

# Klasyfikacja lokalnych rynków pracy z wykorzystaniem przestrzenno-czasowej analizy składowych głównych

Wojciech Łukaszonek<sup>1,2</sup>, Marcin Szymkowiak<sup>2,3</sup>,  
Waldemar Wołyński<sup>4</sup>

<sup>1</sup>Uniwersytet Kaliski

<sup>2</sup>Uniwersytet Ekonomiczny w Poznaniu

<sup>3</sup>Urząd Statystyczny w Poznaniu

<sup>4</sup>Uniwersytet im. Adama Mickiewicza w Poznaniu

02-04.07.2024

- 1 Cel prezentacji
- 2 Opis metody
- 3 Źródło danych i zmienne diagnostyczne
- 4 Wybrane wyniki
- 5 Kluczowe wnioski
- 6 Kierunki dalszych prac badawczych
- 7 Literatura

## Główne cele prezentacji:

- ocena sytuacji na rynku pracy w przekroju powiatów województwa wielkopolskiego,
- omówienie zastosowanego podejścia tj. przestrzenno-czasowej analizy składowych głównych (ang. spatio-temporal principal component analysis - **STPCA**),
- przedstawienie najważniejszych wyników, ich interpretacja oraz sformułowanie wniosków z przeprowadzonej analizy.

# Problem metodologiczny

- Próba składa się z  $n = 35$  obserwacji opisanych  $p = 12$  cechami statystycznymi w  $T = 19$  momentach czasu, zatem każda z obserwacji jest reprezentowana przez macierz o wymiarach  $p \times T = 12 \times 19$ .
- Dokonując wektoryzacji otrzymujemy obiekty będące wektorami o  $pT = 228$  wymiarach.
- **Problem** konstrukcji składowych głównych sprowadza się do estymacji macierzy kowariancji stopnia  $pT = 228$ . Warunkiem poprawności oszacowania jest  $n > pT$ , co w omawianym przypadku nie zachodzi, mamy bowiem:  $35 \ll 12 \cdot 19 = 228$ .
- **Rozwiązanie:** przejście z „przestrzeni pierwotnej” do „przestrzeni danych funkcjonalnych”.
- Operując w konwencji FPCA (Functional Principal Component Analysis) wzbogaciliśmy model o komponent zależności przestrzennych - **STPCA (Spatio-Temporal Principal Component Analysis)**.

# Funkcjonalna analiza składowych głównych (FPCA)

## Dane funkcyjne

Dla  $p$ -wymiarowego procesu stochastycznego:

$$\mathbf{X} = (X_1, X_2, \dots, X_p)^\top$$

zakładamy, że jego  $i$ -ta składowa może być reprezentowana przez skończoną liczbę funkcji bazowych  $\{\varphi_{ij}\}$ :

$$X_i(t) = \sum_{j=1}^{B_i} \alpha_{ij} \varphi_{ij}(t), \quad j = 1, \dots, B_i, \quad i = 1, \dots, p,$$

gdzie  $\alpha_{ij}$  są zmiennymi losowymi.

# Funkcjonalna analiza składowych głównych (FPCA)

## Dane funkcjonalne (cd.)

Używając notacji macierzowej, proces  $\mathbf{X}$  ma postać:

$$\mathbf{X}(t) = \Phi(t) \boldsymbol{\alpha}, E(\boldsymbol{\alpha}) = 0, \text{Var}(\boldsymbol{\alpha}) = \boldsymbol{\Sigma}, \quad (1)$$

gdzie

$$\boldsymbol{\alpha} = (\alpha_{11}, \dots, \alpha_{1B_1}, \dots, \alpha_{p1}, \dots, \alpha_{pB_p})^T,$$

$$\Phi(t) = \text{diag} \left( \phi_1^T(t), \dots, \phi_p^T(t) \right),$$

$$\phi_i^T(t) = (\varphi_{i1}(t), \dots, \varphi_{iB_i}(t)),$$

$$\boldsymbol{\alpha} \in \mathbb{R}^B, B = B_1 + \dots + B_p, i = 1, \dots, p.$$

# Funkcjonalna analiza składowych głównych (FPCA)

## Rozwiązanie

Funkcjonalne składowe główne definiujemy następująco:

$$U = \langle \mathbf{u}, \mathbf{X} \rangle = \int_a^b \mathbf{u}^\top(t) \mathbf{X}(t) dt,$$

gdzie  $\mathbf{X}$  oraz  $\mathbf{u}$  są elementami tej samej przestrzeni  $\mathcal{L}_2^p([a, b])$ .

Przyjmując  $\mathbf{X}$  postaci (1) oraz  $\mathbf{u}$  wyrażone przez:

$$\mathbf{u}(t) = \Phi(t)\gamma, \quad \gamma \in \mathbb{R}^B,$$

otrzymujemy

$$U = \gamma^\top \alpha,$$

$$\text{Var}(U) = \gamma^\top \Sigma \gamma,$$

gdzie  $\Sigma = \text{Var}(\alpha)$ .

# Funkcjonalna analiza składowych głównych (FPCA)

## Rozwiązanie (cd.)

Zakładając, że  $\mathbf{a}$  jest estymatorem (MNK) losowego wektora  $\alpha$ , oraz  $\mathbf{A}_{n \times B} = (\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n)$ , mamy:

$$\mathbf{S} = \frac{1}{n} \mathbf{A}^\top \mathbf{A}, \quad (2)$$

gdzie estymator wariancji FPCA ma postać:

$$\widehat{\text{Var}}(U) = \gamma^\top \mathbf{S} \gamma = \frac{1}{n} \gamma^\top \mathbf{A}^\top \mathbf{A} \gamma.$$

W kolejnym kroku, model uzupełnimy o komponent "przestrzenny".



# Zależności przestrzenne - macierz wag przestrzennych

Rozpatrujemy obserwacje prowadzone w punktach (przestrzeni) oraz relacje pomiędzy parami punktów. Formalnie mamy  $n$  punktów, w których dokonujemy obserwacji  $x_1, \dots, x_n$ . Punkty te wraz z ich zależnościami tworzą sieć opisaną macierzą wag przestrzennych  $\mathbf{W}$ . Wagi  $w_{ij}$ , ( $i, j = 1, 2, \dots, n$ ) wyrażają występowanie zależności pomiędzy punktami  $i$  oraz  $j$ , jak i (opcjonalnie) siłę tej zależności.

W niniejszym opracowaniu rozpatrzono 9 różnych macierzy wag  $\mathbf{W}$ . Najczęściej stosowane wagi  $w_{ij}$  są funkcjami:

- odległości między jednostkami geograficznymi (tu powiatami),
- długości granicy między jednostkami geograficznymi,
- obu wymienionych jednocześnie.

Wyniki STPCA są zależne od wybranej macierzy wag przestrzennych.

## Zależności przestrzenne - zastosowane wagi

No.	Wagi przestrzenne	Parametr
1.	k-Nearest Neighbour Weights (standard form)	$k = 3, 4, 6$
2.	k-Nearest Neighbour Weights (symmetric form)	$k = 3, 4, 6$
3.	Radial Distance Weights	nd.
4.	Power Distance Weights	$\alpha = 1, 2$
5.	Exponential Distance Weights	$\alpha = 1, 2$
6.	Double-Power Distance Weights	$k = 2, 3, 4$
7.	Rook Contiguity Weights	nd.
8.	FPCA	nd.

Wszystkie macierze wag skonstruowano dla 4 wariantów odległości (S-distance, S-time, F-distance i F-time)

distance - odległość geograficzna, time - czas dojazdu,

S - shortest, F - fastest

## Zależności przestrzenne - współczynnik Geary'ego

Konstruując funkcjonalne składowe główne  $U = \gamma^\top \alpha$ , szukamy wektora  $\gamma$ , który maksymalizuje  $\widehat{\text{Var}}(U)$  przy  $\gamma^\top \gamma = 1$ . Dodatkowo, uwzględniamy aspekt przestrzenny, wykorzystując współczynnik Geary'ego  $C$ .

Przyjmijmy, że w  $n$  punktach przestrzeni  $s_1, \dots, s_n$  obserwujemy  $x_1 = x(s_1), \dots, x(s_n)$  oraz niech  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ ,  $y_i = x_i - \bar{x}$ ,  $i = 1, 2, \dots, n$ . Przyjmujemy również, że wagi przestrzenne  $w_{ij}$  spełniają warunki:  $w_{ij} \geq 0$  dla dowolnych  $i \neq j$ ,  $w_{ii} = 0$  dla wszystkich  $i$ .

Niech  $\mathbf{W} = (w_{ij})_{n \times n}$  oznacza macierz wag oraz  $a = \sum_{i=1}^n \sum_{j=1}^n w_{ij}$ . Wówczas współczynnik Geary'ego  $C$  zdefiniowany jest następująco (Geary, 1954):

$$C(y_1, \dots, y_n) = \frac{n-1}{2a} \frac{\sum_{i=1}^n \sum_{j=1}^n w_{ij} (y_i - y_j)^2}{\sum_{i=1}^n y_i^2}.$$

$C$  jest nieujemną miarą autokorelacji przestrzennej.

## Zależności przestrzenne - współczynnik Geary'ego (cd.)

Współczynniki autokorelacji przestrzennej określają, czy i w jakim stopniu obserwacje  $x_i$  oddziałują na siebie wzajemnie poprzez przyjętą strukturę.

Rozważmy  $\mathbf{B} = (B_{ij})$  - macierz stopnia  $n$  o elementach:

$$B_{ij} = \begin{cases} R_i + K_j - 2w_{ij}, & \text{if } i = j, \\ -2w_{ij}, & \text{if } i \neq j, \end{cases}$$

gdzie  $R_i$  jest sumą elementów  $i$ -tego wiersza macierzy  $\mathbf{W}$  oraz  $K_j$  jest sumą elementów  $j$ -tej kolumny tej macierzy.

Współczynnik Geary'ego  $C$  przyjmuje wówczas postać (Jong et al., 1984):

$$C(\mathbf{y}) = \frac{n-1}{2a} \frac{\mathbf{y}^\top \mathbf{B} \mathbf{y}}{\mathbf{y}^\top \mathbf{y}}, \quad (3)$$

gdzie  $\mathbf{y} = (y_1, \dots, y_n)^\top \in \mathbb{R}^n$ .

Niech  $\mathbf{y} = \mathbf{A}\boldsymbol{\gamma}$  będzie rzutem macierzy obserwacji  $\mathbf{A}$  na przestrzeń funkcjonalnych składowych głównych  $U = \boldsymbol{\gamma}^\top \boldsymbol{\alpha}$  określonych przez wektor kierunkowy  $\boldsymbol{\gamma}$ .

Wówczas:

$$C(\mathbf{y}) = C(\mathbf{A}\boldsymbol{\gamma}) = \frac{n-1}{2a} \frac{\boldsymbol{\gamma}^\top \mathbf{A}^\top \mathbf{B} \mathbf{A} \boldsymbol{\gamma}}{\boldsymbol{\gamma}^\top \mathbf{A}^\top \mathbf{A} \boldsymbol{\gamma}}.$$

Chcemy wyznaczyć wektor  $\boldsymbol{\gamma}$  funkcjonalnej składowej głównej  $U = \boldsymbol{\gamma}^\top \boldsymbol{\alpha}$  tak, aby maksymalizował iloczyn wariancji tej składowej oraz współczynnik Geary'ego rzutu macierzy obserwacji na tę składową.

## STPCA (cd.)

Innymi słowy, szukamy wektora  $\gamma$  który spełnia warunek:

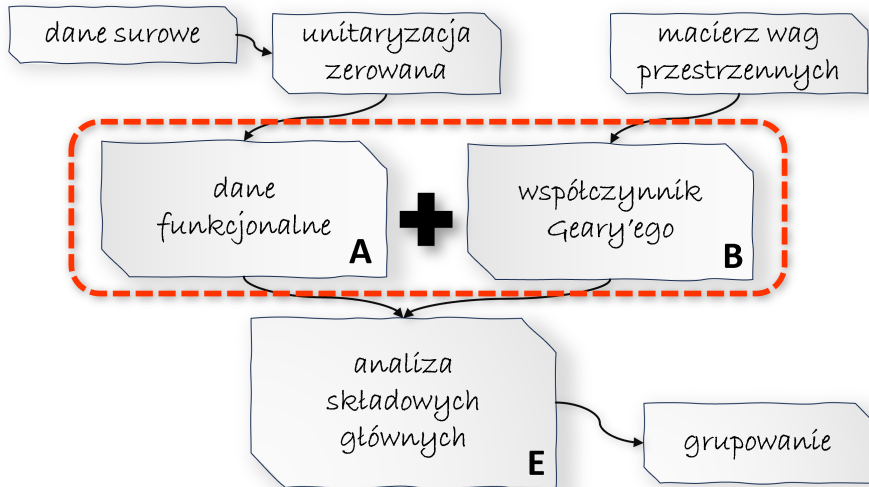
$$\begin{aligned} \arg \max_{\gamma} \widehat{\text{Var}}(U) C(\mathbf{A}\gamma) &= \arg \max_{\gamma} \frac{1}{n} \gamma^{\top} \mathbf{A}^{\top} \mathbf{A} \gamma \frac{n-1}{2a} \frac{\gamma^{\top} \mathbf{A}^{\top} \mathbf{B} \mathbf{A} \gamma}{\gamma^{\top} \mathbf{A}^{\top} \mathbf{A} \gamma} \\ &= \arg \max_{\gamma} \frac{n-1}{2na} \gamma^{\top} \mathbf{A}^{\top} \mathbf{B} \mathbf{A} \gamma = \arg \max_{\gamma} \gamma^{\top} \mathbf{E} \gamma, \end{aligned}$$

przy  $\gamma^{\top} \gamma = 1$ , gdzie

$$\mathbf{E} = \frac{n-1}{2na} \mathbf{A}^{\top} \mathbf{B} \mathbf{A}. \quad (4)$$

Nowy układ współrzędnych (przestrzenno-czasowe składowe główne), wyznaczony jest poprzez wektory własne macierzy  $\mathbf{E}$  odpowiadające jej niezerowym wartościom własnym.

# Schemat badania



# Źródło danych

Wykorzystano zasoby Banku Danych Lokalnych (<https://bd1.stat.gov.pl>) pochodzące z kategorii:

- Rynek pracy (K4),
- Wynagrodzenia i świadczenia społeczne (K40),
- Podmioty gospodarki narodowej, przekształcenia własnościowe i strukturalne (K25),
- Ludność (K3).

Dane zostały poddane wstępnemu przygotowaniu:

- przeliczono wielkości bezwzględne odpowiednio, w relacji do liczby mieszkańców powiatu lub do ogólnej liczby bezrobotnych w powiecie,
- zastosowano unitaryzację zerowaną (z uwzględnieniem stymulant i destymulant).



# Zmienne diagnostyczne

## Badaniem statystycznym objęto:

- wszystkie powiaty województwa wielkopolskiego ( $n = 35$ ),
- dwanaście charakterystyk związanych z rynkiem pracy ( $p = 12$ ),
- okres czasu odnoszący się do lat 2004–2022 ( $T = 19$ ).

ZMN1 –	Odsetek zarejestrowanych bezrobotnych z wyższym wykształceniem (w bezrobotnych ogółem)	ZMN7 –	Przeciętne wynagrodzenie brutto
ZMN2 –	Odsetek zarejestrowanych bezrobotnych pozostających bez pracy dłużej niż 1 rok (w bezrobotnych ogółem)	ZMN8 –	Podmioty wpisane do REGON na 1000 mieszkańców
ZMN3 –	Odsetek nowo zarejestrowanych bezrobotnych (w bezrobotnych ogółem)	ZMN9 –	Podmioty nowo zarejestrowane w REGON na 1000 mieszkańców
ZMN4 –	Odsetek zarejestrowanych bezrobotnych w wieku 24 lata i mniej (w bezrobotnych ogółem)	ZMN10 –	Podmioty wyrejestrowane z REGON na 1000 mieszkańców
ZMN5 –	Odsetek zarejestrowanych bezrobotnych w wieku 55 lat i więcej (w bezrobotnych ogółem)	ZMN11 –	Stopa bezrobocia
ZMN6 –	Odsetek pracujących w ludności w wieku produkcyjnym	ZMN12 –	Oferty pracy na 1000 mieszkańców w wieku produkcyjnym

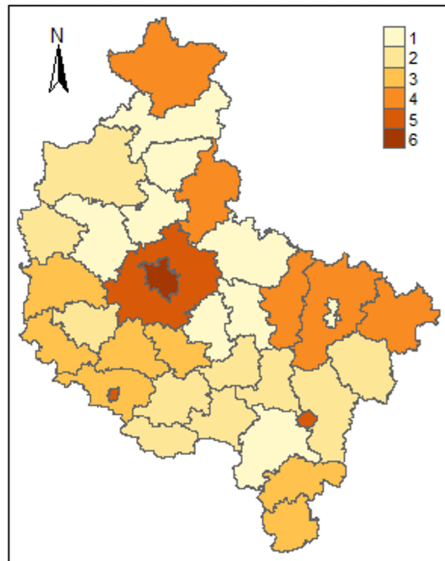
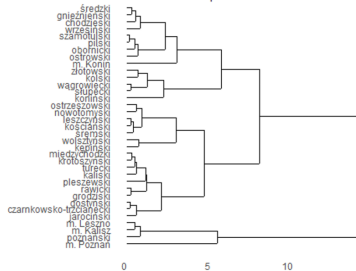
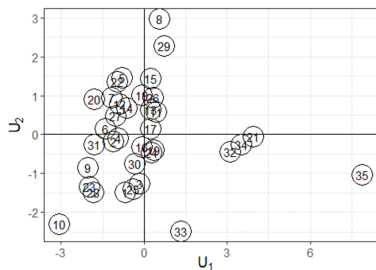
Stymulanta

Destymulanta

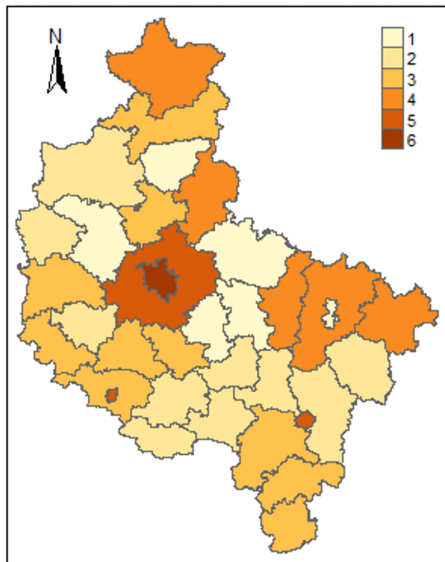
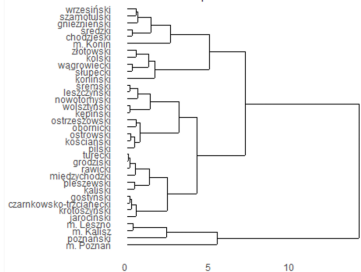
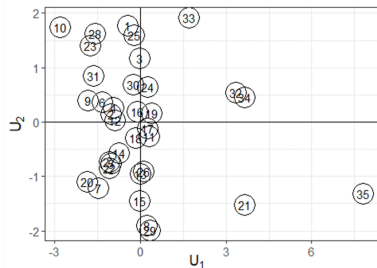
# Wyniki: dla dwóch pierwszych składowych głównych

Model		S-dist	S-time	F-dist	F-time
FPCA		60.13	60.13	60.13	60.13
<i>k</i> -Nearest Neighbour Weights	<i>k</i> = 3	<b>61.39</b>	60.79	60.69	60.32
	<i>k</i> = 4	<b>62.54</b>	58.67	62.33	58.42
	<i>k</i> = 6	61.03	<b>62.42</b>	61.68	61.19
<i>k</i> -Nearest Neighbour Weights (symmetric form)	<i>k</i> = 3	<b>62.11</b>	60.12	60.79	61.69
	<i>k</i> = 4	60.13	59.04	<b>61.87</b>	60.66
	<i>k</i> = 6	61.15	<b>63.18</b>	61.45	61.47
Radial Distance Weights		62.85	<b>63.08</b>	62.70	62.43
Power Distance Weights;	$\alpha$ = 1	61.97	61.98	<b>62.08</b>	61.87
	$\alpha$ = 2	65.10	63.31	<b>65.66</b>	62.98
Exponential Distance Weights	$\alpha$ = 1	62.11	<b>62.57</b>	62.26	62.16
	$\alpha$ = 2	63.25	<b>63.64</b>	63.38	63.34
Double-Power Distance Weights	<i>k</i> = 2	64.34	<b>65.09</b>	64.43	64.82
	<i>k</i> = 3	64.45	<b>65.32</b>	64.57	65.06
	<i>k</i> = 4	64.45	<b>65.41</b>	64.60	65.16
Rook Contiguity Weights		54.69	54.69	54.69	54.69

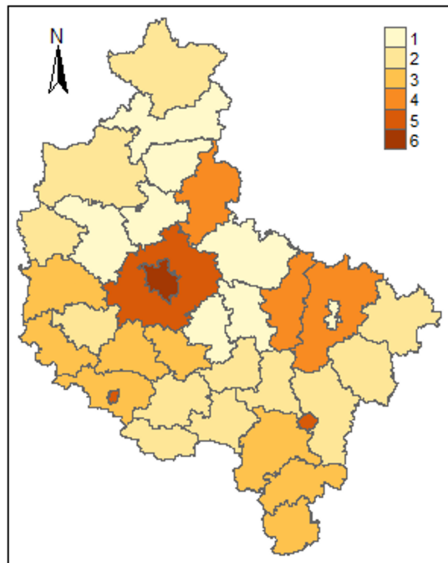
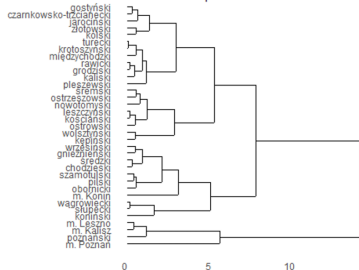
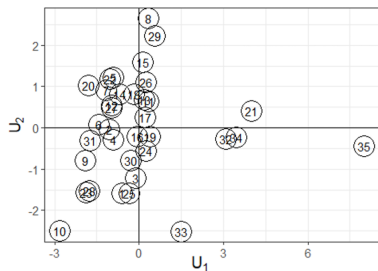
# Wyniki: FPCA (bez komponentu przestrzennego)



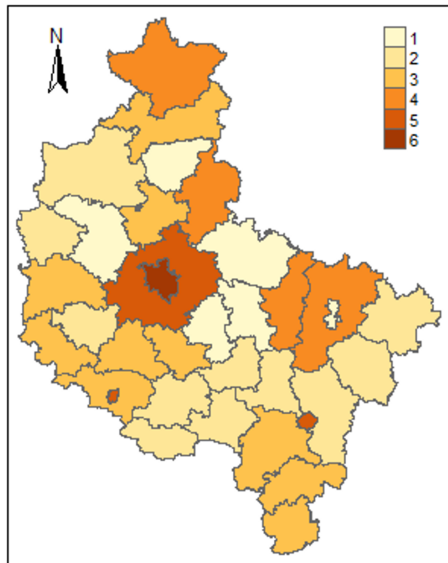
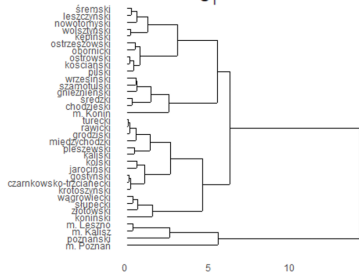
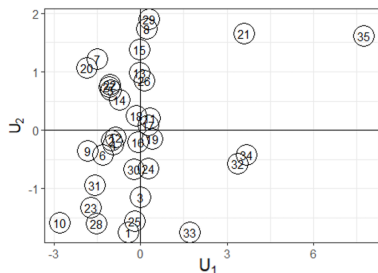
# Wyniki: STPCA S-distance (Power Distance Weights, $\alpha = 2$ )



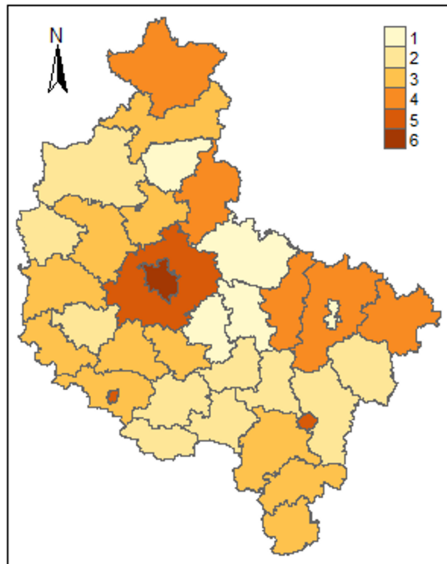
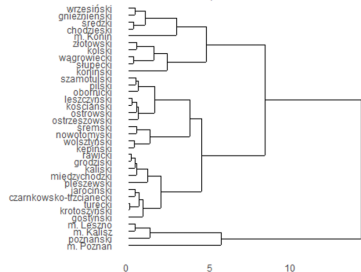
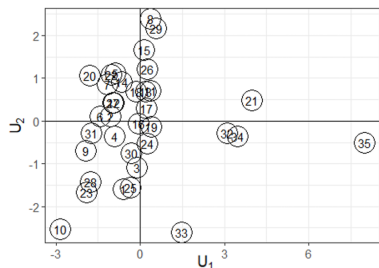
# Wyniki: STPCA S-time (Double-Power Distance Weights, $k = 4$ )



# Wyniki: STPCA F-distance (Power Distance Weights, $\alpha = 2$ )



# Wyniki: STPCA F-time (Power Distance Weights, $\alpha = 2$ )



# Kluczowe wnioski

- STPCA stanowi użyteczną metodę badawczą, w której uwaga koncentruje się na jednoczesnym uwzględnieniu czasowego komponentu cech opisujących jednostki badania wraz z ich przestrzenną lokalizacją.
- Uwzględnienie w badaniu, z wykorzystaniem funkcjonalnej analizy składowych głównych, komponentu przestrzennego w postaci odpowiednio dobranej macierzy wag przestrzennych zwiększa procent wyjaśnianej wariancji.
- Wyniki analiz w kontekście przestrzennego zróżnicowania lokalnych rynków pracy w województwie wielkopolskim są na ogół podobne, choć w niektórych przypadkach (dla pewnych macierzy wag) można jednak zaobserwować nieco odmienne wyniki grupowania powiatów.



## Kierunki dalszych prac badawczych

- badanie dla wszystkich powiatów ( $n = 380$ ),
- badanie na poziomie województw ( $n = 16$ ),
- uwzględnienie w badaniach innych macierzy wag przestrzennych (również innych parametrów  $k, \alpha$ ),
- inne modyfikacje aparatu statystycznego (wprowadzenie w modelach elementów „wygładzania” w wymiarze czasowym).

# Literatura

- Geary, R. (1954), The Contiguity Ratio and Statistical Mapping. The Incorporated Statistician. **5**, 115.
- Górecki, T., Krzyśko, M., Waszak, Ł. & Wołyński, W. (2018), Selected statistical methods of data analysis for multivariate functional data. Statistical Papers. **59**, 153–182.
- Hotelling, H. (1933), Analysis of app. complex of statistical variables into principal components.. Journal Of Educational Psychology. **24**, pp. 417–441.
- Kossowski, T. (2020), Wybrane zagadnienia modelowania matematyczno-statystycznego struktur i procesów przestrzennych. Rozwój Regionalny i Polityka Regionalna., 159–174.
- Krzyśko, M., Nijkamp, P., Ratajczak, W., Wołyński, W., Wojtyła, A. & Wenerska, B. (2023), A novel spatio-temporal principal component analysis based on Geary's contiguity ratio. Computers, Environment And Urban Systems. **103** pp. 101980.
- Ramsay, J., Hooker, G. & Graves, S. (2009), Introduction to Functional Data Analysis. Functional Data Analysis With R And MATLAB. pp. 1–19.

Dziękujemy za uwagę!