

Stanisław Gruszczyński*, Krzysztof Urbański*

Zastosowanie algorytmów interpolacji i sztucznych sieci neuronowych do wyznaczania charakterystyki zawartości chromu w glebach**

1. Wstęp

Potrzeba określania stopnia zanieczyszczenia gleb zachodzi w wielu przypadkach ocen oddziaływania inwestycji na środowisko. Przykładem jest faza tak zwanego screeningu, polegająca na opracowaniu obrazu aktualnego stanu komponentów środowiska. Także przeglądy ekologiczne niekiedy koncentrują się na tej problematyce.

Przed rozpoczęciem prac związanych z oceną stanu zanieczyszczenia gleb np. metalami ciężkimi między innymi wymagane jest podjęcie decyzji dotyczących niżej wymienionych zagadnień.

- 1) Granice opróbowania: należy wskazać linię graniczną, poza którą badania (przynajmniej w pierwszym etapie) nie będą sięgać; jest to oczywisty wymóg wynikający z przesłanek finansowych.
- 2) Techniczne szczegóły pobierania prób: głębokość opróbowania, liczba i miąższość analizowanych warstw, rodzaj prób (złożone czy indywidualne).
- 3) Strategia opróbowania: rozmieszczenie prób losowe, systematyczne czy też nastawione na otoczenie źródeł zagrożeń (celowe) bądź też hybrydowe, czyli łączące cechy wyboru losowego i celowego, lub systematyczne ze zmienną gęstością opróbowania, dostosowaną do lokalnych warunków.
- 4) Przewidywany sposób wizualizacji i interpretacji wyników: mapa ciągła, dyskretna (klasy skażeń), interpretacja statystyczna (parametry statystyczne dla wybranych z góry jednostek: gmin, wsi, sposobów użytkowania, jednostek typologicznych) – wszystko to musi być rozważone łącznie z punktem 3.

* Akademia Górniczo-Hutnicza, Wydział Geodezji Górniczej i Inżynierii Środowiska

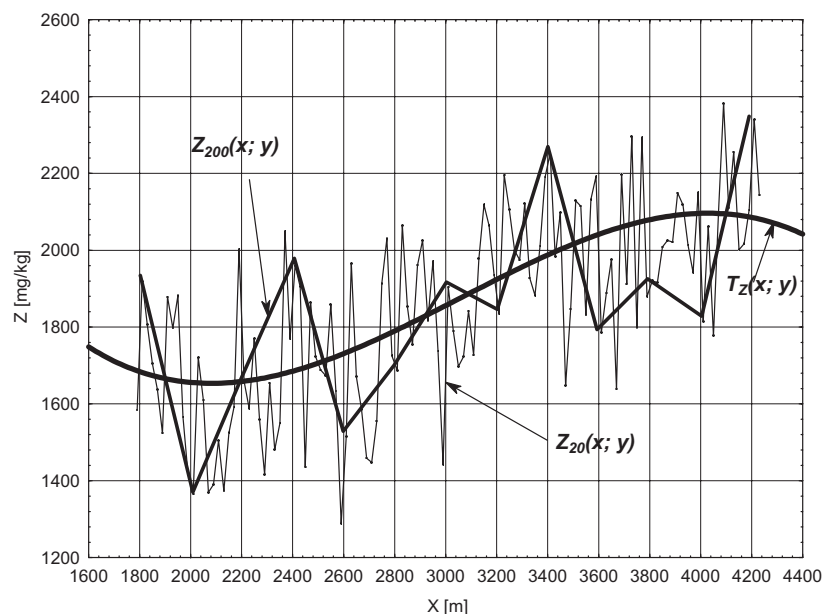
** Praca badawcza finansowana przez Komitet Badań Naukowych w ramach projektu 8 T 12E 007 20 „Technologia prowadzenia badań i kartograficznego opracowania wyników, dotyczących skażenia gleb w rejonach przemysłowych i górniczych”

- 5) Każda z decyzji podjęta w ramach punktu 3. i 4. pociąga konieczność ustalenia dalszych szczegółów: sposobu rozszerzenia obserwacji na powierzchni między punktami opróbowania, interpretacji wyników (rozkład koncentracji, interpretacja statystyczna, wyłonienie głównych składowych generujących zmienność, wykrycie trendów przestrzennych).
- 6) Sposób wizualizacji i waloryzacji konturów (który „lepszy”, który „gorszy”), zwłaszcza przy skażeniu wieloczynnikowym.

Po latach doświadczeń związanych z badaniami zanieczyszczenia gleb metalami ciężkimi w dalszym ciągu występują problemy metodyczne związane z opracowaniem wiarygodnego obrazu ich koncentracji w glebach. Szczególną klasę problemów stwarzają tereny bardzo silnego zanieczyszczenia gleb. Niniejsza praca prezentuje niektóre uwarunkowania związane z interpretacją obserwacji koncentracji metali w glebach stref silnie obciążonych.

2. Sformułowanie problemu

Główną niedogodnością towarzyszącą próbom utworzenia wiarygodnego obrazu stopnia obciążenia gleb metalami ciężkimi jest brak widocznych przesłanek poprawnego rozmieszczenia punktów poboru prób w terenie. Inaczej niż w zadaniach klasycznych pomiarów sytuacyjno-wysokościowych, brakuje tu wiedzy o położeniu „linii szkieletowych” hipotetycznej powierzchni reprezentującej przestrzenny rozkład koncentracji składnika zanieczyszczającego w glebach. Podobna sytuacja występuje także w badaniach geologicznych, gdzie pośrednie wskazówki do lokalizacji wierceń są bardzo słabe. W tych okolicznościach uwaga koncentruje się na odpowiednim zagęszczeniu punktów obserwacyjnych. Problem polega na doborze gęstości punktów obserwacji stężeń, które upoważniałoby do interpolacji wartości między nimi. Występują tu różne uwarunkowania sensownego poziomu zagęszczenia punktów. Elementarny model następstw wyboru gęstości sieci opróbowania gleb schematycznie ilustruje rysunek 1 prezentujący zmienność zanieczyszczenia wzdłuż hipotetycznej linii obserwacyjnej biegnącej w obszarze o dużej zmienności zawartości pierwiastka stanowiącego zanieczyszczenie. Pogrubiona linia ciągła odzwierciedla trend, czyli średni, najbardziej prawdopodobny przebieg poziomu akumulacji, linia łamana oscylująca w jej pobliżu przedstawia prawdziwe stężenie pierwiastka w glebie oznaczone w bardzo gęsto rozlokowanych odkrywkach glebowych. Kluczowa jest relacja między lokalną zmiennością koncentracji pierwiastka a gęstością opróbowania. W wielu możliwych przypadkach mogą pojawić się poważne trudności nawet z dostatecznie dokładnym odzwierciedleniem linii trendu. Z drugiej jednak strony dążenie do maksymalnego zagęszczenia opróbowania, poza znacznym wzrostem kosztów obserwacji, może spowodować, że obraz rozkładu zanieczyszczeń stanie się zupełnie nieczytelny.



Rys. 1. Ilustracja zależności między zmiennością przestrzenną koncentracji zanieczyszczeń w glebach, gęstością opróbowania i dokładnością odwzorowania zagrożeń. Gruba linia ciągła reprezentuje średnią wartość koncentracji (hipotetyczny trend $T_Z(x|y = \text{const})$) obserwowany wzdłuż linii obserwacyjnej. Łamana przedstawia koncentrację zanieczyszczenia $Z_{20}(x|y = \text{const})$ obserwowaną co 20 m, zaś grubsza łamana oznacza koncentrację $Z_{200}(x|y = \text{const})$ obserwowaną co 200 m. Łatwo dostrzec, że w tych warunkach interpolacja bazująca na interwale 200 m prowadzi do znacznego odchylenia od stanu realnego

Tradycyjnie narzędziem analizy zmienności przestrzennej zjawisk jest wariogram, niezbędny do realizacji jednego z najpowszechniej stosowanego algorytmu interpolacji, tak zwanego krigingu. Wariogram daje istotne informacje dotyczące przestrzennej natury zjawiska.

W warunkach silnego zróżnicowania przestrzennego cech glebowych dokładność interpolacji musi być rozpatrywana pod kątem celu, któremu posłuży. Nawet jeżeli założy się, że uzyskany obraz w granicach jakiegoś terenu dobrze odzwierciedla trend, to trzeba pamiętać, że jest to tylko oszacowanie statystyczne (z reguły oszacowanie wartości oczekiwanej, nieuchronnie związane z błędem), zaś lokalne odchylenia od niego mogą być duże. Bardzo przydatne jest uzupełnienie modelu odzwierciedlającego zróżnicowanie cechy o informację dotyczącą stopnia niepewności związanego z jego akceptacją. Niektóre wyniki obserwacji sygnalizują większą, niż się powszechnie zakłada, zmienność koncentracji metali ciężkich w glebach [11]. Oznacza to, że lokalne odchylenia od wartości przyjmowanych za średnie mogą być bardzo duże.

W badaniach zanieczyszczeń gleb stosuje się obecnie różne gęstości opróbowania. W większości przypadków autorzy opracowań skłaniają się do sferycznego modelu wariogramu, z poziomem stabilizacji wariancji od kilkudziesięciu do kilkuset metrów. Można przypuszczać, że zróżnicowanie koncentracji zanieczyszczeń powinno także zależeć od lokalnych czynników mających naturę dyskretną, na przykład od rozmiarów działek warunkujących odmienne warunki akumulacji, pobierania składników oraz innych zjawisk zależnych od użytkownika. Wnioski wyprowadzane zatem z badania obszarów, gdzie dominują tereny dużej własności, nie muszą być w pełni adekwatne dla powierzchni o własności bardzo rozdrobnionej. To samo można powiedzieć o zróżnicowaniu wywołanym bliskim położeniem zakładów emitujących zanieczyszczenia, w szczególności w przypadku niskich emitorów, gdzie ogólny poziom zawartości zanieczyszczeń w glebach, a także zróżnicowanie koncentracji mogą być bardzo wysokie. Poważny problem wynika też z pionowego, silnie zarysowującego się gradientu koncentracji zanieczyszczeń w glebach, zwiększającego znacznie ryzyko błędów opróbowania w wyniku małej dokładności głębokości pobierania próby.

Świadomość istotności wielu czynników kształtujących obraz przestrzennego rozkładu zanieczyszczeń gleb wymaga wielu analiz metodycznych, mających na celu ustalenie poziomu niepewności towarzyszącej ocenie zanieczyszczeń gleb. Uzasadnia to podjęcie szczegółowych badań dotyczących kartograficznego dokumentowania przestrzennej zmienności zanieczyszczeń w rejonach silnie obciążonych emisją metali ciężkich.

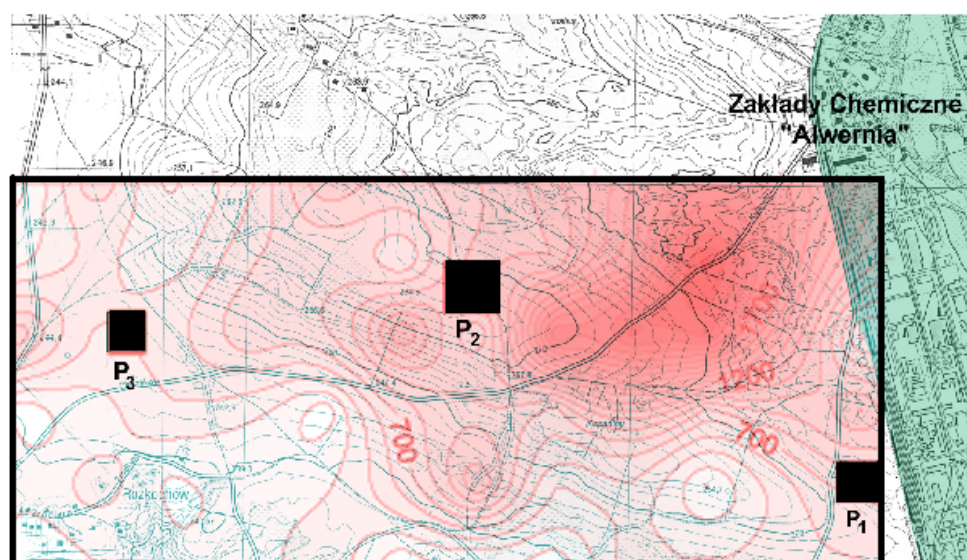
3. Metodyka i rejon badań terenowych

Celem podjętych badań było doświadczalne określenie wpływu zmienności przestrzennej na obraz trendu zanieczyszczenia gleb chromem generowany różnymi technikami, w rejonie o dużym zagrożeniu tym pierwiastkiem [16].

Rejon badań (rys. 2) położony jest na zachód od zakładów chemicznych w Alwerni. Tworzy go prostokąt o rozmiarach 1,0 km w kierunku NS oraz 2,5 km w kierunku WE. Zakłady położone w obrębie granic gminy Alwernia zlokalizowane są w wąskiej ukierunkowanej południkowo dolinie rzeki Regulki. Wzgórza otaczające zakład od wschodu wznoszą się stromo, a różnica poziomów dochodzi do 150 m. Łagodny, falisty teren otwiera się w kierunku zachodnim, natomiast od północy dolinę rzeki zamyka porośnięty lasem Garb Tencyński. W bezpośrednim sąsiedztwie zakładów przeważają lasy (głównie siedliska borowe) z dominującym udziałem sosny w drzewostanie. Średnia suma opadów rocznych wynosi 600±700 mm. Okresem najbardziej nasłonecznionym jest późne lato i wczesna jesień. Średnia temperatura lipca wynosi około 18,5°C. Średnia temperatura stycznia kształtuje się na poziomie -4°C. Średnioroczna temperatura wynosi 7°C. Zaznaczające się cechy

klimatu kontynentalnego to przede wszystkim znaczne amplitudy temperatur. Pokrywa śnieżna zalega przez 80 dni w roku, dominują wiatry zachodnie (około 60% wszystkich wiatrów).

Zakłady chemiczne w Alwerni powstały na przełomie lat 1923/24. Pierwszymi produktami wytwarzanymi przez zakład były kwas solny i dwuchromian sodu. Obecnie Zakłady Chemiczne „Alwernia” SA specjalizują się w trzech profilach produkcyjnych: związkach chromu (zielony tlenek chromu, dwuchromian sodu, dwuchromian potasu, zasadowy siarczan chromu – chromal, bezwodnik kwasu chromowego), związkach fosforu (kwas fosforowy, trójpolifosforan sodu, pirofosforan dwusodu, pirofosforan czterosodowy, fosforan sodu) oraz fosforanach paszowych, siarczanach sodu i magnezu. Rozbudowa zakładu spowodowała wzmoczoną emisję zanieczyszczeń do środowiska, w tym emisję związków chromu, co zdecydowało o umieszczeniu zakładu w latach 90. na liście 80 najbardziej uciążliwych w kraju firm przemysłowych. Pyły z kominów, jak też niezorganizowana emisja pyłów ze skarp osadników błota pochromowego, przyczyniały się do skażenia terenów wokół zakładów. Po przeprowadzonej w ostatnich latach modernizacji procesów technologicznych, zamontowaniu urządzeń ochronnych na emitorach oraz wykonaniu zabiegów rekultywacyjnych na osadnikach Zakłady Chemiczne „Alwernia” SA zmniejszyły ilości emitowanych pyłów o ponad 75% w porównaniu z sytuacją sprzed 6 lat.



Rys. 2. Fragment mapy topograficznej obszaru badań z zaznaczonymi rejonami zagęszczenia punktów opróbowania P_1 , P_2 , P_3 . W części wschodniej zabudowa przemysłowa zakładów chemicznych w Alwerni

Badania obejmowały opróbowanie gleb w węzłach regularnej siatki kwadratów o boku 200 m. Położenie punktów węzłowych siatki ustalono przy użyciu odbiornika GPS z dokładnością do ok. 5 m. Próbkę została pobrana z warstwy gleby o miąższości 0÷30 cm. Zgromadzone w ten sposób dane z opróbowania tworzą zbiór o nazwie *OBSZAR*. W kolejnych cyklach opróbowania zostały wybrane trzy rejony, w których siatka została zagęszczona do wymiarów 20 na 20 m. Wybrane rejony, położone na wschodniej granicy obszaru badań, w jego części centralnej oraz zachodniej (rys. 2), dostarczyły danych do zbiorów oznaczanych odpowiednio: P_1 , P_2 , P_3 . Dane z wszystkich rejonów tworzą zbiór o nazwie *CALOSC*. W pobranych próbach oznaczono metodą ASA całkowitą zawartość chromu, jak również szereg innych właściwości gleb. W niniejszej pracy przeanalizowano jedynie koncentrację chromu w pobranych próbach.

3.1. Wykorzystane oprogramowanie

Omówione w dalszej części pracy wyniki zostały uzyskane przy wykorzystaniu różnego rodzaju oprogramowania statystycznego, symulacyjnego i służącego wizualizacji danych.

Ze względu na wymogi licencji oraz prawa autorskie, wśród wykorzystanych pakietów należy wymienić:

- pakiety statystyczne: Statistica (firma StatSoft), Stagraphics Plus v. 5 (firma Manugistic);
- pakiety GIS: Surfer (firma Golden Software), Idrisi32 (firma ClarkLabs), GRASS (pakiet podlegający licencji GNU GPL);
- pakiety algorytmów adaptacyjnych: Statistica Neural Networks v. 3 i v. 4 (firma StatSoft), NeuroSolutions v. 4.32 (firma NeuroDimensions) oraz SNNS (pakiet podlegający licencji GNU GPL), zbiór makr programu MATLAB, noszący nazwę Netlab, opracowany w Neural Computing Research Group Information Engineering Uniwersytetu Aston w Birmingham, realizujący różne hybrydowe algorytmy sieci neuronowych;
- różne pakiety pomocnicze: MS Excell, MS Access, kompilator Delphi (firma Borland), służące do konwersji i przechowywania danych.

4. Wyniki badań

W problematyce związanej z dokumentowaniem stanu środowiska często istotne są ekstremalne wartości jakiejś cechy. Do tego typu zagadnień należy ocena koncentracji szkodliwych dla ludzi i roślin składników w glebach. Jeżeli zakłada się związek między zawartością składnika w glebach a jego koncentracją w plo-

nach, z możliwością dalszej propagacji do kolejnych receptorów, to obok tendencji centralnej ważne jest też ustalenie obszarów lub prawdopodobieństwa przekroczenia wartości granicznych bądź określenie wartości ekstremalnych stężeń. Możliwe jest tutaj wystąpienie jednej z dwu klas przypadków.

W pierwszej z nich koncentracja szkodliwego składnika wynika głównie z naturalnej jego obecności w komponentach środowiska, na przykład w tworzywie mineralnym gleb lub w wodach podziemnych i powierzchniowych. Wpływ człowieka sprowadza się do drobnych modyfikacji koncentracji wywołanych na przykład odmiennościami użytkowania gleb, nawożenia, uprawy itp. Rozważając problem dokumentowania obciążenia gleb szkodliwymi składnikami w takich warunkach, w granicach terenu o jednorodnej budowie geologicznej, możemy założyć losowy charakter tego zjawiska. Oznacza to, że poszukiwanie deterministycznej zależności między współrzędnymi a koncentracją składnika jest w praktyce bezpodstawne. W takich okolicznościach uzasadniony jest dyskretny obraz obciążenia gleb zanieczyszczeniami, jak również posłużenie się w tym celu odpowiednimi miarami statystycznymi. Łącznie z informacją dotyczącą kształtu rozkładu statystycznego opisanej cechy wystarcza to do określenia potencjalnego zagrożenia.

W drugiej klasie przypadków obok mechanizmu losowego w kształtowaniu się zjawiska występują także składniki deterministyczne, związane z ukierunkowaniem akumulacji w glebach, położeniem źródeł emisji, warunkami propagacji i klimatem. Składnik deterministyczny, wiążący oczekiwaną wartość koncentracji zanieczyszczeń ze współrzędnymi, może być wykrywany, opisywany, jak również uwidaczniany za pomocą metod służących tym celom, zaś obraz ten może być uzupełniony o składową losową.

Podsumowując, w praktyce mogą być brane pod uwagę możliwości postępowania wymienione poniżej.

- Potraktowanie zjawiska jako w pełni losowego. Nie można przekreślać całkowicie takiego podejścia nawet w warunkach występowania trendu przestrzennego, zwłaszcza w peryferyjnych partiach obszaru oddziaływania przedsięwzięcia. Oznacza to konieczność ustalenia parametrów rozkładu statystycznego badanej cechy, z konsekwencjami dotyczącymi wyznaczenia stanu zagrożenia jednolitego dla większego obszaru.
- Interpolacja wartości cechy w oparciu o jej rejestrację w punktach obserwacji. Oznacza to potraktowanie obserwowanej cechy analogicznie jak morfologii terenu i skonstruowanie „powierzchni” odzwierciedlającej poziom koncentracji obserwowanego składnika. Dostępny jest bardzo obszerny zbiór metod interpolacji różniących się założeniami dotyczącymi sposobów wagi obserwacji w celu ustalenia przebiegu izolinii. Wśród nich dostępne są także metody zapewniające zerowy poziom odchyłań interpolowanej

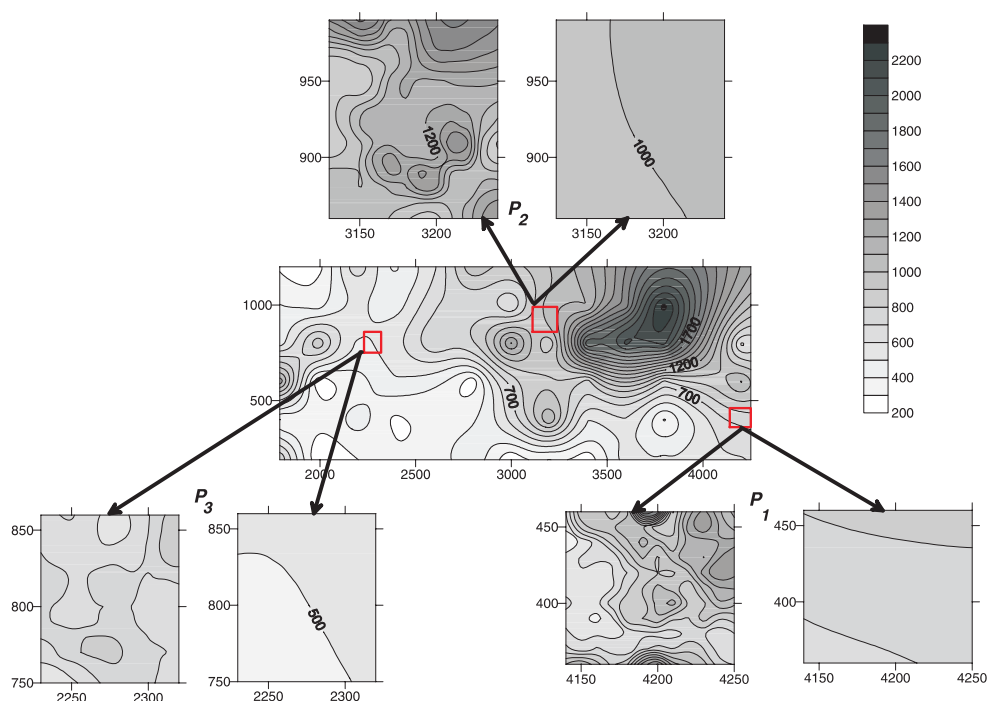
„powierzchni” od „rzędnych” w punktach obserwacji. Taki model postępowania może być dość zawodny, gdy występują bardzo duże różnice stężeń składnika na bardzo małych odległościach, przy ogólnie znacznie rzadszej sieci obserwacyjnej. Warto też zwrócić uwagę, że przebieg tak wyznaczonej „powierzchni” może się bardzo różnić nawet od „powierzchni” obrazującej trend zjawiska. Odwzorowanie takie odbiegałoby zatem i od poprawnego obrazu statystycznego, i od deterministycznego zróżnicowania koncentracji składnika.

- Zaakceptowanie dualizmu zjawiska poprzez odizolowanie części deterministycznej i losowej. Problemem jest w tym wypadku dobór reprezentacji funkcji odzwierciedlającej część deterministyczną (trend zjawiska). Można tu rozpatrywać stosowanie funkcji elementarnych lub wielomianów niskiego stopnia jako funkcji aproksymujących trend, jakkolwiek trudno z góry przewidzieć, jaki model będzie odpowiedni do tego celu. Można też konstruować modele udostępnione przez algorytmy adaptacyjne: sieci neuronowe i rozmyto-neuronowe. Ich przydatność jako uniwersalnych aproksymatorów jest znana od wielu lat [6, 19]. Sieci neuronowe z jedną warstwą ukrytą mogą aproksymować każdą funkcję ciągłą. Nie stanowi przeszkody w ich zastosowaniu do modelowania trendu przestrzennego brak wygodniejszej, analitycznej formy, ponieważ celem ich wykorzystania w tym przypadku jest wyłącznie obliczenie wartości trendu w dowolnym punkcie obszaru. Zastosowanie sztucznych sieci neuronowych (podobnie jak równań regresji) jako uniwersalnego aproksymatora do omawianego celu wymaga istnienia deterministycznego związku między położeniem punktu a koncentracją szkodliwego składnika; należy także przestrzegać restrykcji dotyczących topologii sieci, jej rozmiaru oraz liczebności danych uczących. Przy ustaleniu jako kryterium optymalizacji sieci minimalizacji błędu średniokwadratowego wyniki ich obliczeń można traktować analogicznie jak modele statystyczne. Różnice między obliczonym modelem a wartościami obserwowanymi dają podstawę do oszacowania składnika losowego związanego ze zjawiskiem akumulacji szkodliwych składników w glebach.

Głównym zagrożeniem środowiska w otoczeniu zakładów chemicznych w Alwerni (ZCh Alwernia) jest nadmierna koncentracja związków chromu w glebach. Nie ma też wątpliwości co do źródła tego zagrożenia i jego położenia. Stopień zanieczyszczenia gleb potwierdzają jednoznacznie statystyczne parametry poszczególnych wyróżnionych prób (tab. 1): P_1 , P_2 , P_3 , *OBSZAR* oraz sumaryczne dla całego zbioru danych. Przestrzenną zmienność koncentracji obrazują izoliny stężenia Cr w glebach, wykonane z użyciem programu Surfer (rys. 3).

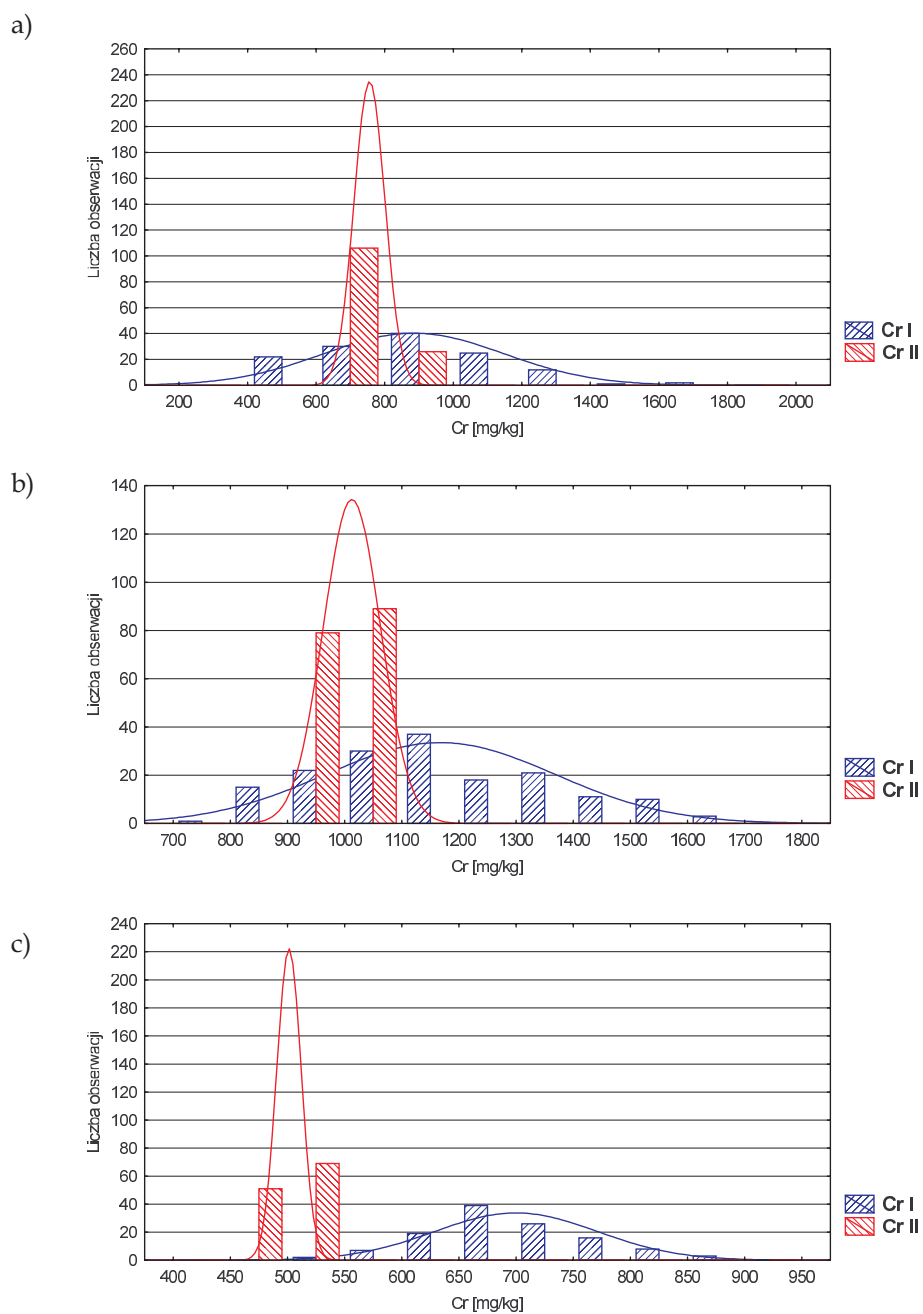
Tabela 1. Statystyczna charakterystyka koncentracji Cr [mg/kg gleby] w warstwie 0÷30 cm, w rejonie badań w okolicach Alwerni

Rejon	Liczba prób	Średnia	Odch. stand.	Minimum	Maksimum
P_1	122	885	332	380	2850
P_2	42	1154	273	680	1740
P_3	30	699	100	520	900
OBSZAR	63	699	477	212	2335
CAŁOŚĆ	257	861	380	212	2850



Rys. 3. Izolinie koncentracji Cr [mg/kg gleby] w glebach w rejonie badań. Zaznaczono położenie rejonów P_1 , P_2 i P_3 . Powiększono fragmenty w celu porównania rozkładu koncentracji ustalonego na podstawie zagęszczonej sieci punktów (rysunki lewostronne) oraz interpolowanego wyłącznie na podstawie obserwacji w siatce 200 m × 200 m (rysunki prawostronne)

Rysunek warstwiczny w rejonach, gdzie gęstość opróbowania została zwiększona, wskazuje na znaczne zróżnicowanie koncentracji Cr nawet przy niewielkich odległościach między obserwowanymi miejscami. Zwracają uwagę różnice obrazu skażeń wygenerowanego na podstawie danych silniej rozproszonych oraz zagęszczonych; potwierdzeniem są histogramy rozkładu koncentracji chromu dla tych samych rejonów (rys. 4); występują tu istotne statystycznie różnice między obu ocenami.

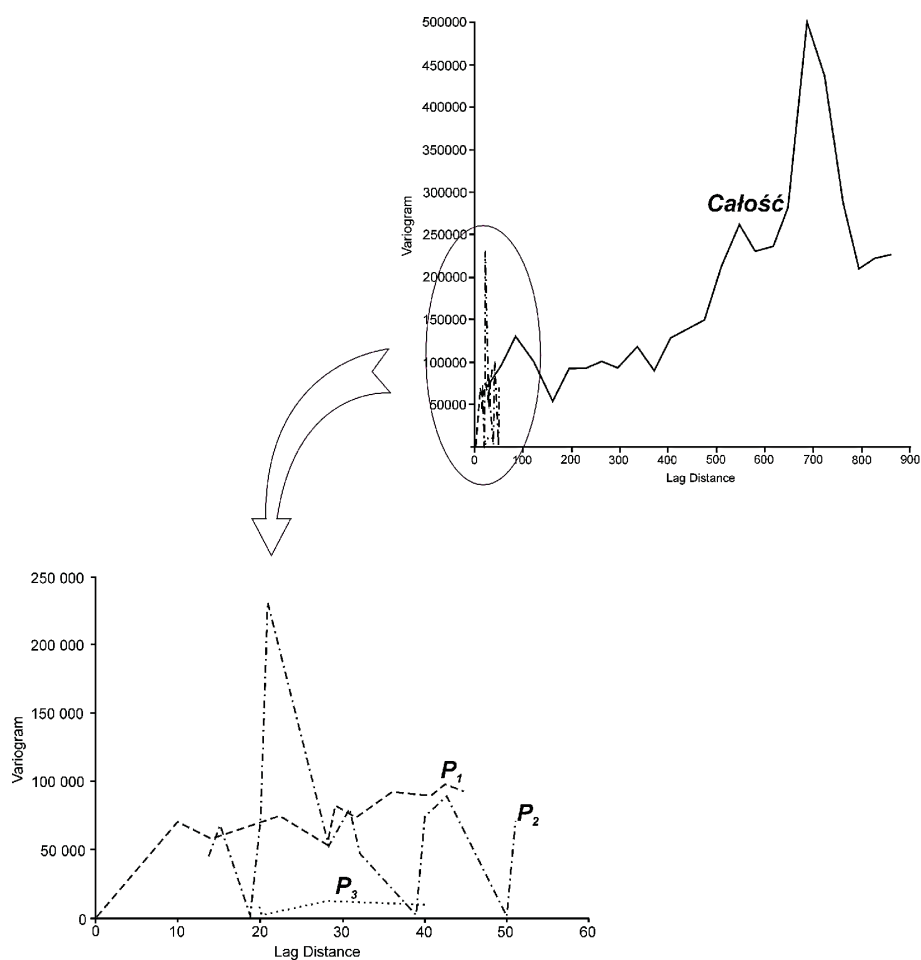


Rys. 4. Histogramy oszacowań zawartości Cr w glebach dla rejonów P_1 (a), P_2 (b) i P_3 (c) sporządzone na podstawie danych z opróbowania w siatce zagęszczonej (Cr I) i na podstawie interpolacji w siatce 200×200 m (Cr II)

W szczególności na histogramach widoczne są znaczne dysproporcje w rozproszeniu wyników oszacowań, w tym wyraźne zawężenie przedziałów zmienności koncentracji chromu w przypadku obrazu konstruowanego na podstawie sieci rozrzedzonej. Anomalie pojawiające się w formie nieregularnego rysunku warstwicowego w rejonach zagęszczenia opróbowania wskazują, że zastosowanie skrajnie gęstej siatki dostarcza mało czytelnego obrazu zanieczyszczenia gleb. Dokładniejsza analiza poziomów rozproszenia losowego poszczególnych zbiorów danych (P_1 , P_2 , P_3 , OBSZAR i sumy tych zbiorów) wskazuje, że w obrębie stosunkowo niewielkich fragmentów terenu wariancja jest tylko nieznacznie mniejsza od wariancji liczonej dla całego obszaru objętego badaniami. Równocześnie jednak na obrazie sporządzonym wyłącznie na podstawie danych z siatki podstawowej generalizacja zacierza istotne informacje. Wynika to z faktu, że w rozkładzie przestrzennym koncentracji Cr na tym terenie występują silne oscylacje nawet w obrębie stosunkowo małych powierzchni. Ważnym wnioskiem wynikającym z analizy obrazu koncentracji zanieczyszczeń są: symptomy występowania trendu przestrzennego oraz mozaikowatość zmienności koncentracji zanieczyszczeń w rejonach o zagęszczonym opróbowaniu. Trend nie ma raczej charakteru liniowego, ponieważ najwyższa koncentracja Cr występuje w pewnym oddaleniu od granic obszaru badań. Pewnym problemem pozostaje w dalszym ciągu pogodzenie wymogów czytelności obrazu zanieczyszczeń z uwzględnieniem silnego zróżnicowania przestrzennego, nie uwidocznionego na zgeneralizowanym odwzorowaniu.

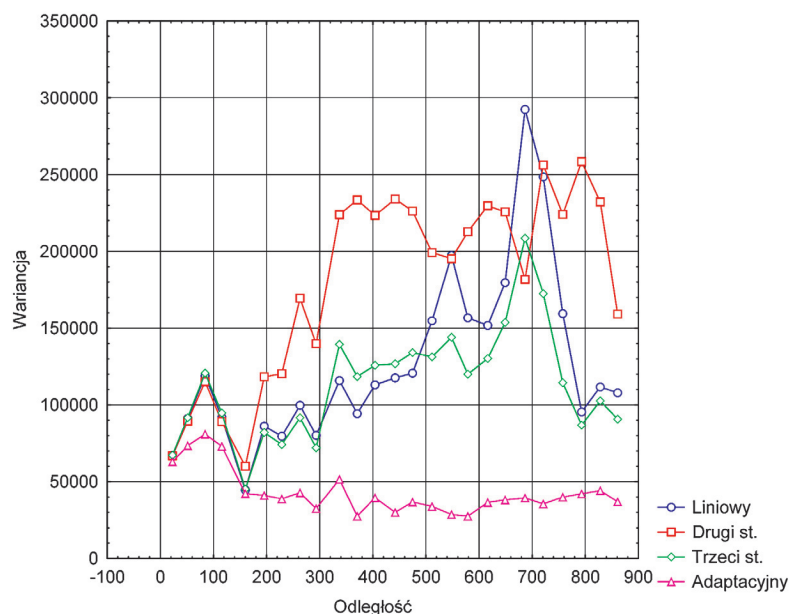
Informacji o przestrzennych zależnościach zachodzących między koncentracją Cr w poszczególnych punktach dostarczają wariogramy sporządzone dla całości danych oraz dla wyodrębnionych pól (rys. 5). Widoczne są znaczne różnice między wariogramami uzyskanymi dla różnych fragmentów obszaru badań. Bardzo istotnie różnią się na przykład wariogramy dla obszarów z zagęszczonym opróbowaniem. Sygnalizują one bardzo dużą wariancję lokalną, nawet przy bardzo bliskich odległościach między punktami. Wariogram uzyskany na podstawie sumy wszystkich danych ma charakter liniowy. Jego liniowość jest pośrednim potwierdzeniem obecności trendu przestrzennego, wyrażającego się między innymi brakiem wyraźnego poziomu stabilizacji wariancji stężeń Cr, nawet przy dużych odległościach między miejscami poboru prób. Anizotropowość, w połączeniu z wyraźnymi różnicami nachylenia wariogramów dla fragmentów gruntu w różnych odległościach od źródła, wskazuje, że nie jest w pełni uprawniony jeden model zmienności wariancji przyjęty dla całego obszaru. Przeciwnie, obecność wyraźnego trendu w przestrzennym rozkładzie zanieczyszczenia oraz dysproporcje wariancji koncentracji w różnych odległościach od domniemanego źródła emisji uzasadniają potrzebę jego wyizolowania. Z dostępnych danych można przy tym wnioskować, że oprócz wyraźnego trendu przestrzennego, będącego głównym źródłem zmienności, rozproszenie koncentracji Cr jest zwiększane przez znaczne wahania stężeń

pierwiastka na stosunkowo niewielkich odległościach. Jest to szczególnie istotne, gdy w tworzeniu przestrzennego obrazu zanieczyszczeń korzystamy z modelu krigingu, który zwykle zdaje egzamin w przypadku występowania trendu przestrzennego, lecz może dawać gorsze rezultaty w wyniku skokowych różnic między wariogramami z różnych fragmentów obszaru.



Rys. 5. Wariogramy empiryczne dla danych z rejonu Alwerni. Powiększono fragment wykresu zawierający wariogramy dla rejonów P_1 , P_2 i P_3

Nie ma jednoznacznych przesłanek do wybrania konkretnego modelu trendu, służącego do wyodrębnienia stałego składnika zmienności przestrzennej. We wstępnych badaniach zastosowano model liniowy, wielomiany drugiego i trzeciego stopnia, a także model adaptacyjny wybrany na podstawie procedury iteracyjnej z sieci wielowarstwowych realizujących zadania regresyjne.



Rys. 6. Wariogramy reszt odchyleń od powierzchni modeli trendu koncentracji Cr w glebach

Wariogramy sporządzone dla reszt pozostałych po usunięciu poszczególnych rodzajów trendu (rys. 6) wskazują na słabości rozwiązania analitycznego. Sygnalizują one w dalszym ciągu obecność resztowego trendu. Najlepszy skutek dało zastosowanie modelu adaptacyjnego, który na tle pozostałych charakteryzuje się stosunkowo małą wariacją resztową oraz właściwie poziomym przebiegiem wariogramu (niezależność wariacji od odległości między obserwowanymi punktami powyżej pewnej odległości). Anomalię, polegającą na ujawnianiu się wyższych wariacji przy małych odległościach punktów opróbowania, tłumaczy mała elastyczność modelu neuronowego, wynikająca z konieczności ograniczenia rozmiaru topologii sieci w celu uniknięcia zjawiska nadmiernego dopasowania (*overfitting*) zagrażającego przy małej liczności zbioru uczącego i zbyt dużej liczbie swobodnych parametrów ustalanych w procesie treningu.

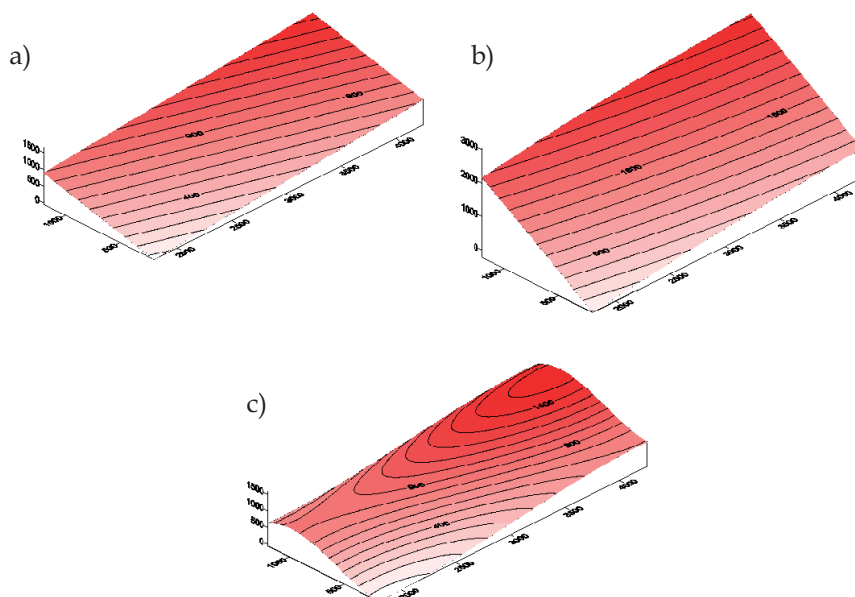
4.1. Porównanie skutków działania algorytmów

W warunkach quasi-punktowego źródła emisji zagadnienie ustalenia przestrzennego trendu zanieczyszczenia gleb jest zadaniem z zakresu regresji. Poszukiwana jest funkcja $f(\mathbf{x})$ aproksymująca nieznaną funkcję $\varphi(\mathbf{x})$ z zadowalającą praktycznie dokładnością. Jednakże zadanie to może być realizowane na wiele różnych sposobów, udostępnionych zarówno przez tradycyjne metody modelowania statystycznego, jak i współczesne techniki eksploracji danych.

Algorytm regresji wielomianowej

Do tradycyjnych technik wykrywania trendu przestrzennego z pewnością należy zaliczyć statystyczną analizę regresji. W hipotetycznym modelu zależności $t(\mathbf{x}) = \varphi(\mathbf{x})$, \mathbf{x} oznacza wektor współrzędnych decydujący o lokalnym rozkładzie badanej cechy. Poszukiwana jest najbardziej prawdopodobna wartość $t(\mathbf{x})$ utożsamiana zwykle ze średnią. W tradycyjnym podejściu zakłada się, że obserwacje zjawiska są obciążone szumami losowymi (zakłóceniami), w związku z tym regresja $f(\mathbf{x})$ aproksymująca nieznaną funkcję $\varphi(\mathbf{x})$ reprezentuje lokalną wartość zjawiska obciążoną wariacją $v(\mathbf{x})$, jednakową w całym zakresie zmienności \mathbf{x} . Oczwistym problemem jest trudność w zaproponowaniu reprezentatywnego dla zjawiska kształtu funkcji $f(\mathbf{x})$. Jak wiadomo, nie istnieje teoria dotycząca rozkładu imisji zanieczyszczeń w glebach, zatem musi być ona poszukiwana w drodze prób i błędów. W badaniach trendów przestrzennych zjawiska tego rodzaju nie ma podstaw do wyboru żadnej szczególnej funkcji. W większości wypadków [9] stosowane są wielomiany różnych stopni.

Uzyskane w toku badań terenowych i laboratoryjnych dane wskazują na istnienie trendu przestrzennego zawartości Cr w glebach, polegającego – ogólnie rzecz biorąc – na spadku koncentracji Cr wraz z odległością od źródła emisji (rys. 3). Rysunek 7 przedstawia pseudo-trójwymiarowy obraz trendu obliczonego jako wielomiany pierwszego, drugiego i trzeciego stopnia.



Rys. 7. Pseudo-trójwymiarowa wizualizacja trendu rozkładu Cr w glebach, obliczonego jako wielomian pierwszego (a), drugiego (b) i trzeciego (c) stopnia

Tabela 2. Niektóre statystyki rozkładu wartości resztowych oszacowania trendu zawartości Cr w glebach; wielomiany pierwszego, drugiego i trzeciego stopnia. Zawartość chromu wyrażona w mg/kg gleby

Statystyka	I stopnia	II stopnia	III stopnia
Minimum	-683,2	-890,7	-865,9
Maksimum	1896,6	1955,2	1886,8
Średnia	9,9	59,8	-33,3
Mediana	-29,2	9,5	-87,5
Pierwszy kwartyl	-210,1	-139,3	-230,3
Trzeci kwartyl	147,3	223,6	129,8
Błąd standardowy	19,7	19,30	19,4
Odchylenie standardowe	313,1	298,8	307,6

Istotnych informacji na temat stopnia niepewności związanego z aprobatą tych modeli dostarcza analiza ich wartości resztowych (tab. 2). Przede wszystkim wynika z niej, że reszty charakteryzują się asymetrią rozkładu, zarówno w odniesieniu do wartości średniej, jak i mediany; asymetria ta jest dodatnia, co oznacza, że w stosunku do hipotetycznej, modelowej wartości średniej należy oczekiwać raczej dużych przekroczeń. Informacje te potwierdzają niezwykle dużą przestrzenną zmienność koncentracji chromu, aczkolwiek odległość międzykwartylowa (obejmująca połowę przypadków) wynosi około 350 mg/kg gleby.

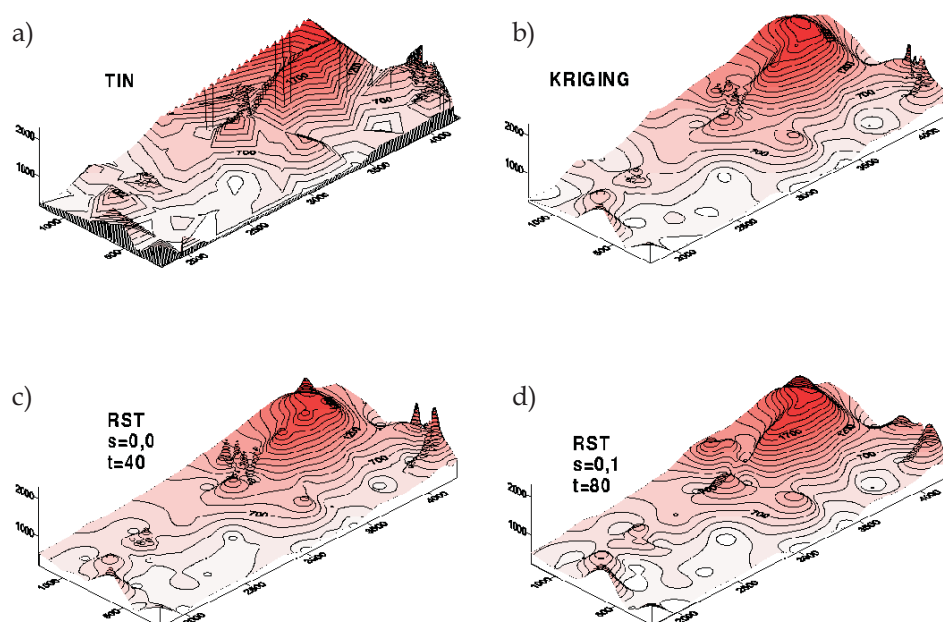
Algorytmy interpolacyjne

Algorytmy interpolacyjne stanowią szczególny przypadek aproksymacji, zaś ich jakość, obok innych cech, jest oceniana pod kątem odległości punktów eksperymentalnych od powierzchni aproksymującej. Wynika to wprost z potrzeb geodezyjnych – głównego pola zastosowań interpolacji w naukach przyrodniczych, w których potencjalny błąd obserwacji jest mały, zaś rozmieszczenie punktów rejestracji rzędnych jest zgodne z zasadami sztuki nakazującymi odzwierciedlanie przebiegu linii szkieletowych. Postulowana precyzja algorytmów interpolacyjnych, powszechnie udostępnianych w pakietach systemów informacji terenowej (SIT), uzasadnia ich wykorzystanie jako bazowej metody porównawczej w stosunku do konkurencyjnych algorytmów aproksymujących.

Rysunek 8 ilustruje pseudo-trójwymiarowe obrazy trendu uzyskane przy użyciu różnych algorytmów interpolacyjnych dostępnych w pakietach SIT.

Algorytm TIN (*Triangulation Irregular Network*), jeden z podstawowych algorytmów interpolacyjnych, buduje powierzchnię aproksymowaną przez sieć nieregularnych trójkątów, z węzłami położonymi w punktach z zaobserwowanymi wartościami rzędnych. Oznacza to, że lokalny kształt siatki jest warunkowany tylko najbliższym sąsiedztwem. Algorytm zapewnia wysoką zgodność przebiegu

powierzchni z obserwacjami i bardzo dobrze zdaje egzamin w małych skalach. Przy dużych skalach widoczny staje się sztuczny, nie przystający do zmienności obserwowanej w realnym świecie rysunek warstwiczny tworzony przez linie łamane. Algorytm charakteryzuje się brakiem możliwości ekstrapolacji i przy jego stosowaniu należy pamiętać o konieczności rozmieszczenia punktów obserwacji także poza obszarem zainteresowania.



Rys. 8. Pseudo-trójwymiarowa wizualizacja trendu rozkładu Cr w glebach wyznaczonego przy użyciu algorytmów interpolacyjnych: *Triangular Irregular Network* (TIN) (a), kriging (b) i *Regularized Spline with Tension* (RST) z dwoma wersjami parametrów (c i d)

Kriging jest powszechnie znanym algorytmem stosowanym w zadaniach kartograficznego dokumentowania zjawisk z nieznanym przebiegiem „linii szkieletowych”, na przykład w zastosowaniach geologicznych, z definicji bazujących na regularnej sieci wierceń. Stosunkowo chętnie jest on także używany w zadaniach z zakresu skażeń gleb. Centralnym narzędziem krigingu jest wariogram służący do ustalania sposobu wagowania wpływu otoczenia na rzędne węzłów regularnej siatki, najczęściej sieci kwadratów lub prostokątów. W przeciwieństwie do TIN nie ma tu przeszkód do ekstrapolacji, aczkolwiek miejscami generuje to poważne błędy aproksymacji powierzchni. Liczne zalety algorytmu ujawniają się w warunkach pewnej jednorodności rzeźby, to jest w przypadkach, kiedy wariogram algorytmu jest reprezentatywny dla całej powierzchni. Może on zawodzić w przypadkach znacznych odstępstw od regularności rzeźby, to znaczy wówczas, gdy charakter powierzchni wskazywałby na potrzebę zróżnicowania wariogramów.

Najbardziej złożony i czasochłonny pod względem obliczeniowym jest algorytm RST (*Regularized Spline with Tension* nazywany też *Completely Regularized Spline*). Wykorzystuje on lokalne funkcje sklejące [10], zaś w implementacji zrealizowanej w pakiecie GRASS (*s.surf.rst*) wyróżnia się wyjątkowo obszerną listą parametrów sterujących jego zachowaniem. W zależności od zadeklarowanych parametrów, algorytm z większą lub mniejszą tolerancją podchodzi do problemu przejścia powierzchni przez punkty danych. Autorzy algorytmu zachęcają do eksperymentowania w celu otrzymania zadowalających rezultatów. Przedstawione na rysunku 8 zobrazowania interpolacji metodą RST różnią się dwoma parametrami. Parametr wygładzania w pierwszym przypadku ustalono na $s = 0,0$ (brak tolerancji w prowadzeniu powierzchni przez punkty obserwacyjne), zaś w drugim $s = 0,1$ (słaba tolerancja). Odpowiednio parametr napięcia (ang. *tension*), decydujący o elastyczności i ciągłości powierzchni, w pierwszym przypadku został przyjęty z wartością domyślną $t = 40$, zaś w drugim $t = 80$ (zmniejszona sztywność). Rezultaty zmiany obu parametrów są widoczne jako utrata szczegółowości i zanik ostrych deniwelacji powierzchni w drugim przypadku.

Ze względu na istotne różnice między algorytmami na wykresach (rys. 8) wyraźnie widoczne są ich efekty. Wrażenie najbardziej nienaturalnego przebiegu sprawia wynik interpolacji metodą TIN. Widoczne są charakterystyczne dla tego algorytmu ostre krawędzie i łamany przebieg izolinii. Pozostałe wizualizacje nie różnią się znacznie od siebie, aczkolwiek z pewnością obie wersje interpolacji wygenerowanej przez RST robią wrażenie ściślej dopasowanych do danych od wariantu powstałego przy użyciu krigingu. W celu porównania wyników tego podejścia do problemu trendu przestrzennego należy przeanalizować rozkłady reszt będących różnicami obserwowanych wartości i rzędnych interpolowanej powierzchni. Zestawienie głównych statystycznych parametrów rozkładu różnic zawiera tabela 3.

Tabela 3. Niektóre statystyki rozkładu wartości resztowych oszacowania trendu zawartości Cr w glebach przy użyciu algorytmów interpolacyjnych stosowanych w pakietach GIS: TIN, kriging oraz RST, z dwoma zestawami parametrów: RST(d) – mniejsza tolerancja wobec zgodności powierzchni z punktami obserwacji, RST(e) – większa tolerancja. Zawartość Cr w mg/kg gleby

Statystyka	TIN	Kriging	RST(d)	RST(e)
Minimum	-570,5	-497,0	-465,5	-455,6
Maksimum	1116,9	1135,6	1380,7	1665,6
Średnia	2,2	0,4	-2,0	-1,0
Mediana	-1,7	-2,2	-4,9	-6,4
Pierwszy kwartyl	-48,3	-44,7	-62,6	-101,5
Trzeci kwartyl	45,8	32,1	24,2	60,6
Błąd standardowy	9,6	8,5	10,8	13,4
Odchylenie standardowe	144,5	135,1	166,2	205,9

Porównanie wyników analizy reszt po zastosowaniu algorytmów interpolacyjnych z wynikami klasycznej regresji wielomianowej nie pozostawia wątpliwości co do formalnej przewagi algorytmów SIT. Rozmiar reszt jest tutaj znacznie mniejszy, odległość między pierwszym i trzecim kwartyłem, niezależnie od metody interpolacji, nie przekracza 165 mg/kg, zaś manipulacja parametrami przetwarzania może doprowadzić do odległości nie przekraczającej 80 mg/kg. Należy tu jednak podkreślić, że silne strony algorytmów interpolacyjnych są jednak, w tym konkretnym przypadku, także ich wadami. Pomijając nawet, w dalszym ciągu dość znaczne, maksymalne odchylenia obserwacji od wartości modelowych (ponad 1000 mg/kg), za istotny problem należy uznać także fakt, że uzyskana pozorna dokładność bazuje na danych obserwacyjnych obciążonych szumami losowymi. Podejście to jest o tyle niebezpieczne, że przestrzenna zmienność zjawiska jest dużo większa, niż można byłoby przypuszczać, zaś potencjalnych źródeł zakłóceń obserwacji jest bardzo dużo. Oznacza to, że przeprowadzenie interpolowanej powierzchni ściśle przez punkty obserwacyjne, niezależnie od sposobu ustalania wartości zjawiska między tymi punktami, jest obciążone nieznanym błędem nieadekwatności położenia punktów rejestracji w stosunku do potrzeb odzwierciedlenia rzeczywistej zmienności koncentracji C_r na powierzchni terenu.

Sztuczne sieci neuronowe

Charakterystyczną cechą obserwowanego zjawiska jest dość znaczna nieliniowość, trudna do analitycznego wyrażenia. Charakteryzuje ona obserwacje zarówno na większych dystansach, jak i na mniejszych odległościach, aczkolwiek w jakimś stopniu odstępstwa od regularności stanowią szum losowy, generujący pewien stopień niepewności obrazu. Oczywiście możliwe jest żmudne poszukiwanie akceptowanego modelu analitycznego, wprowadzanie do niego coraz wyższych stopni wielomianu lub kolejnych funkcji aproksymujących, jednakże takie działanie nie jest oparte na umocowanej teoretycznie analizie mechanizmu zjawiska. Uzasadnia to, w niektórych sytuacjach, rezygnację z formalnej elegancji modelu na rzecz praktycznej skuteczności współcześnie dostępnych metod inteligencji obliczeniowej i tak zwanego *soft computingu*, nazwanego przetwarzaniem miękkim lub ewolucyjnym [18].

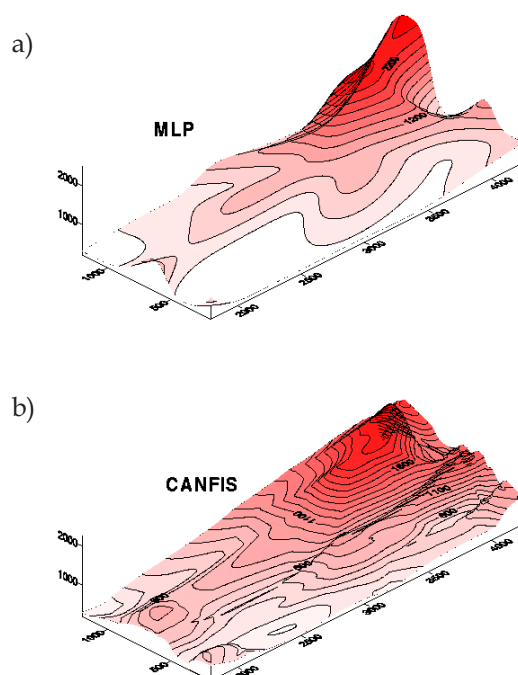
Istnieje wiele algorytmów ewolucyjnych służących do rozwiązywania zadań regresyjnych. W wyniku ich stosowania, podobnie jak w przypadku regresji badanej metodami statystycznymi, otrzymujemy oszacowanie warunkowej wartości oczekiwanej, uwarunkowanej przez stan wektora wejściowego [2, 5, 15, 17]. Wektorem wejściowym x jest para współrzędnych x, y , aczkolwiek można sobie wyobrazić sytuacje, w których może on być rozszerzony o inne zmienne.

Najpowszechniej stosowanym algorytmem budowy neuronowego modelu regresji nieliniowej jest przypuszczalnie perceptron wielowarstwowy MLP (*Multi*

Layer Perceptron), którego architektura, właściwości, techniki uczenia i zastosowania są opisane w wielu publikacjach [6, 7, 15, 19]. Ze względu na swoje właściwości jest on uniwersalnym aproksymatorem funkcji ciągłych. Niezależnie od struktury problemu regresji, przy błędzie średniokwadratowym jako kryterium optymalizacji, perceptron dostarcza oszacowania warunkowej wartości oczekiwanej zmiennych zależnych.

Zbliżonym algorytmem ewolucyjnym jest CANFIS (*Co-Active Neuro-Fuzzy Inference System*), udostępniany w pakiecie NeuroSolutions. Topologia sieci neuronowej jest wzbogacona o moduły rozmywania (fuzyfikacji) i wyostrzenia (defuzyfikacji) przekształcających wejście i wyjście sieci. MLP i CANFIS reprezentują algorytmy sztucznych sieci neuronowych z nielokalnymi [4] funkcjami transferu, charakteryzujące się na ogół dobrymi właściwościami generalizacyjnymi. Oszacowanie wartości funkcji $f(x_n)$ w punkcie x_n jest odpowiedzią sieci na przetworzenie wartości wejściowych, w tym przypadku pary współrzędnych x, y . Odpowiedź ta jest ważoną sumą odpowiedzi jednostek przetwarzających, z reguły wyposażonych w funkcje logistyczne lub tangensoidalne [4, 19].

Rysunek 9 obrazuje pseudo-trójwymiarowy obraz powierzchni regresji trendu zawartości Cr w glebach wyznaczony z zastosowaniem obu algorytmów.



Rys. 9. Pseudo-trójwymiarowa wizualizacja trendu rozkładu Cr w glebach wyznaczonego przy użyciu perceptronu wielowarstwowego MLP (a) oraz algorytmu CANFIS (b)

Obrazy te można uznać za pośrednie między modelami wielomianowymi i interpolacyjnymi: lepiej niż regresja wielomianowa odzwierciedlają one nielinowości, są natomiast pozbawione ekstremów, charakterystycznych dla algorytmów interpolacyjnych. Ze względu na funkcje celu wykorzystane w procesie treningu (minimalizacja błędu średniokwadratowego) reprezentują one przypuszczalnie, lepiej niż interpolacja, najbardziej prawdopodobny przebieg zjawiska. Oczywiście ceną tej zmiany jest wzrost niepewności i wartości odchyleń resztowych.

Tabela 4. Niektóre statystyki rozkładu wartości resztowych oszacowania trendu zawartości Cr w glebach przy użyciu algorytmów ewolucyjnych: MLP i CANFIS. Zawartość Cr w mg/kg gleby

Statystyka	MLP	CANFIS
Minimum	-502,2	-577,4
Maksimum	1866,9	1673,7
Średnia	72,5	-1,5
Mediana	44,9	-27,5
Pierwszy kwartył	-82,7	-131,6
Trzeci kwartył	188,7	104,6
Błąd standardowy	14,9	14,5
Odchylenie standardowe	235,9	231,1

Tabela 4 wskazuje, że oba zastosowane algorytmy dają zbliżone wyniki, przy pewnej przewadze modelu CANFIS. Nie jest wykluczone, że dalsze próby ze zmienionymi architekturami sieci MLP doprowadziłyby do lepszych rezultatów. Należy pamiętać, że w pewnym sensie oszacowania te są bardziej wiarygodne od interpolacji, ponieważ odsiewają one szumy nieuchronnie występujące w obrazie interpolowanym.

Zadania regresyjne mogą wypełniać także sieci neuronowe wyposażone w lokalne funkcje transferu. W odróżnieniu od perceptronów, wyniki przetwarzania tego rodzaju algorytmów dotyczą małych wycinków wielowymiarowej przestrzeni cech, przez co są one uznawane za słabiej generalizujące i pozbawione możliwości ekstrapolacji wyników [4].

Sieci wyposażone w kołowe funkcje bazowe, sieci RBF (*Radial Basis Function*), są chyba najbardziej znaną alternatywą wobec algorytmu MLP metodą, zarówno w zastosowaniach do klasyfikacji, jak i regresji. Istnieje wiele metod optymalizacji tego rodzaju sieci, jak również sposobów wyliczania, kluczowej w tego rodzaju algorytmie, odległości od punktów eksperymentalnych w przestrzeni cech [1, 4, 19], kształtującej stopień wpływu określonego wzorca na odpowiedź sieci.

Wyspecjalizowanym algorytmem regresyjnym jest algorytm sieci neuronowej realizującej regresję uogólnioną GRNN (*Generalized Regression Neural Network*). Jest

to sieciowa, opracowana przez Spechta [8, 14], implementacja dużo wcześniejszej propozycji Parzena [12], dotyczącej regresji jądrowej realizowanej za pośrednictwem tak zwanych „okienek” Parzena. „Okienka”, poprzez zastosowanie symetrycznej funkcji (najczęściej Gaussa), wyznaczają granice superpozycji wzorcowych wartości funkcji w otoczeniu badanego punktu, wagowanych przez lokalną wartość funkcji jądrowej. Zaletą algorytmu GRNN jest krótki czas optymalizacji (trening jednokrokowy). Jego wadą są wysokie wymagania co do pamięci operacyjnej systemu przetwarzającego.

Klasyfikacyjną wersją algorytmu Spechta jest probabilistyczna sieć neuronowa PNN (*Probabilistic Neural Network*) [13]. Zasadniczy algorytm przetwarzania jest podobny jak w sieci GRNN, z tym że odpowiedź sieci jest oszacowaniem prawdopodobieństwa $p(k_c | \mathbf{x})$ wystąpienia w konfiguracji \mathbf{x} klasy k_c . Oznacza to, że wyjściem sieci jest wektor wartości prawdopodobieństw wystąpienia rozróżnianych przez sieć klas. Ze względu na tę właściwość, warto przeanalizować wyniki przetwarzania tej sieci, przy niewielkich modyfikacjach zbioru treningowego. Został on uzupełniony o zmienną nazywaną klasą skażenia K_s , powstałą po podzieleniu rozstępu zawartości Cr w glebach na przedziały o szerokości 250 mg/kg. Klasy $K_0 \dots K_9$ obejmowały tym samym pełny zakres zmienności koncentracji Cr w glebach.

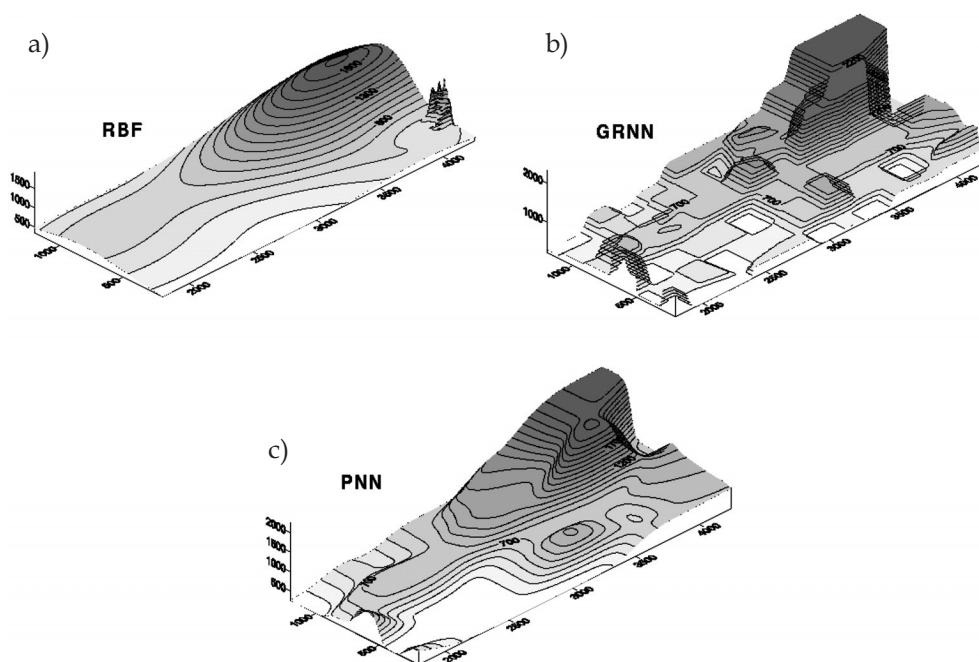
Niektóre statystyczne parametry rozkładu wartości resztowych modeli (rys. 10) zawiera tabela 5. Wynika z niej, że uzyskane w tej próbie wyniki są zbliżone do wcześniej omówionych rezultatów algorytmów ewolucyjnych.

Tabela 5. Niektóre statystyki rozkładu wartości resztowych oszacowania trendu zawartości Cr w glebach przy użyciu algorytmów RBF, GRNN i PNN. Zawartość Cr w mg/kg

Statystyka	RBF	GRNN	PNN
Minimum	-985,6	-346,8	-530,3
Maksimum	1532,8	1822,1	1912,9
Średnia	1,7	1,4	-22,7
Mediana	-6,0	-0,2	-43,3
Pierwszy kwartyl	-125,8	-137,9	-209,9
Trzeci kwartyl	106,2	93,7	122,6
Błąd standardowy	14,9	14,5	17,0
Odchylenie standardowe	236,7	230,9	269,4

Elastyczność algorytmów ewolucyjnych wymaga bardzo ostrożnego traktowania. Jest oczywiste, że głównym ograniczeniem inwentaryzacji stanu zanieczyszczenia gleb na większych obszarach są koszty związane z wydatkami na prace terenowe oraz na oznaczenia laboratoryjne. Przy ograniczonych nakładach (co jest raczej regułą) występuje naturalna skłonność do oszczędnego gospodarowania

punktami rozpoznania, co poprzez zmniejszenie liczby prób prowadzi także do redukcji prac laboratoryjnych. Tendencja ta jest w konflikcie z wielkimi potrzebami algorytmów ewolucyjnych pod względem ilości i reprezentatywności danych.



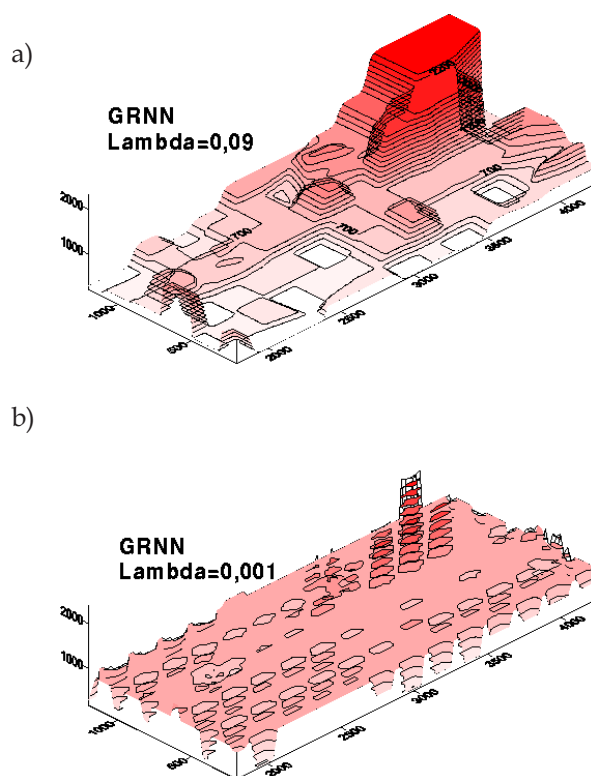
Rys. 10. Pseudo-trójwymiarowa wizualizacja trendu rozkładu Cr w glebach wyznaczonego przy użyciu sieci z radialnymi funkcjami bazowymi RBF (a), realizującej uogólnioną regresję GRNN (b) oraz probabilistycznej PNN (c)

Z reguły zatem liczba danych będzie stanowić górne ograniczenie rozmiarów topologii sieci neuronowych, a co za tym idzie – liczby swobodnie optymalizowanych parametrów. Zawsze jednak będzie istnieć sprzeczność między naturalnym dążeniem do formalnie możliwie najdokładniejszego odwzorowania dostępnych danych a akceptacją jakiegoś poziomu nieokreśloności zjawiska, nakładającej węższy lub szerszy przedział ufności wokół uzyskanego oszacowania. Ten konflikt można zaprezentować na przykładzie sterowania współczynnikiem λ , będącego parametrem kształtu funkcji jądrowej w sieciach PNN i GRNN.

Decyduje on o superpozycji poszczególnych „okienek”, wyznaczanych przez funkcję

$$f(x) \cong f_n(x) = \frac{1}{n\lambda} \sum_{i=1}^n g\left[\frac{x - x(i)}{\lambda}\right] \quad (1)$$

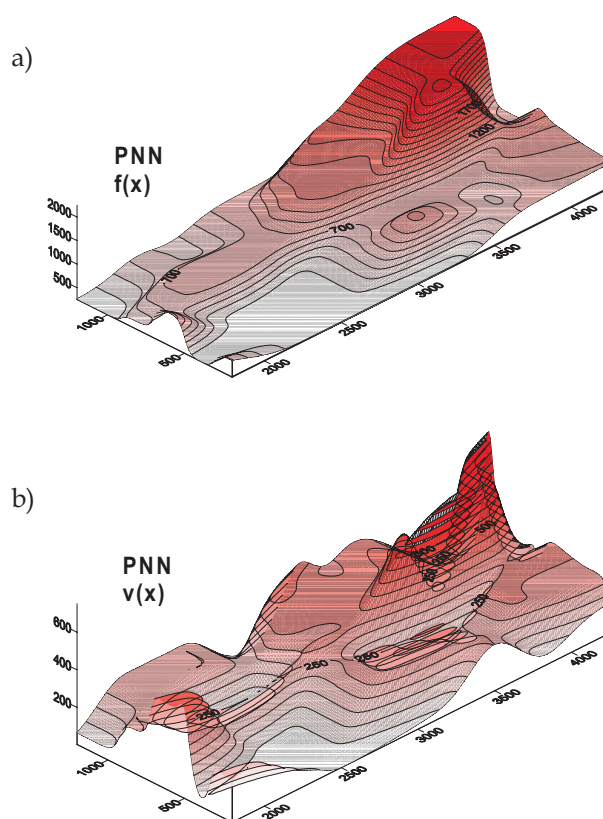
Unimodalna, symetryczna funkcja $g(\cdot)$ (najczęściej eksponencjalna funkcja Gaussa) ma zasięg w przestrzeni wielowymiarowej kształtowany przez λ : duża jej wartość wyznacza większy zasięg jądra, mniejsza – odpowiednio mniejszy. Przy bardzo małej wartości parametru λ algorytm wykazuje formalnie wysoką zgodność z danymi eksperymentalnymi (na przykład wysoki współczynnik korelacji prognozy sieci GRNN z danymi uczącymi), natomiast oczywista jest znikoma lub żadna zdolność interpolacyjna („okienka” nie stykają się). Skutek prezentuje wykres (rys. 11) obrazujący pseudo-trójwymiarowy rzut powierzchni trendu przy dwu różnych wartościach λ .



Rys. 11. Pseudo-trójwymiarowa wizualizacja trendu rozkładu Cr w glebach wyznaczonego przy użyciu sieci GRNN z różnymi wartościami parametru λ : a) $\lambda = 0,09$; b) $\lambda = 0,001$

Obydwa wykresy sygnalizują pewną sztuczność odwzorowania trendu, co wynika głównie z regularności sieci opróbowania w obrębie obszaru o bardzo dużej zmienności przestrzennej obserwowanej cechy, jednakże przypadek skrajnie niskiej wielkości λ pokazuje słabość rozwiązania, przy którym algorytm utracił całkowicie zdolność do interpolacji. Uzasadnia to potrzebę wykonywania wielu eksperymentów przy wykorzystaniu tego algorytmu.

Zastosowanie sieci probabilistycznej PNN, pomimo pewnych trudności interpretacyjnych, daje dodatkowo niezwykle użyteczną informację o lokalnej zmienności przestrzennej cechy, co może być przydatne przy szacowaniu gęstości opróbowania terenu. Problem jest szczególnie widoczny właśnie w terenach z wyraźnie zarysowującym się trendem przestrzennym zanieczyszczenia, zwłaszcza gdy można podejrzewać (jak w tym przypadku) heteroscedastyczność zależności, to znaczy zależność odchylenia o wartości oczekiwanej od wartości wektora wejściowego (w tym przypadku położenia punktu obserwacji). Interpretując wyjściowy wektor sieci PNN jako lokalne prawdopodobieństwa poszczególnych stanów natury (klas zanieczyszczeń) odpowiadające wektorowi wejściowemu x , można tym samym dokonać zgrubnego oszacowania zmienności lokalnej. Przykład takiej interpretacji prezentuje rysunek 12. Informacja taka pozwala na sterowanie gęstością opróbowania na obszarach o dużej przestrzennej zmienności obciążenia gleb polutantami.



Rys. 12. Pseudo-trójwymiarowa wizualizacja trendu rozkładu Cr w glebach wyznaczonego przy użyciu sieci PNN (a) wraz z oszacowaniem lokalnego odchylenia standardowego przestrzennej zmienności (b). Skale pionowe wykresów różne

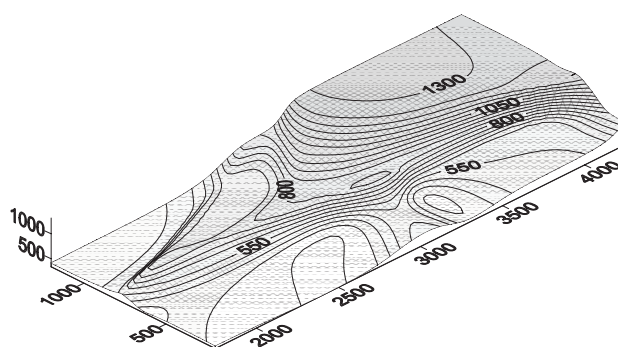
Należy zwrócić uwagę, że zamiar utrzymania jednakowej (co najmniej zbliżonej) dokładności oszacowania stanu zanieczyszczenia gleb na tym terenie uzasadniałby zróżnicowanie gęstości sieci opróbowania, tak by liczba prób z terenów poszczególnych poziomów zmienności zapewniała wymaganą precyzję.

Algorytmem zaprojektowanym do modelowania warunkowych cech rozkładu (warunkowy rozkład, średnia i wariancja) jest rozwiązanie hybrydowe noszące nazwę MDN (*Mixture Density Network*), przedstawione przez Bishopa [2]. Algorytm ten łączy ideę mieszanego modelu gaussowskiego GMM (*Gaussian Mixture Model*) z siecią optymalizacją jego parametrów wejściowych. Optymalizacja polega na minimalizacji błędu w drodze manipulacji udziałami, średnimi i lokalnymi wariancjami modelu GMM. Procedura jest szczególnie cennym uzupełnieniem listy modeli regresyjnych o postępowanie w przypadkach, w których odchylenia od średniej modelu mają rozkład asymetryczny lub są wielomodalne [3]. Uzyskiwany model pozwala na oszacowanie warunkowej wartości średniej oraz warunkowego rozkładu – zależnego od wejścia do modelu – i warunkowej wariancji. Udostępnianie takich informacji bardzo pogłębia znajomość modelowanych zjawisk, zwłaszcza wtedy, gdy potrzebna jest wiedza dotycząca stopnia ryzyka przekroczenia określonych wartości granicznych, zaś założenie dotyczące homoscedastyczności nie jest spełnione. Implementacja programowa – między innymi – tego algorytmu jest udostępniona w postaci zbioru makr o nazwie *Netlab*, przeznaczonych do wykorzystania w środowisku programu Matlab. W podjętych próbach stosunkowo najkorzystniejsze wyniki, jeśli chodzi o liczbę swobodnych parametrów sieci, dała sieć MDN licząca 6 jednostek w warstwie ukrytej (MLP) z modelem gaussowskim złożonym z 8 centrów trenowanych metodą *k*-średnich. Syntetyczną informację dotyczącą oszacowania statystyk wartości resztowych modelu MDN zawiera tabela 6.

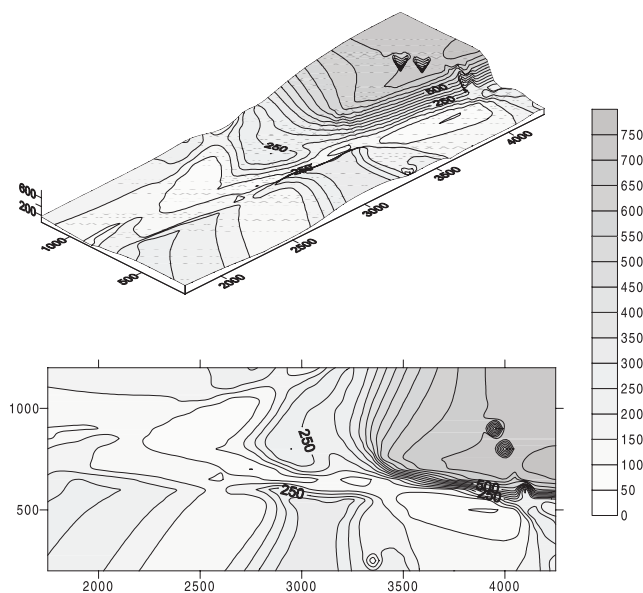
Tabela 6. Niektóre statystyki rozkładu wartości resztowych oszacowania trendu zawartości Cr w glebach przy użyciu algorytmu MDN; zawartość Cr w mg/kg gleby

Statystyka	MLP
Minimum	-703,2
Maksimum	1788,3
Średnia	25,9
Mediana	-16,2
Pierwszy kwartył	-132,6
Trzeci kwartył	126,1
Błąd standardowy	17,1
Odchylenie standardowe	272,4

Wyniki te są zbliżone do rezultatów pozostałych algorytmów z lokalnymi funkcjami transferu. Rysunek 13 przedstawia oszacowanie wartości średnich C_r w omawianym rejonie uzyskane przy użyciu sieci MDN. Rysunek 14 przedstawia warunkowe wartości odchylenia standardowego oszacowania zawartości C_r w glebach. Podobnie jak w przypadku modelu PNN obserwuje się tutaj związek między średnią koncentracją C_r a wahaniami stężenia. Oczywiście wyższe wartości odchylenia standardowego oznaczają większy stopień niepewności oszacowania lokalnej wartości cechy.

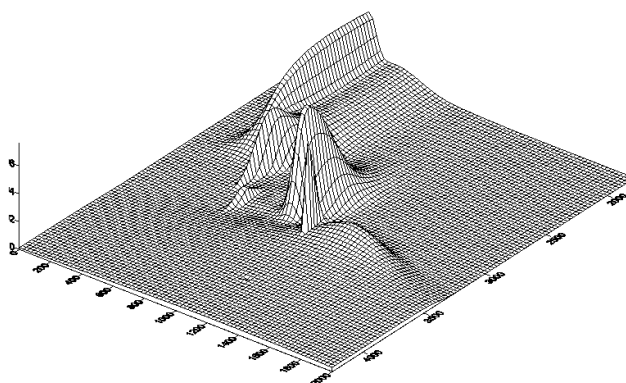


Rys. 13. Pseudo-trójwymiarowa wizualizacja trendu rozkładu C_r w glebach wyznaczonego przy użyciu sieci MDN



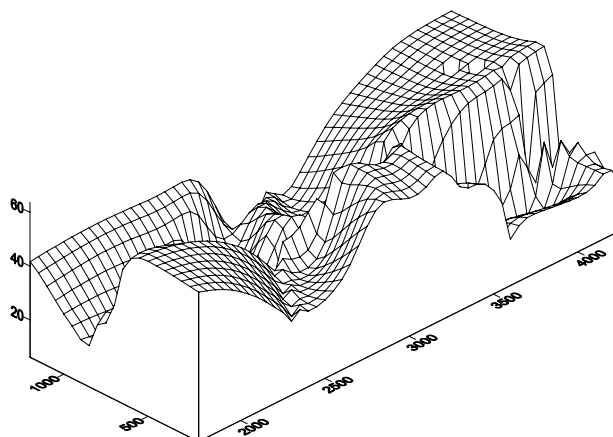
Rys. 14. Wizualizacja oszacowania wartości warunkowego odchylenia standardowego uzyskana przy użyciu sieci MDN

Wyobrażenie o skali zmienności koncentracji Cr daje rysunek 15 prezentujący otrzymane z modelu MDN oszacowanie kształtów jej rozkładu warunkowego obserwowanego wzdłuż wybranego profilu o przebiegu równoleżnikowym. Spłaszczenia rozkładu oznaczają oczywiście potencjalnie znaczniejszy błąd oszacowania średniej, zaś lokalne maksima oznaczają stosunkowo mniejsze ryzyko popełnienia istotnego błędu.



Rys. 15. Pseudo-trójwymiarowa wizualizacja rozkładu warunkowego zawartości Cr w glebach, obserwowanego wzdłuż przekroju równoleżnikowego, wzdłuż północnej granicy obszaru badań. Rozkłady obserwowane od wschodniej strony obszaru badań

Rysunek 16 prezentuje pseudo-trójwymiarowy obraz zróżnicowania współczynnika zmienności, informującego o zakresie zmienności cechy w obszarze opracowania. Jego znajomość umożliwi bezpośrednie oszacowanie gęstości opróbowania terenu zapewniającej uzyskanie wymaganej dokładności.



Rys. 16. Pseudo-trójwymiarowa wizualizacja współczynnika zmienności v [%], oszacowania zawartości Cr w glebach

5. Podsumowanie i wnioski

Rozważając potencjalne cele badania trendu w rozkładzie przestrzennym zanieczyszczenia gleb, można przyjąć, że zazwyczaj mniej istotne jest wykrycie lub uwypuklenie regularności związanej z mechanizmem zjawiska imisji, zaś ważniejsze jest zwykle określenie precyzji informacji oraz granic błędu popełnianego przez akceptację jakiegoś oszacowania zanieczyszczenia. Obserwacje stanu gleb w rejonie silnego skażenia gleb związkami chromu wskazują, że przynajmniej przy dużej koncentracji tego metalu w glebach założenia o dostatecznej gęstości siatki opróbowania rzędu 100 czy 200 metrów są dość optymistyczne, zwłaszcza jeżeli narzędziem określania trendu przestrzennego są tradycyjne, geodezyjne algorytmy interpolacji udostępniane przez pakiety SIT. Zdają one bardzo dobrze egzamin w przypadkach celowego rozmieszczania punktów rejestracji obserwowanej zmiennej w punktach załamania lub odwrócenia trendu, mogą jednak zawodzić, gdy zgromadzone dane są zbierane bez znajomości konfiguracji interpolowanej powierzchni. W takich przypadkach, gdy nie ma istotnych przesłanek do zakładania regularności rozkładu zanieczyszczeń w przestrzeni, dobrą alternatywą wobec interpolacji metodą mogą być algorytmy regresji, w tym metody regresji udostępniane przez sztuczne sieci neuronowe.

Przeprowadzone badania, w warunkach silnej koncentracji zanieczyszczeń i przy oczekiwanej regularności ich rozkładu, uzasadniają wyprowadzenie następujących wniosków:

1. Obserwacje wskazują, że krótkodystansowa zmienność przestrzenna koncentracji Cr w rejonie badań jest bardzo duża. Trudno jest wskazać przyczyny tego stanu rzeczy; mogą nimi być większe niż w Europie Zachodniej rozdrobnienie działek (a zatem także odmiennie użytkowanych pól terenu), nieuchronne błędy opróbowania – zwłaszcza określania głębokości odspojenia próby (przy silnie zarysowującym się spadkowym trendzie zawartości zanieczyszczenia wraz z głębokością) czy też może natura zjawiska akumulowania i wymywania zanieczyszczeń. Niezależnie od przyczyny obraz zanieczyszczeń pozbawiony oszacowania jego precyzji jest co najmniej niepełny, a można go także uznać za fałszywy wtedy, gdy narzędzie generujące go zapewnia pozornie wielką dokładność. Problem tkwi także w tym, że oprócz finansowego, istnieje także obiektywne górne ograniczenie gęstości siatki opróbowania, związane z czytelnością i wiarygodnością rysunku izolinii.
2. Istnieje oczywisty odwrotny związek między precyzją regresji mierzoną średnim kwadratem odchyłań lub entropią a rozmiarem wariancji resztowej. W tym przypadku konieczne jest znalezienie rozsądnej równowagi między pozorną precyzją formalną geodezyjnych algorytmów interpolacyjnych

a ogólnikowością modelu regresji liniowej czy wielomianowej. Oprócz elastyczności wyrażającej się naturalną przydatnością do odzwierciedlenia zależności nieliniowych modele uzyskiwane dzięki algorytmom ewolucyjnym posiadają zaletę łatwości generowania wyniku oszacowania w dowolnym punkcie powierzchni.

3. Przy pewnych założeniach modele sieciowe pozwalają na odzwierciedlenie lokalnej niepewności oszacowania zanieczyszczenia mierzonej lokalną wariancją. Elementarną realizację tego zadania umożliwiają sieci probabilistyczne oraz realizujące model MDN. Możliwa jest także bardziej zaawansowana analiza polegająca na szacowaniu rozkładu prawdopodobieństwa zależnego od wejścia. Jest ona tym istotniejsza, że obok oczekiwanej wartości stężenia zanieczyszczenia jako funkcji położenia punktu w stosunku do źródła emisji może dostarczyć przybliżonej informacji dotyczącej lokalnego kształtu rozkładu niepewności związanej z przestrzennym zróżnicowaniem stężeń, a tym samym także zróżnicowania ocen ryzyka dotyczącego przekroczenia wartości granicznych.

Literatura

- [1] Bishop C.: *Improving the generalization properties of radial basis function neural networks*. *Neural Computation*, 3 (4), 1991, 579–588
- [2] Bishop C.M.: *Mixture density networks*. Raport instytutowy NCRG/94/004, 1994. Dostępne na stronie: citeseer.ist.psu.edu/bishop94mixture.html
- [3] Cornford D., Nabney I.T., Bishop C.M.: *Neutral network based wind vector retrieval from satellite scatterometer data*. *Neural Computing and Application*, 8, 1999, 206–217
- [4] Duch W., Jankowski N.: *Survey of neutral transfer functions*. 1999. Dostępne na stronie: citeseer.ist.psu.edu/duch99survey.html
- [5] Goldberg P.W., Williams C.K.I., Bishop C.M.: *Regression with input-dependent noise: A gaussian process treatment*. [w:] Jordan M.I., Kessarns M.J., Solla S.A. (red.), *Advances In Neural Information Processing Systems*, vol. 10, The MIT Press 1998. Dostępne na stronie: citeseer.ist.psu.edu/article/goldberg-98regression.html
- [6] Hecht-Nielsen R.: *Neurocomputing*. Reading: Addison-Wesley Pub. Co., Reading
- [7] Kerlirzin P., Vallet F.: *Robustness in multilayer perceptrons*. *Neural Computation*, 5(3), 473–482
- [8] Masters T.: *Sieci neuronowe w praktyce. Programowanie w języku C++*. Warszawa, Wydawnictwo Naukowo-Techniczne 1996

- [9] Mitasova H., Brown M., Hofierka J.: *Multidimensional dynamic cartography*. Dostępne na stronie: citeseer.ist.psu.edu/mitasova95multidimensional.html
- [10] Mitasova H., Mitas L.: *Interpolation by regularized spline with tension: I, theory and implementation*. Dostępne na stronie: citeseer.ist.psu.edu/mitasova93interpolation.html
- [11] Nowicka E., Opryszek Z.: *Ocena metod opróbowania gruntów w celu określenia obciążenia gleb metalami ciężkimi*. Kraków, Akademia Górniczo-Hutnicza 2000 (praca magisterska)
- [12] Parzen E.: *On the estimation of a probability density function and mode*. *Annals of Mathematical Statistics*, 33, 1962, 1065–1076. Dostępne na stronie: citeseer.ist.psu.edu/parzen62estimation.html
- [13] Specht D.F.: *Probabilistic neutral networks*. *Neutral Networks*, 3 (1), 1990, 109–118
- [14] Specht D.F.: *A generalized regression neutral network*. *IEEE Transactions on Neutral Networks*, 5, November 1991, 568–576
- [15] Tadeusiewicz R.: *Sieci neuronowe*. W serii: *Problemy współczesnej nauki i techniki, Informatyka*, Warszawa, Akademicka Oficyna Wydawnicza RM 1993
- [16] Trafas M.: *Technologia prowadzenia badań i kartograficznego opracowania wyników dotyczących skażenia gleb w rejonach przemysłowych. Raport końcowy z realizacji projektu badawczego KBN 8 T 12E 007 20*, Kraków, Akademia Górniczo-Hutnicza 2004
- [17] Weigend A., Nix D.: *Predictions with confidence intervals (local error bars)*. 1994. Dostępne na stronie: citeseer.ist.psu.edu/weigend94predictions.html
- [18] Zadeh L.A.: *Fuzzy sets*. *Information and Control*, 8 (3), June 1965, 338–353
- [19] Żurada J., Barski M., Jędruch W.: *Sztuczne sieci neuronowe*. Warszawa, Wydawnictwo Naukowe PWN 1996