

Adam Sagan, Mariusz Łapczyński

Uniwersytet Ekonomiczny w Krakowie
e-mails: sagana@uek.krakow.pl; lapczynml@uek.krakow.pl

**MODELE HYBRYDOWE SEM-TREE
W BADANIACH RÓŻNICOWANIA SIĘ
PREFERENCJI CZŁONKÓW POLSKICH
GOSPODARSTW DOMOWYCH**

**SEM-TREE HYBRID MODELS
IN THE PREFERENCES ANALYSIS
OF THE MEMBERS OF POLISH HOUSEHOLDS**

DOI: 10.15611/pn.2018.508.20

JEL Classification: C51, B54

Streszczenie: Celem artykułu jest identyfikacja wymiarów kształtowania się preferencji względem strategii alokacji zasobów stosowanych przez członków polskich gospodarstw domowych, aby dokonać ich segmentacji. Wymiary te wyodrębniono na podstawie ogólnopolskich danych zebranych na reprezentatywnej próbie 1100 respondentów w 410 gospodarstwach domowych. W analizie wyników wykorzystano modele hybrydowe SEM-Tree, stanowiące połączenie konfirmacyjnych modeli strukturalnych z eksploracyjnymi i predykcyjnymi modelami drzew klasyfikacyjnych i regresyjnych. Pozwala to na zastosowanie rozpoznawczego podejścia do budowy modeli strukturalnych dla heterogenicznych populacji i ocenę wpływu zmiennych klasyfikacyjnych na identyfikację segmentów, w których występują możliwe jednorodne parametry modelu strukturalnego (SEM). Podejście to łączy zalety podejścia modelowego (na etapie budowy hipotez dotyczących relacji strukturalnych i specyfikacji modeli pomiarowych) i oparteo na danych podejścia eksploracyjnego.

Słowa kluczowe: postawy konsumenckie, modelowanie strukturalne, drzewa klasyfikacyjne, modele hybrydowe.

Summary: The purpose of the paper is to identify the dimensions of the strategy of resources allocation of Polish households members. These dimensions were identified on the basis of nationwide empirical data gathered on a representative sample of 1100 respondents nested in 410 households. SEM-Tree hybrid models are used in the analysis of the results, which combine the confirmatory structural equation models with exploratory and predictive classification and regression trees. This allows to apply structural modeling for the study of heterogeneous populations and to assess the hierarchical impact of exogenous predictors on the identification of segments with separate and unique model structural parameters. The approach combines the advantages of a model approach (at the stage of constructing hypotheses

on structural relationships and specifications of measurement models) and exploration-based data (at the stage of recursive division of the sample).

Keywords: customers' attitudes, structural equation models, decision trees, hybrid models.

1. Wstęp

Stosowane metody analizy danych dzieli się często na metody confirmacyjne, mające zastosowanie w przypadku dobrze ustrukturalizowanych problemów badawczych i służące do weryfikacji (falsyfikacji) modeli teoretycznych, oraz metody eksploracyjne, które związane są z rozwiązywaniem problemów nieustrukturalizowanych i formułowaniem propozycji teoretycznych „wyłaniających” się z analizowanych danych. Podział ten jest uzupełniany o podejścia hybrydowe stanowiące złożenie metod często należących do różnych tradycji analitycznych, korzystających z innych założeń metodologicznych.

Tego typu metody hybrydowe wykorzystywane są zarówno w takich podejściach confirmacyjnych, jak modelowanie strukturalne ze zmiennymi ukrytymi (*Structural Equation Modelling* – SEM), jak i w podejściach eksploracyjnych, do których należy metoda drzew klasyfikacyjnych. Przykładami podejść hybrydowych w modelowaniu strukturalnym są tzw. automatyzowane modele strukturalne, które stanowią połączenie modeli strukturalnych i heurystycznych procedur numerycznych. W ich ramach modele SEM są łączone z heurystycznymi algorytmami mrówkowymi (*ant colony optimization* – ACO-SEM), algorytmami genetycznymi (*genetic algorithms* – GA-SEM), poszukiwaniem tabu (*tabu search* – TS-SEM), procedurą niszczenia i odtwarzania (*ruin-and-recreate* – R&R-SEM), czy symulacyjnym wyżarzaniem (*simulated annealing* – SA-SEM) [Marcoulides, Ing 2012; Sagan, Perek-Białas 2016]. Ich zastosowanie wiąże się z próbami heurystycznego poszukiwania specyfikacji modelu (*specification search*), które wynikają z eksploracyjnego charakteru procesu poszukiwania specyfikacji, złożoności modelu, zmiennych i możliwych kombinacji liczby potencjalnych zależności między nimi. Bez uwzględnienia jakichkolwiek założeń teoretycznych liczba możliwych modeli SEM zbudowanych na danej macierzy kowariancji wynosi bowiem $n = 4^{p(p-1)/2}$.

W podejściach eksploracyjnych, takich jak analiza drzew klasyfikacyjnych, modele hybrydowe stanowią połączenie analizy CHAID i CART oraz modeli regresyjnych (modeli logitowych i probitowych budowanych w liściach drzewa). Według podobnego schematu łączy drzewa klasyfikacyjne z analizą skupień (np. metodą *k*-średnich). Warto dodać, że stosuje się także algorytmy eksploracji danych, np. SOM, GHSOM [Łapczyński 2016].

Jedną z hybrydowych metod pozwalających na łączenie zalet confirmacyjnego podejścia SEM i eksploracyjnych drzew klasyfikacyjnych jest model SEM-Tree.

Celem artykułu jest identyfikacja wymiarów kształtowania strategii alokacji zasobów stosowanych przez członków polskich gospodarstw domowych i dokonanie ich segmentacji. Podejście hybrydowe pozwala na wyodrębnienie segmentów gospodarstw domowych na podstawie układów zależności między strategiami podejmowania decyzji w gospodarstwie (altruistyczną a rywalizacyjną) a preferencjami w zakresie zasobów konsumenta (reputacji, czasu i pieniądza).

2. Charakterystyka podejścia SEM-Tree

SEM-Tree to wielowymiarowa metoda statystyczna łącząca podejście konfirmacyjne (modele strukturalne SEM) i eksploracyjne (podział rekurencyjny znany z drzew klasyfikacyjnych i regresyjnych) [Brandmaier i in. 2013a]. Graficzna postać modelu SEM-Tree przypomina strukturę drzewa klasyfikacyjnego, w którego węzłach znajdują się oszacowane modele strukturalne. Podział węzłów jest dokonywany na podstawie kowariant, które w drzewach klasyfikacyjnych traktowane są jako zmienne niezależne, i nie są wykorzystywane do budowy modeli strukturalnych. Budowa modeli SEM w węzłach potomnych jest oparta na podzbiorach zbioru danych.

Algorytm budowy modelu SEM-Tree jest następujący [Brandmaier i in. 2013a, s. 75]:

1. Dopasowanie parametrycznego modelu strukturalnego (tzw. modelu szablonowego) do bieżącego zbioru obserwacji.
2. Binarny podział zbioru danych na wszystkie możliwe sposoby z wykorzystaniem do tego celu wszystkich kowariant. W wydzielonych podzbiorach zbudowanie modeli strukturalnych i porównanie dopasowania modelu złożonego (*compound model*) z dopasowaniem modelu szablonowego.
3. Wybór tego z modeli złożonych, który najlepiej opisuje dane z punktu widzenia przyjętego kryterium. Jeżeli ten model jest dopasowany lepiej niż model szablonowy, powtórzenie procedury od kroku pierwszego. W innym wypadku zakończenie procedury optymalizacji.

Model strukturalny M zbudowany na podstawie wszystkich przypadków (w węźle początkowym drzewa) nazywany jest albo modelem szablonowym (*template model*), albo modelem sprzed podziału (*pre-split model*), albo modelem bazowym (*base model*). Powstaje z wykorzystaniem wskaźnika dopasowania modelu, którym zazwyczaj jest wskaźnik najwyższej wiarygodności (*maximum likelihood index*). Zbiór danych \mathbf{D} jest reprezentowany przez macierz o wymiarach $n(k+1)$, gdzie n oznacza liczbę przypadków, k – liczbę obserwowanych zmiennych a l – liczbę kowariant nieuwzględnionych w modelu strukturalnym. Podział na zmienne obserwowane i towarzyszące oznacza, że macierz \mathbf{D} zostaje podzielona na odpowiednio: podmacierz \mathbf{D}_k i podmacierz \mathbf{D}_l .

Kowarianty są wykorzystywane do podziału węzłów, a ich wartości lub kategorie decydują o tym, do którego węzła potomnego zostanie przydzielony obiekt ze zbioru obserwacji. W celu włączenia do analizy kowariant ciągłych (ilościowych), porząd-

kowych i nominalnych, wszystkie wielowartościowe i wielowariantowe zmienne są zamieniane na zbiory kowariant dychotomicznych. Sposób ich przekształcenia jest zależny od typu (poziomu pomiaru). Niech N oznacza liczbę kategorii albo wartości zmiennych towarzyszących. Zgodnie z zaproponowaną procedurą binaryzacji, kowarianty ciągłe i porządkowe mogą być zdychotomizowane na $N - 1$ sposobów, kowarianty kategoryjne zaś na $2^{(N-1)} - 1$ sposobów. Dla każdego możliwego podziału w węzłach potomnych buduje się modele strukturalne, które są nazywane podmodelami (*submodels*). Ich suma jest następnie określana terminem „model po podziale” (*post-split model*) albo terminem „model złożony” (*compound model*).

Ponieważ model złożony i model szablonowy (bazowy) są zagnieżdżone, stosuje się test ilorazu wiarygodności w celu określenia, czy model może zostać podzielony na podmodele. Iloraz wiarygodności ma rozkład chi-kwadrat przy hipotezie zerowej, która mówi, że kowarianta „nie wpływa” na model (dopasowanie modelu szablonowego nie różni się istotnie od dopasowania modelu złożonego). Na każdym poziomie drzewa wybiera się kowariantę z najwyższą wartością logarytmu ilorazu wiarygodności. Ten sposób postępowania jest kontynuowany rekurencyjnie w kolejnych etapach podziału drzewa. W procedurze występują naturalne kryteria stopu [Brandmaier i in. 2013b, s. 105]:

- 1) brakuje kowariant, które mogłyby dzielić węzły drzewa,
- 2) liczba obserwacji w liściu (węzle końcowym) jest niższa od progowej – ustalonej przez badacza,
- 3) osiągnięto pożądaną głębokość drzewa (liczbę podziałów),
- 4) najlepszy podział nie jest wystarczająco dobry.

Czwarte kryterium jest stosowane, żeby zapobiec nadmiernemu dopasowaniu się modelu do danych.

Formalna ocena podziału węzła wygląda następująco:

1. Model szablonowy M charakteryzuje się funkcją wiarygodności dla zestawu parametrów θ z m wolnymi parametrami i zbiorem danych \mathbf{D} .

2. Otrzymuje się zestaw parametrów θ_F dla całego zbioru obserwacji poprzez minimalizację.

3. Dla danego potencjalnego podziału węzła można przyjąć, że zbiór danych \mathbf{D} zostanie podzielony na rozłączne podzbiory $\mathbf{D}_1, \dots, \mathbf{D}_k$. Ponieważ podzbiory są rozłączne, więc parametry modeli strukturalnych $\theta_1, \dots, \theta_k$ są szacowane niezależnie poprzez minimalizację funkcji wiarygodności.

4. Model złożony ze wszystkich pod modeli jest od teraz określany symbolem M_{SUB} . Model szablonowy M jest zagnieżdżony w obrębie M_{SUB} , ponieważ M odpowiada M_{SUB} z n dodatkowymi liniowymi ograniczeniami równości parametrów θ .

5. Biorąc pod uwagę tę zagnieżdżoną strukturę, można sformułować hipotezę zerową, która mówi, że dopasowanie modelu szablonowego M nie różni się istotnie od dopasowania modelu złożonego M_{SUB} , czyli $H_0: \theta_1 = \theta_2$.

6. Logarytm ilorazu wiarygodności pomiędzy M i M_{SUB} wyraża się wzorem: $\Lambda = -2LL(\textit{presplit}) - \sum 2LL(\textit{postsplit})$

7. Ze względu na to, że przyjmuje rozkład chi-kwadrat z $(k-1)$ stopniami swobody i biorąc pod uwagę wszystkie kowarianty, ocenie podlegają wszystkie możliwe podziały węzła. Model z najwyższym przyrostem dopasowania jest porównywany z poprzednio wybranym progmem wyznaczonym przez poziom istotności α . Jeżeli podział jest statystycznie istotny, procedura budowy drzewa jest kontynuowana.

3. Charakterystyka danych

Model SEM-Tree został zbudowany na podstawie danych pochodzących z ogólnopolskich badań ankietowych dotyczących analizy struktury preferencji i wartości polskich gospodarstw domowych w zakresie alokacji zasobów na konsumpcję, oszczędzanie i inwestowanie. Badania zostały przeprowadzone w 2013 roku na próbie 1100 respondentów pochodzących z 410 gospodarstw domowych (wywiady przeprowadzono w rodzinach z ojcem, matką i przebywającym w domu najstarszym dzieckiem powyżej 16 roku życia).

Zmienne wybrane do modelu SEM stanowią wskaźniki trzech konstruktów związanych z postawami wobec alokacji zasobów i strategią podejmowania decyzji. Pierwszy dotyczy strategii podejmowania decyzji alokacyjnych w rodzinie względem wymiaru altruistycznego (dbanie o dobro wspólne rodziny) i rywalizacyjnego (posiadanie własnych, niezależnych „budżetów” przez członków rodziny). Pomiar był dokonywany za pomocą skal Likerta i następujących stwierdzeń dotyczących altruistycznego bieguna wymiaru (A):

- p23 – „Rodzina powinna ograniczać wydatki na indywidualne potrzeby na rzecz zaspokojenia wspólnych”;
- p24 – „Dobro wspólne całej rodziny jest ważniejsze niż członka rodziny z osobna”;
- p25 – „Poczucie spełnienia dają dobra przeznaczone dla całej rodziny”;
- p26 – „Radość życia czerpie się w rodzinie bardziej z dóbr, które służą wszystkim”.

Dwie pozostałe zmienne ukryte dotyczyły kompromisów wynikających z relacji między reputacją a dochodami oraz czasem wolnym a dążeniem do wyższych dochodów i bezpieczeństwa materialnego (dylemat „twarz” a „pieniądz” i „czas” a „pieniądz”).

1. Y1 – „czas” a „pieniądz”:
 - p51 – „Zarabianie i wydawanie pieniędzy jest ważniejsze niż czas wolny”;
 - p52 – „Zarabianie na bezpieczne jutro jest ważniejsze niż czas wolny”.
2. Y2 – „twarz” a „pieniądz”:
 - p54 – „Cel uświęca środki – pieniądze są ważniejsze niż własna reputacja”;
 - p55 – „Można stracić dobre imię, byle za bezpieczną finansową przyszłość”.

Dodatkowymi zmiennymi towarzyszącymi, które uczestniczyły w budowie drzewa SEM-Tree (predyktorami), były: m1 – płeć, m2 – wiek, m3 – wykształcenie, m5 – subiektywna sytuacja materialna.

4. Wyniki

4.1. Model CART na wartościach czynnikowych

W modelu strukturalnym występowały trzy zmienne (A, Y1 i Y2), co wymusiło zbudowanie trzech modeli drzew dla trzech różnych zmiennych zależnych. Ponieważ zmienne A, Y1 i Y2 są ilościowe, a ich wartościami są wartości czynnikowe, więc powstały tutaj modele regresyjne ze zbiorem zmiennych niezależnych (kowariant) obejmującym: płeć, wiek, wykształcenie i ocenę sytuacji materialnej respondenta¹. Do budowy modeli drzew regresyjnych wykorzystano algorytm C&RT z minimalną liczebnością liścia równą 25 i przycinaniem według wariancji. Ostatecznie, znając rodzaj i kierunek zależności pomiędzy A, Y1 i Y2, zdecydowano się na wybór zmiennej zależnej Y1. Ze względu na dużą liczbę węzłów końcowych w drzewie optymalnym zredukowano je do dwóch poziomów, gdyż taka właśnie ograniczona postać jest zalecana podczas hybrydyzacji tego narzędzia analitycznego z innymi metodami [Steinberg, Cardell 1998]. Zredukowany model drzewa regresyjnego przedstawia rys. 1. Warto dodać, że wariancja wyjaśniona przez model optymalny była bardzo niska (3,2%), a po redukcji obniżyła się do poziomu 1,6%. Nie jest to w tym miejscu kryterium przesądzające o jakości rozwiązania całej hybrydy SEM Tree, a świadczy jedynie o tym, że zestaw kowariant w bardzo niewielkim stopniu wyjaśnia wariancję zmiennej Y1. Po redukcji modelu można zauważyć, że ujemne wartości czynnikowe zmiennej Y1 odnoszą się do kobiet z przedziałów wieku 18-34 i 50-64 (liść nr 4), dodatnie zaś wartości czynnikowe do mężczyzn w tym wieku (liść nr 5) oraz do wszystkich osób (bez względu na płeć) z przedziałów wieku: 35-49 i 65+ (liść nr 3).

Następny etap procedury hybrydyzacyjnej obejmuje budowę modelu szablonowego w węzle początkowym drzewa (nr 1) i modelu złożonego we wszystkich liściach drzewa. Parametry modelu szablonowego są przedstawione w tab. 1.

Model ścieżkowy jest nasycony ($\chi^2 = 0,00$ ($df = 0$), $p = 1,00$) i interpretacji podlegają jedynie parametry ścieżkowe. Wynika z nich, że strategia altruistyczna sprzyja silniej skłonności do poświęcenia czasu wolnego w celu uzyskania wyższych dochodów niż poświęcenia własnej reputacji w tym celu (aczkolwiek zależności wskazują, że dbanie o rodzinę jako całość jest silniej związane z „poszukiwaniem dochodów” niż w przypadku występowania autonomii dochodowej poszczególnych członków rodziny).

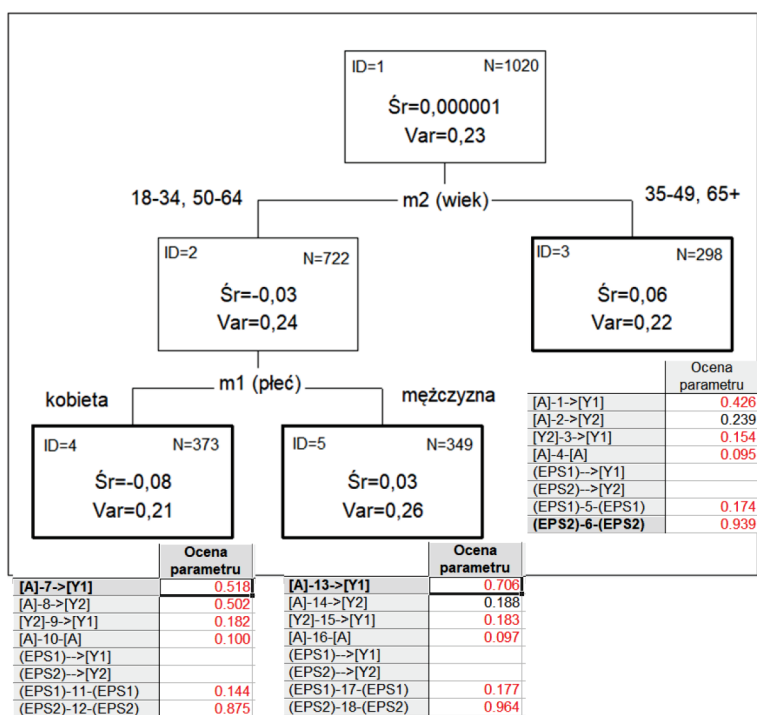
¹ Zmienna „wiek” posiadała 5 kategorii: 18-24, 25-34, 35-49, 50-64 i 65+. Zmienna „wykształcenie” posiadała 6 kategorii: podstawowe, gimnazjalne, zasadnicze zawodowe, średnie, wyższe 1. stopnia i wyższe 2. stopnia. Zmienna „sytuacja materialna” posiadała 5 kategorii opisujących stopień zamożności badanego: bardzo dobra, dobra, przeciętna, zła i bardzo zła.

Tabela 1. Parametry szablonowego modelu ścieżkowego

Regressions:				
	Estimate	Std.Err	z-value	P(> z)
Y1 ~				
A	0.563	0.041	13.722	0.000
Y2 ~				
A	0.324	0.096	3.367	0.001
Y1 ~				
Y2	0.175	0.013	13.181	0.000
Variances:				
	Estimate	Std.Err	z-value	P(> z)
.Y1	0.167	0.007	22.583	0.000
.Y2	0.925	0.041	22.583	0.000

Źródło: wydruk biblioteki *lavaan* programu R.

W kolejnym etapie zbudowane zostały modele ścieżkowe w poszczególnych liściach drzewa klasyfikacyjnego.



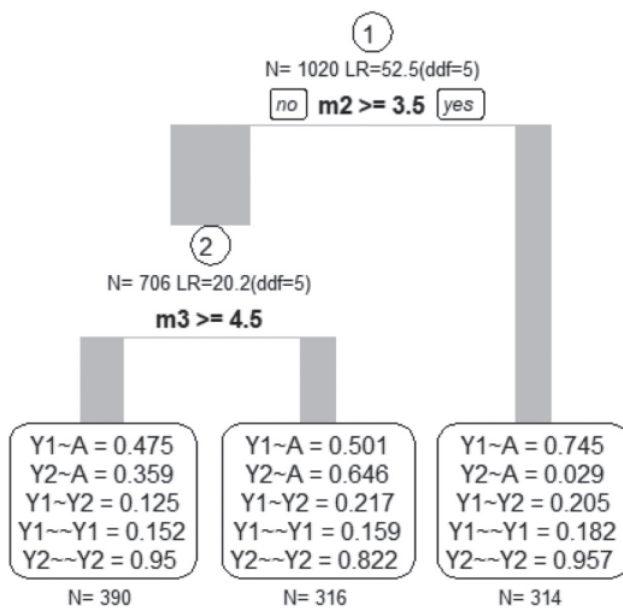
Rys. 1. Zredukowany model drzewa regresyjnego C&RT ze zmienną zależną Y1 i modelami ścieżkowymi

Źródło: opracowanie własne na podstawie wyników modelowania w programie Statistica.

Z analizy submodeli SEM wynika, że kobiety z przedziałów wiekowych 18-34 i 50-64 charakteryzuje podobna, wysoka i jednakowo silna zależność między wy- miarem altruistycznym a preferencją pieniądza nad czasem wolnym i preferencją pieniądza a reputacją, mężczyźni w podobnym wieku cechuje silna zależność między preferencją pieniądza nad czasem wolnym, a wszystkie osoby (bez względu na płeć) z przedziałów wiekowych: 35-49 i 65+ skłonne są raczej wraz ze wzrostem strategii altruistycznej do preferencji pieniądza nad czasem wolnym.

4.2. Model SEM-Tree

Wadą modeli SEM budowanych w liściach drzewa klasyfikacyjnego jest konieczność uwzględnienia predefiniowanej zmiennej zależnej w procesie budowy drzewa. W prezentowanym na wstępie algorytmie podziały obejmują wszystkich kombinacje zbinaryzowanych predyktorów bez konieczności uwzględniania zmiennej zależnej. Rysunek 2 przedstawia finalny model hybrydowy SEM-Tree obliczony za pomocą biblioteki *SEM-Tree* programu R.



Rys. 2. Model SEM-Tree

Źródło: opracowanie własne na podstawie biblioteki *SEM-Tree* programu R.

Z modelu wynika, że najważniejszymi predyktorami okazały się: wiek i poziom wykształcenia. Na podstawie kryterium alfa ($\alpha = 0,05$) i minimalnej liczebności w liściach ($n = 250$), zostały wyodrębnione 3 segmenty członków rodzin. Najmniej licze-

ny segment reprezentują osoby powyżej 49 lat ($m_2 > 3.5$), które charakteryzują się silnym związkiem między strategią altruistyczną a preferencją pieniądza nad czasem wolnym (nie ma zależności między strategią altruistyczną a preferencją pieniądza nad reputacją). Drugim segmentem są osoby o wykształceniu wyższym ($m_3 > 4,5$) do 49 lat ($m_2 < 3.5$), wśród których występuje nieco silniejszy związek między strategią altruistyczną a preferencją pieniądza nad reputacją niż preferencją pieniądza nad czasem wolnym. Ostatni, najliczniejszy segment, to osoby o wykształceniu średnim ($m_3 > 4,5$) do 49 lat ($m_2 < 3.5$), które cechują się względnie silniejszym związkiem między strategią altruistyczną a preferencją pieniądza nad czasem wolnym. Zmienna „stopień zamożności” nie pojawiła się w drzewie, ponieważ inne predyktory w wyższym stopniu redukowały wariancje w węzłach potomnych.

5. Zakończenie

Modele hybrydowe SEM-Tree pozwalają na eksploracyjną analizę zależności ścieżkowych w heterogenicznej populacji. W odróżnieniu od klasycznych drzew klasyfikacyjnych z modelami SEM w liściach, umożliwiają wyodrębnienie predyktorów związanych ze zróżnicowaniem nie tylko poziomu zmiennej zależnej, lecz także relacji ścieżkowych między zmiennymi. Oba rodzaje modeli potwierdzają hipotezę o wpływie strategii podejmowania decyzji (altruistycznej) na poświęcanie wolnego czasu i reputacji dla uzyskiwania wyższych dochodów w gospodarstwie domowym.

Podsumowując wyniki modelu, należy podkreślić, że wybór strategii altruistycznej (dbanie o „wspólne dobro” z silną kontrolną funkcją głowy rodziny) silnie skłania członków gospodarstwa domowego do przedkładania wartości pieniądza („walki o byt”) nawet kosztem czasu wolnego i utraty własnej reputacji. Można więc przyjąć, że decentralizacja budżetów, większa autonomiczność i tym samym poczucie bezpieczeństwa członków rodzin skłaniają silniej do ujawniania wyższej preferencji czasu wolnego, dbania o własną reputację i tym samym bardziej harmonijnego życia indywidualnego.

Literatura

- Brandmaier A.M., Oertzen T., McArdle J.J., Lindenberger U., 2013a, *Structural Equation Model Trees*, Psychological Methods, vol. 18, no. 1, s. 71-86.
- Brandmaier A.M., Oertzen T., McArdle J.J., Lindenberger U., 2013b, *Exploratory Data Mining with Structural Equation Model Trees*, [w:] McArdle J.J., Ritschard G. (red.), *Contemporary Issues in Exploratory Data Mining in the Behavioral Sciences*, Routledge, New York, s. 96-127.
- Łapczyński M., 2016, *Modele hybrydowe w marketingu relacji*, Wydawnictwo UEK, Kraków.
- Marcoulides G.A., Ing M., 2012, *Automated Structural Equation Modeling Strategies*, [w:] Hoyle R.H., (red.), *Handbook of Structural Equation Modeling*, Guilford Press, New York.

- Sagan A., Perek-Białas J., 2016, *Eksplozacyjne podejścia w budowie modeli strukturalnych*, [w:] Grześkowiak A., Mazurek-Łopacińska K., Sobocińska M., Stanimir A. (red.), *Metody badań marketingowych: modelowanie, technologia, wizualizacja*, Wydawnictwo Uniwersytetu Ekonomicznego we Wrocławiu, Wrocław.
- Steinberg D., Cardell N.S., 1998, *The Hybrid Cart-logit Model in Classification and Data Mining*, Eighth Annual Advanced Research Techniques Forum, American Marketing Association, Salford Systems, s. 1-7.