recurrence plots for the analysis of complex systems*, „Physics Reports” 2007, nr 438, s. 237–329.
14. Mori U., Mendiburu A., Lozano J. A., *Distance measures for time series in R: The TSdist Package*, „The R Journal” 2016, nr 2, s. 451–459.
15. Rabiner L., Juang B.-H., *Fundamentals of Speech Recognition*, New Jersey 1993.
16. R Core Team, *R: A language and environment for statistical computing*. R Foundation for Statistical Computing, Vienna 2022, <https://www.R-project.org/> (on-line 16.02.2024).
17. Sankoff D., Kruskal J., *Time Warps, String Edits, and Macromolecules. The Theory and Practice of Sequence Comparison*, 1999.
18. Todeschini R., Ballabio D., Consonni V., Manganaro A., Mauri A., *Canonical measure of correlation (CMC) and canonical measure of distance (CMD) between sets of data*.

Part 1. Theory and simple chemometric applications, „Analytica Chimica Acta” 2009, nr 648, s. 45–51.

19. Thunor, <https://www.thunor.net> (on-line: 16.02.2024).

20. Venables W. N., Ripley B. D., *Modern Applied Statistics with S*, New York 2002.

21. Wilczek P., *An application of the local binary pattern algorithm and its uniform variant to improve the recurrence and cross-recurrence quantification analyses of the pharmacologically important time series* [w:] *Recent Advances in Computational Oncology and Personalized Medicine. Vol. 2. The Challenges of the Future !*, red. K. Krukiewicz, M. Marczyk, M. Bugdol, S. Bajkacz, Z. Ostrowski, Gliwice 2022, s. 128–152.



Zezwala się na korzystanie z *Dwa nowe algorytmy klasyfikacji farmakokinetycznych szeregów czasowych* na warunkach licencji Creative Commons Uznanie autorstwa 4.0 (znanej również jako CC-BY), dostępnej pod adresem <https://creativecommons.org/licenses/by/4.0/deed.pl> lub innej wersji językowej tej licencji lub którejkolwiek późniejszej wersji tej licencji, opublikowanej przez organizację Creative Commons