

KAMIL JODŹ

Stochastyczne modelowanie intensywności zgonów na przykładzie Polski

Streszczenie

W artykule zostaną przedstawione różne sposoby stochastycznego modelowania polskiej intensywności zgonów. Zaprezentowane zostaną dwa, popularne w ostatnich latach podejścia: model Lee-Cartera oraz model Renshaw-Habermana. Za pomocą powyższych modeli można uchwycić efekt starzenia się populacji związany zarówno z latami kalendarzowymi, jak i z generacją (kohortą), do której należą badane osoby. Zaletą tych modeli jest również możliwość ekstrapolacji poszczególnych oszacowanych składników. Pozwala to na stworzenie dla Polski prognoz intensywności zgonów. W artykule zostanie również zaprezentowane porównanie powyższych modeli oszacowanych na podstawie polskich danych oraz analiza dalszego trwania życia. Badania zostały przeprowadzone dla obu płci.

1. Wstęp

W krajach wysoko rozwiniętych bardzo wyraźnie widać proces wydłużania się wieku obywateli. Z punktu widzenia ludzi (z oczywistych powodów) jest to zjawisko pożądane, lecz niesie ze sobą również wiele problemów natury społeczno-gospodarczej. Systemy emerytalny oraz służby zdrowia wymagają coraz większych nakładów finansowych, aby zapewnić bezpieczeństwo starzejącej się populacji. Pieniądze na te cele są pozyskiwane z obowiązkowych składek. W tej sytuacji w naturalny sposób pojawia się pytanie: jak określić wysokość takiej składki? Problem ten jest nierozdzielnie związany z długością życia ubezpieczonego, dlatego tak istotne jest znalezienie efektywnego (trafnego) sposobu prognozowania dalszego trwania życia.

W literaturze można znaleźć bardzo wiele różnych podejść do modelowania umieralności oraz dalszego trwania życia. Są to metody zarówno analityczne, jak i oparte na tablicach trwania życia. Do najstarszych i najprostszych metod analitycznych można zaliczyć sposoby modelowania umieralności zaproponowane przez de Moivre'a, Weibulla, Gompertza lub Makehama. W modelach tych (dwa ostatnie) natężenie zgonów μ ma postać funkcji wykładniczej:

- prawo de Moivre'a:

$$\mu_t = \frac{1}{\omega - t},$$

- prawo Weibulla:

$$\mu_t = A \cdot t^n,$$

- prawo Gompertza:

$$\mu_t = A \cdot B^t,$$

- prawo Makehama:

$$\mu_t = C + A \cdot B^t,$$

gdzie A , B , C są parametrami modelu, ω jest wiekiem granicznym, t zaś oznacza czas. Korzenie szacowania umieralności metodą tablic trwania życia sięgają XVII wieku. Metoda w pewnym rozszerzeniu jest stosowana do dzisiaj np. przez Główny Urząd Statystyczny.

Praca zawiera opis metody prognozowania przeciętnego dalszego trwania życia mającej swoje początki w pracach Lee i Cartera (LC). Jest to metoda o charakterze stochastycznym. Liczba zgonów jest w tym przypadku modelowana rozkładem Poissona. Model ten, w przeciwieństwie do poprzednich, nie tylko opisuje spadek umieralności związany z wiekiem i rokiem kalendarzowym, ale również uwzględnia w pewnym stopniu wpływ faktu przynależenia do konkretnego pokolenia (tzw. efekt kohortowy). W dalszej części zostanie również omówiony model Renshaw-Habermana (RH), który jest pewnym rozszerzeniem modelu Lee-Cartera.

2. Modele LC i RH oraz ich estymacja

2.1. Dane oraz równania bilansujące

Rozważamy współczynniki umieralności mieszkańców Polski, bez podziału na kobiety i mężczyzn. Wykorzystamy dane o wielkości populacji, liczbie zgonów oraz urodzeń zgromadzone w „Human Mortality Database”. Przy tworzeniu modelu bierzemy pod uwagę dane roczne, począwszy od 1960 roku. Wprowadźmy następujące oznaczenia: $\mu_{x,t}$ – intenswność umieralności, $P_{x,t}$ – liczba ludności w wieku x na początku roku kalendarzowego t , B_t – liczba urodzeń w roku t , $D_{x,t}$ – liczba zgonów osób w wieku x w roku kalendarzowym t , $I_{x,t}$ – saldo migracji osób w wieku x w roku kalendarzowym t , $t = 1, 2, \dots, T$, $x = 0, 1, \dots, x_{\max}$. Załóżmy, że natężenie zgonów jest stałe w ciągu całego roku kalendarzowego t oraz saldo emigracji jest równomiernie rozłożone (van Imhoff, 1990). Przy tych założeniach równania bilansowe są następującej postaci:

$$P_{0,t} = [1 - \exp(-\mu_{0,t})] \cdot [B_t + I_{0,t}] / \mu_{0,t},$$

dla $x = 0$;

$$P_{x,t} = P_{x-1,t-1} \cdot \exp(-\mu_{x,t}) + [1 - \exp(-\mu_{x,t})] \cdot I_{x,t} / \mu_{x,t},$$

dla $0 < x < x_{\max}$;

$$P_{x_{\max},t} = [P_{x_{\max}-1,t-1} + P_{x_{\max},t-1}] \cdot \exp(-\mu_{x_{\max},t}) + \\ + [1 - \exp(-\mu_{x_{\max},t})] \cdot I_{x_{\max},t}/\mu_{x_{\max},t},$$

dla $x = x_{\max}$.

Powyższe równania są nieliniowe ze względu na μ , dlatego do wyznaczenia intensywności zgonów wykorzystamy metody numeryczne.

2.2. Klasyczny model Lee-Cartera

Lee oraz Carter w swojej pionierskiej pracy modelują logarytm rocznego współczynnika umieralności $m_{x,t}$ (roczny współczynniki zgonów osób w wieku x w roku kalendarzowym t) jako następującą funkcję:

$$\ln(m_{x,t}) = \alpha_x + \beta_x \cdot \kappa_t + \varepsilon_{x,t},$$

gdzie:

- $m_{x,t}$ – roczne współczynniki natężenia zgonów,
- α_x – średnie względem czasu poziomy umieralności, estymowane w następujący sposób: $\hat{\alpha}_x = \frac{1}{T} \sum_t \ln(m_{x,t})$,
- κ_t – opisuje zmiany poziomów umieralności w czasie,
- β_x – modyfikacja wartości κ_t w zależności od wieku x ,
- $\varepsilon_{x,t}$ – niezależne zmienne losowe o średniej zero oraz stałej wariancji.

Uwzględnia się również warunki zapewniające identyfikację modelu:

$$\sum_t \kappa_t = 0, \quad \sum_x \beta_x = 1.$$

Lee oraz Carter (1992) w swojej pracy do estymacji parametrów modelu zaproponowali metodę odwołującą się do rozkładu wartości osobliwych macierzy (SVD). Minusem tej oryginalnie zastosowanej metody jest przyjęcie założenia o homoskedastyczności składnika losowego $\varepsilon_{x,t}$. Założenie to jednak nie ma potwierdzenia w badaniach empirycznych – natężenie zgonów w starszych wiekowo grupach podlega większym zmianom, niż to ma miejsce w grupach młodych osób (Brouhns, Denuit, Vermunt, 2002). W literaturze można znaleźć różne sposoby na rozwiązanie tego problemu (Brouhns, Denuit, Vermunt, 2002; Koissi, Shapiro, 2006). Jedną z proponowanych metod wykorzystywanych przy estymacji parametrów modelu Lee-Cartera jest ważona metoda najmniejszych kwadratów. W tej pracy zostanie jednakże zastosowana estymacja metodą największej wiarygodności. Metoda ta opiera się na założeniu, że liczba zgonów osób w wieku x i roku kalendarzowym t posiada rozkład Poissona, tzn.

$$D_{x,t} \sim \text{Poisson}(\hat{m}_{x,t} E_{x,t}),$$

gdzie:

$$\hat{m}_{x,t} = \exp(\hat{\alpha}_x + \hat{\beta}_x \hat{\kappa}_t),$$

$E_{x,t}$ – centralna liczba osób narażana na ryzyko zgonu.

Istotą tej metody jest maksymalizowanie funkcji wiarygodności (dokładnie logarytmu funkcji wiarygodności). Maksymalizowana funkcja w tym przypadku ma postać:

$$\ln L = \sum_x \sum_t [D_{x,t} \ln(m_{x,t} E_{x,t}) - E_{x,t} \exp(\alpha_x + \beta_x \kappa_t) - \ln(D_{x,t}!)] .$$

W równaniu występuje iloczyn szacowanych składników β , κ , przez co nie jesteśmy w stanie ich oszacować standardowymi metodami. Jednym z alternatywnych sposobów rozwiązania tego problemu jest zastosowanie podejścia iteracyjnego. W każdym kroku iteracyjnym szacowany jest jeden z parametrów (przy niezmiennych wartościach pozostałych – wcześniej oszacowanych – parametrów) zgodnie z ogólną zasadą:

$$\hat{\theta}^{v+1} = \hat{\theta}^v - \frac{\partial L^v / \partial \theta}{\partial^2 L^v / \partial \theta^2} .$$

Dla konkretnych parametrów wyrażenie to przyjmuje postać:

- $\hat{\alpha}^{v+1} = \hat{\alpha}^v + \frac{\sum_t D_{xt} - \hat{D}_{xt}^x}{\sum_t \hat{D}_{xt}^x}$,
- $\hat{\kappa}^{v+1} = \hat{\kappa}^v + \frac{\sum_x D_{xt} - \hat{D}_{xt}^x \hat{\beta}_x^v}{\sum_t \hat{D}_{xt}^x (\hat{\beta}_x^v)^2} \hat{\kappa}^v$,
- $\hat{\beta}^{v+1} = \hat{\beta}^v + \frac{\sum_x D_{xt} - \hat{D}_{xt}^x \hat{\kappa}_x^v}{\sum_t \hat{D}_{xt}^x (\hat{\kappa}_x^v)^2}$.

Procedura jest stopowana w momencie, gdy przyrost wartości funkcji wiarygodności nie przewyższa jakiejś niskiej, z góry narzuconej wartości (przy szacowaniu parametrów dla polskich danych została przyjęta wartość 10^{-7}). Jednym z najważniejszych powodów, dla których model został stworzony i oszacowany, jest możliwość dokonania na jego podstawie prognoz. Do tego konieczna jest ekstrapolacja κ_t , parametru opisującego zmiany w umieralności w czasie. Załóżmy, że κ_t ma charakter błędzenia losowego z dryfem, tzn. $\kappa_t = \kappa_{t-1} + c + \xi_t$, gdzie ξ_t są niezależnymi składnikami losowymi o jednakowych rozkładach normalnych $N(0, \sigma^2)$. Estymacja parametrów c oraz σ^2 została przeprowadzona analogicznie jak w pracach Li, Lee, Tuljapurkara (2004) i Bijaka, Więckowskiej (2008). Estymatory powyższych parametrów mają postać:

$$\hat{c} = (\kappa_T - \kappa_1) / (T - 1)$$

oraz

$$\hat{\sigma}^2 = 1 / (T - 1) \sum_{t=2}^T (\kappa_t - \kappa_{t-1} - \hat{c})^2 .$$

Mając powyższe oszacowania, możemy dokonać ekstrapolacji parametru κ na chwilę $T + s$:

$$\kappa_{T+s} = \kappa_T + (\hat{c} + sc \cdot \eta) \cdot s + \hat{\sigma} \cdot \sum_{\tau=T+1}^{T+s} \xi_{\tau},$$

gdzie η ma rozkład normalny $N(0, 1)$,

$$sc \approx \hat{\sigma} / \sqrt{(T - 1)}$$

jest zaś błędem oszacowania stałej c .

2.3. Model Renshaw-Habermana

Model ten jest rozszerzeniem klasycznego modelu Lee-Cartera. Zawiera dodatkowy składnik $\beta_x \gamma_{t-x}$, który ma za zadanie jeszcze lepiej uchwycić efekt kohortowy. Model ten można zapisać w następującej postaci:

$$\ln(\widehat{m}_{x,t}) = \alpha_x + \beta_x^0 \gamma_{t-x} + \beta_x^1 \kappa_t.$$

Podobnie jak w modelu Lee-Cartera zakłada się dodatkowo warunki identyfikujące model:

- $\sum_x \beta_x^0 = 1$,
- $\sum_x \beta_x^1 = 1$.

Za autorami modelu (Renshaw, Haberman, 2006) przyjmujemy dodatkowo, że estymatory powyższego modelu wyliczane są w następujący sposób:

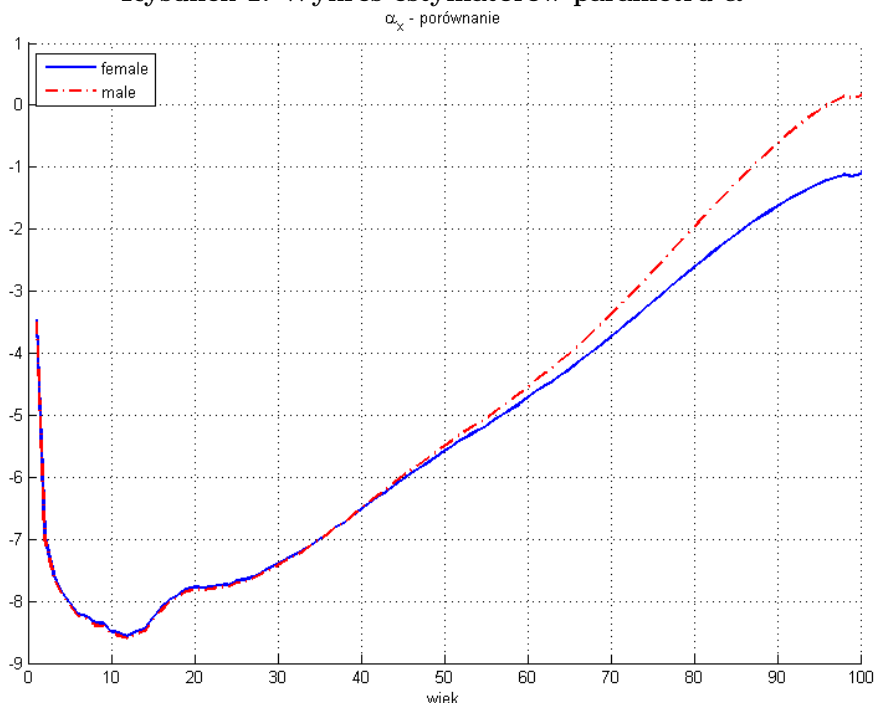
- $\hat{\gamma}_z^{v+1} = \hat{\gamma}_z^v + \frac{\sum_{z=t-x} D_{xt} - \hat{D}_{xt}^x \hat{\beta}_x^v}{\sum_{z=t-x} \hat{D}_{xt}^x (\hat{\beta}_x^v)^2}$,
- $\hat{\beta}^0^{v+1} = \hat{\beta}^0^v + \frac{\sum_x D_{xt} - \hat{D}_{xt}^x \hat{\gamma}_{t-x}^v}{\sum_t \hat{D}_{xt}^x (\hat{\gamma}_{t-x}^v)^2}$,
- $\hat{\kappa}^{v+1} = \hat{\kappa}^v + \frac{\sum_x D_{xt} - \hat{D}_{xt}^x \hat{\beta}_x^v}{\sum_t \hat{D}_{xt}^x (\hat{\beta}_x^v)^2}$,
- $\hat{\beta}^1^{v+1} = \hat{\beta}^1^v + \frac{\sum_x D_{xt} - \hat{D}_{xt}^x \hat{\kappa}_x^v}{\sum_t \hat{D}_{xt}^x (\hat{\kappa}_x^v)^2}$.

Przeprowadzone badania (Plat, 2009; Renshaw, Haberman, 2006) pokazują, że parametry kohortowe są o wiele istotniejsze dla danych dotyczących umieralności w starszych grupach wiekowych.

3. Wyniki oszacowań oraz prognozy

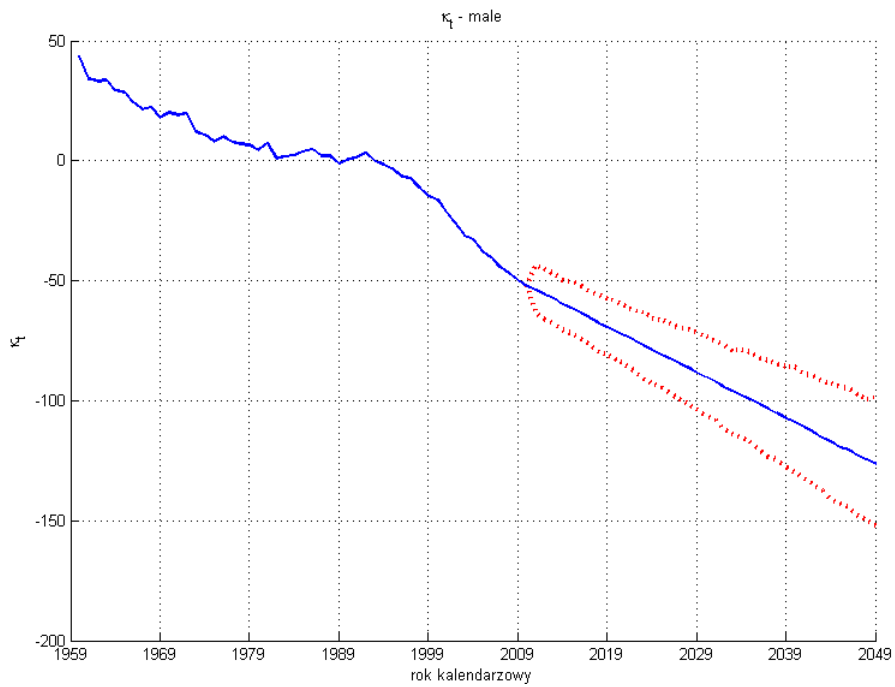
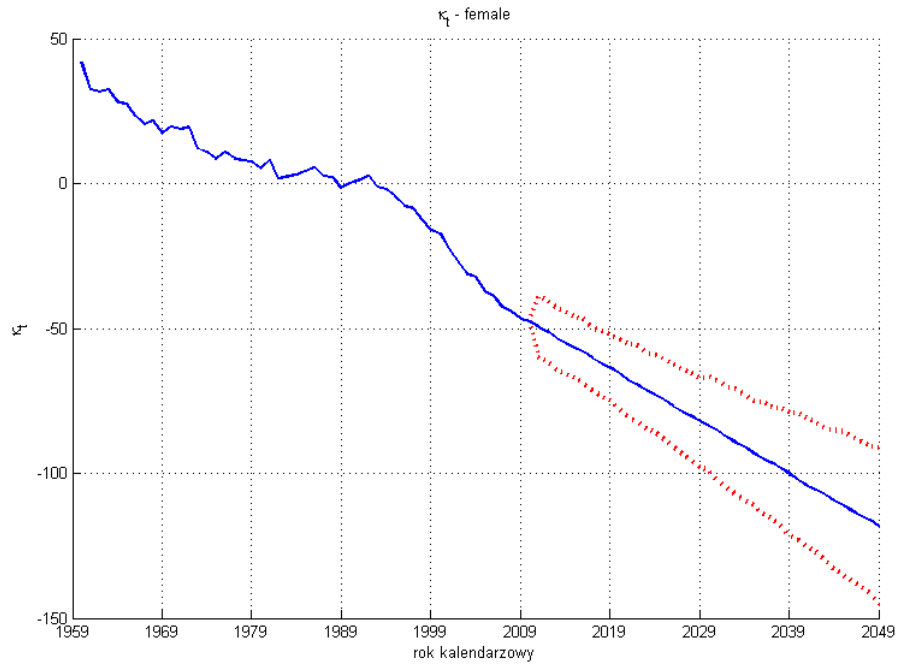
Na rysunku 1 przedstawione są estymatory parametru α_x dla kobiet i mężczyzn.

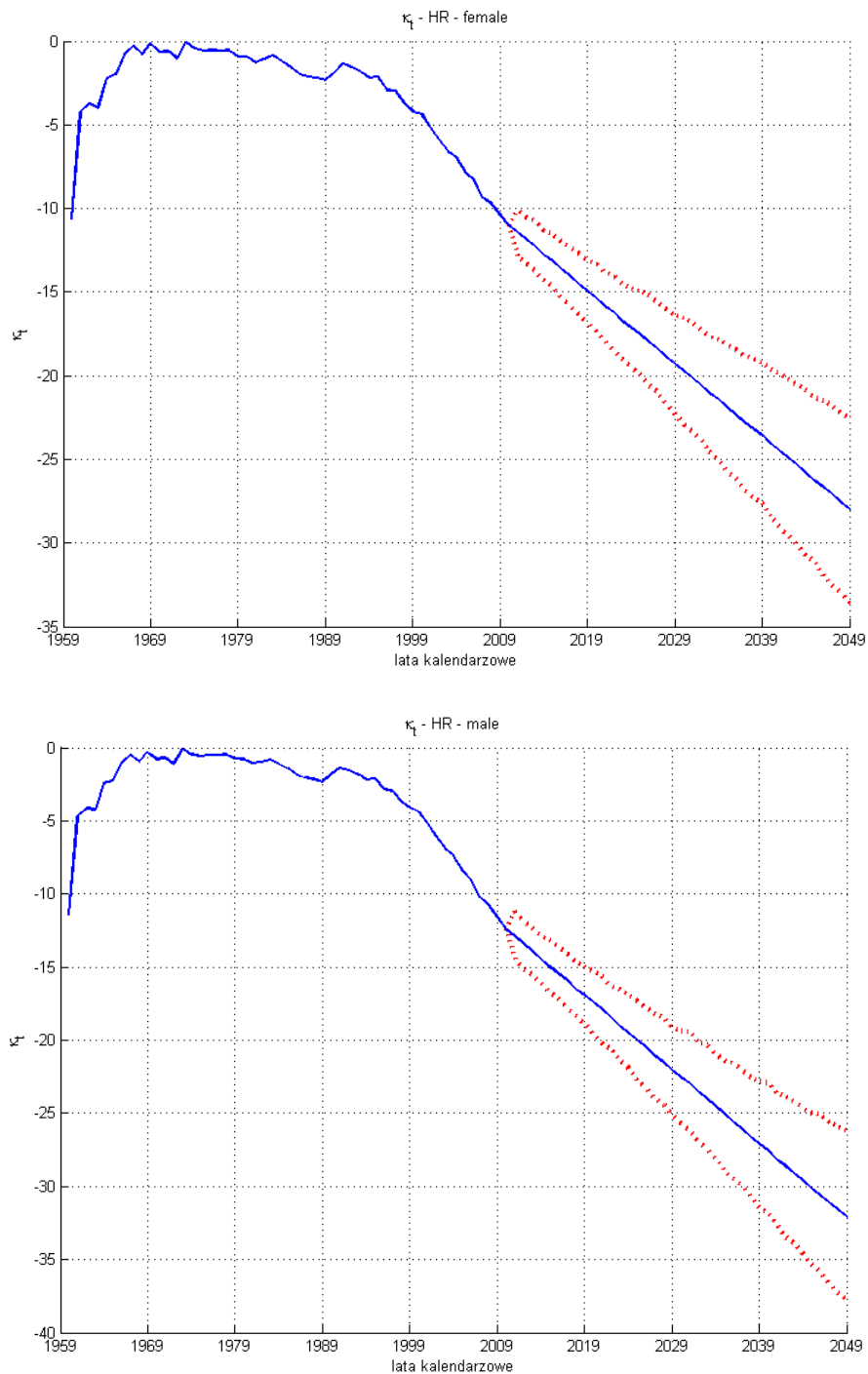
Rysunek 1. Wykres estymatorów parametru α

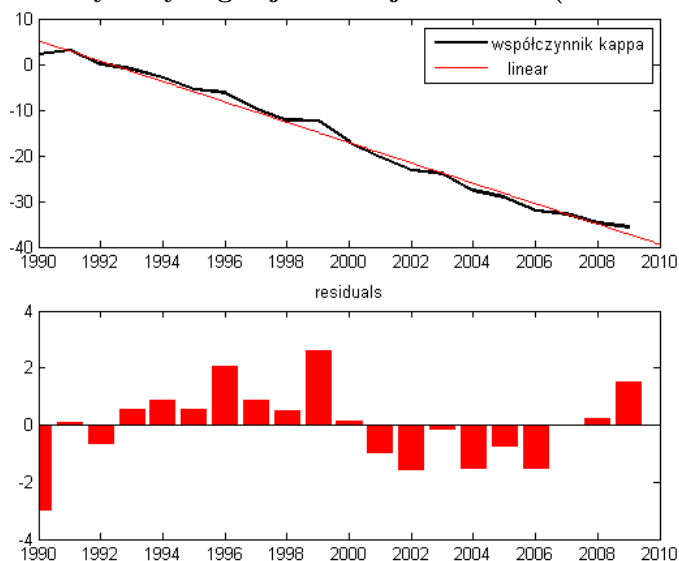


Nietrudno zauważyć w młodych wiekowo grupach (od 0 do 2 lat) wyższe wartości współczynnika α . Związane jest to z wysoką śmiertelnością wśród noworodków. Wyraźnie widać także, że od 50 roku życia intensywność zgonów wśród mężczyzn, wraz z kolejnymi latami, jest coraz wyższa. Na rysunku 2 znajdują się wykresy estymatorów parametrów κ . Kształty obu wykresów są bardzo podobne, z widoczną malejącą tendencją. Parametr κ , tak jak to było powiedziane wcześniej, opisuje zmiany umieralności w czasie. Odnośnie do polskich danych możemy wnioskować, że umieralność w Polsce w kolejnych latach maleje. Jest to efekt bardzo dynamicznego postępu, jaki dokonuje się w dziedzinie medycyny.

Począwszy od 1990 r., krzywe mają wygląd zbliżony do funkcji prostej, co sugerowałoby wykorzystanie regresji liniowej do ekstrapolacji parametru κ . Na rysunku 4 znajduje się wynik zastosowania regresji liniowej do danych wyestymowanych w modelu Lee-Cartera. Wykres reszt wskazuje na możliwość występowania autokorelacji składnika losowego.

Rysunek 2. Wykres estymatorów parametrów κ w modelu LC

Rysunek 3. Wykres estymatorów parametrów κ w modelu HR

Rysunek 4. Wykresy regresji liniowej oraz reszt (κ z modelu LC)

Występowanie autokorelacji potwierdzają również przeprowadzone testy (rysunek 5), dlatego ekstrapolując parametr κ w rozważanych modelach, zakładamy, że jest to proces błędzenia losowego z dryfem.

Rysunek 5. Wynik testu LM na autokorelację rzędu 1

Model 1: Estymacja KMNK, wykorzystane obserwacje 1990–2010 (N = 21)
Zmienna zależna (Y): ogon

	współczynnik	błąd standardowy	t-Studenta	wartość p
const	94,3086	3,98240	23,68	1,45e-015 ***
time	-2,80520	0,0960893	-29,19	3,01e-017 ***

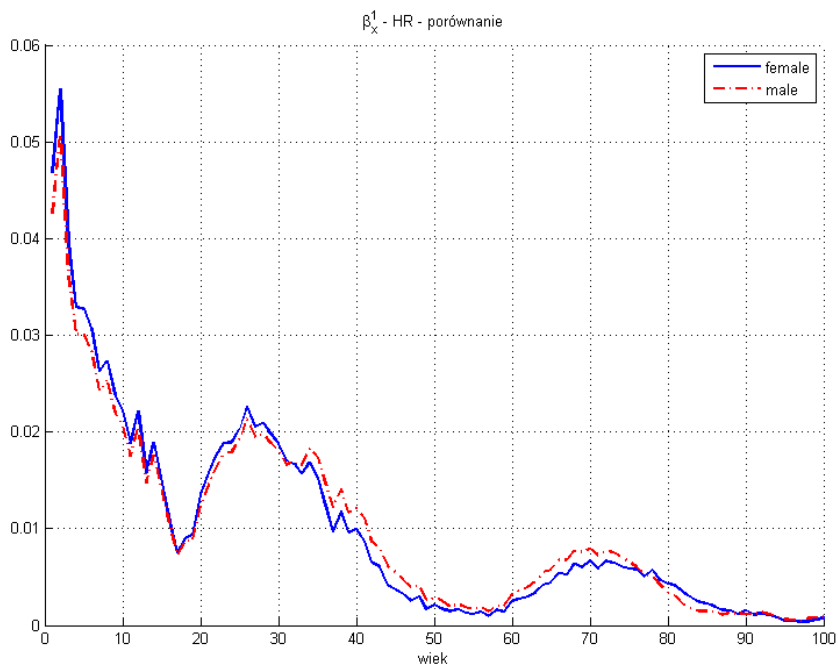
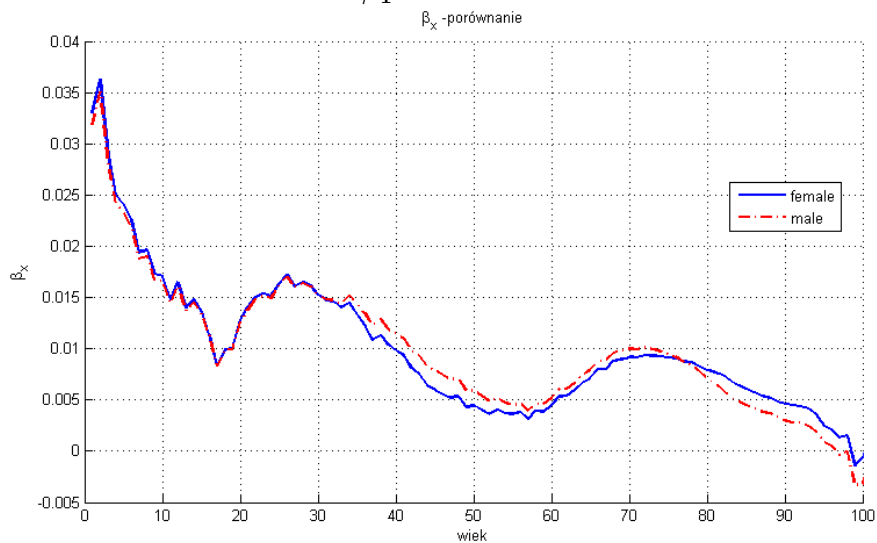
Średn. aryt. zm. zależnej	-20,70445	Odch. stand. zm. zależnej	17,59873
Suma kwadratów reszt	135,0811	Błąd standardowy reszt	2,666370
Wsp. determ. R-kwadrat	0,978193	Skorygowany R-kwadrat	0,977045
F(1, 19)	852,2682	Wartość p dla testu F	3,01e-17
Logarytm wiarygodności	-49,34191	Kryt. inform. Akaike'a	102,6838
Kryt. bayes. Schwarza	104,7729	Kryt. Hannana-Quinna	103,1372
Autokorel. reszt - rho1	0,532403	Stat. Durbin-Watsona	0,554977

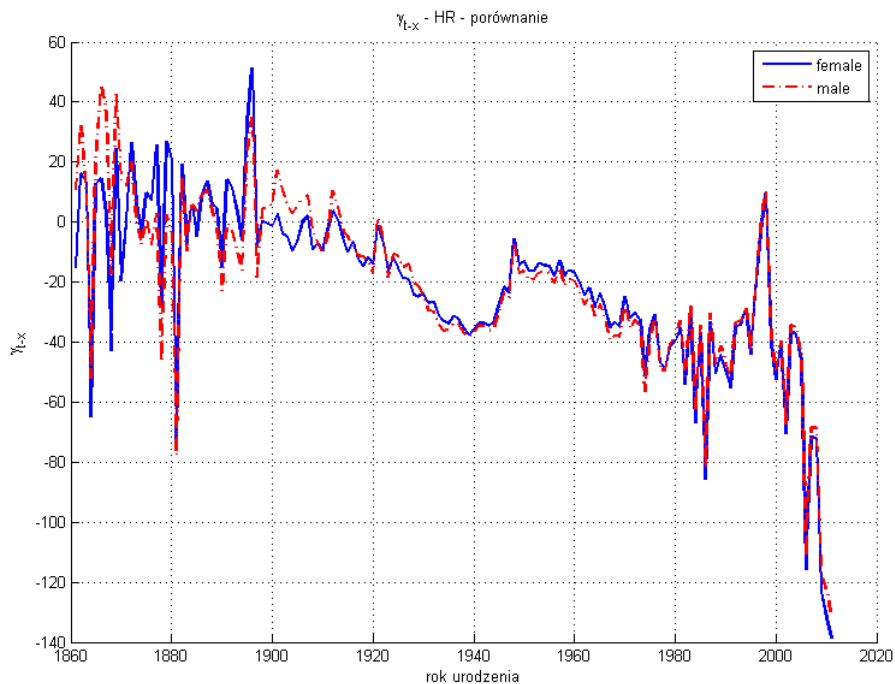
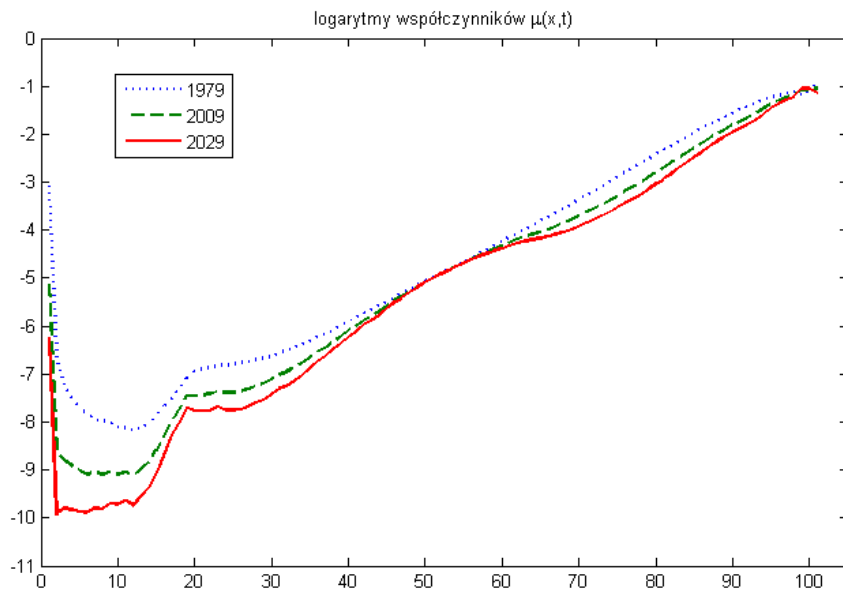
Test LM na autokorelację rzędu 1 -
Hipoteza zerowa: brak autokorelacji składnika losowego
Statystyka testu: LMF = 18,7553
z wartością p = P(F(1,48) > 18,7553) = 7,53145e-005

Na rysunkach 6 i 7 znajdują się wykresy estymatorów pozostałych parametrów: β (w modelu LC), β_1 (w modelu HR) oraz γ_{t-x} (w modelu HR). Interesujące jest zachowanie się parametru γ . Duże wahania wartości tego współczynnika na początku i na końcu wykresu wynikają z tego, że nie dysponujemy pełną informacją o zgonach w skrajnych kohortach. Na przykład, nie wiemy, jak będzie wyglądał rozkład zgonów w kohorcie osób urodzonych w 2005 r., bo te osoby nadal żyją. Na wykresie można również zauważyć osobliwe (zmiana tendencji z malejącą

na rosnącą) zachowanie się parametru γ w okresie 1940–1950, co zapewne jest efektem zmian w polskiej populacji, jakie nastąpiły w czasie II wojny światowej oraz po jej zakończeniu.

Rysunek 6. Wykres estymatorów parametrów β w modelu LC oraz β_1 w modelu HR



Rysunek 7. Wykres estymatora parametru γ_{t-x} w modelu HRRysunek 8. Wykres współczynników umieralności $\mu_{x,t}$ dla kobiet i mężczyzn z Polski (łącznie)

Do porównania oszacowanych modeli wykorzystano bayesowskie kryterium informacyjne oraz kryterium informacyjne Akaikiego:

$$BIC = -\frac{2L\hat{\theta}}{N} + \frac{K \ln N}{N},$$

$$AIC = -\frac{2L\hat{\theta}}{N} + \frac{2K}{N}.$$

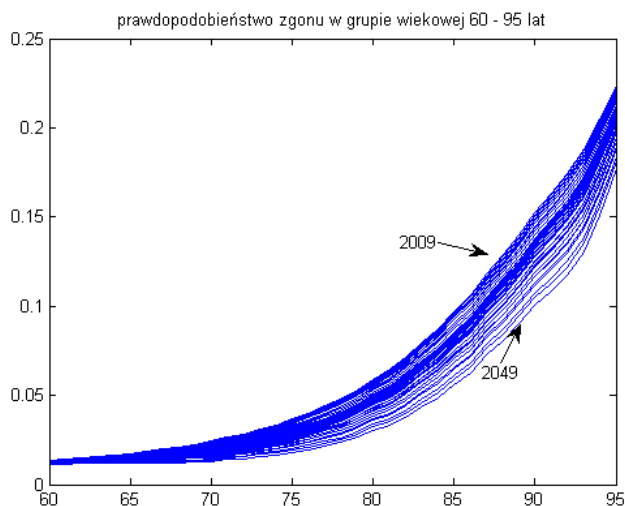
Otrzymane wyniki:

- LC: m – $BIC = AIC = 11786$, f – $BIC = AIC = 13440$,
 - HR: m – $BIC = AIC = 11776$, f – $BIC = AIC = 13434$,
- sugerują wybór modelu HR jako lepszego w stosunku do modelu LC.

Oszacowane parametry zostały wykorzystane do stworzenia prognoz współczynników umieralności oraz przeciętnego dalszego trwania życia (przekrojowego oraz kohortowego). Przypuszczamy, że podobnie jak w innych krajach wysoko rozwiniętych tak i w Polsce będzie następował spadek umieralności. Zmiany te są wyraźnie widoczne na rysunku 8, który zawiera wykresy logarytmów $\mu_{x,t}$.

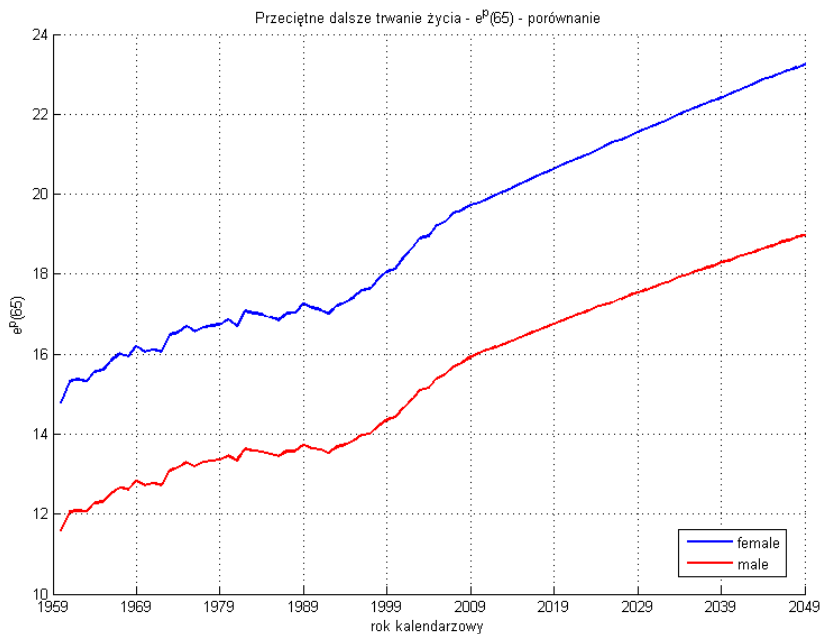
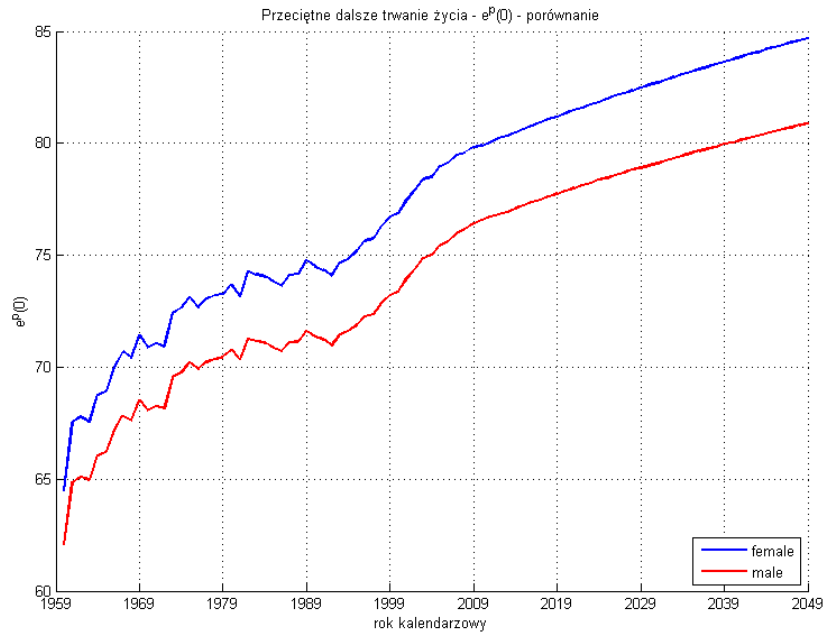
Największy spadek umieralności jest widoczny w najmłodszej grupie wiekowej (osoby do 20 roku życia). Fakt ten można wyjaśnić postępowaniem, który ma miejsce w medycynie. Coraz więcej chorób jest diagnozowanych już u noworodków, co umożliwia wczesne zastosowanie odpowiedniej terapii leczniczej. Spadek, choć nie tak znaczący, można również zaobserwować w grupie osób po 60 roku życia. Czynniki medyczne także w tym przypadku ma decydujące znaczenie. Powyższe wnioski potwierdza również wykres prognozowanych prawdopodobieństw zgonów.

Rysunek 9. Wykres prawdopodobieństw zgonów w grupie wiekowej 60–95 lat dla kobiet i mężczyzn z Polski (łącznie)

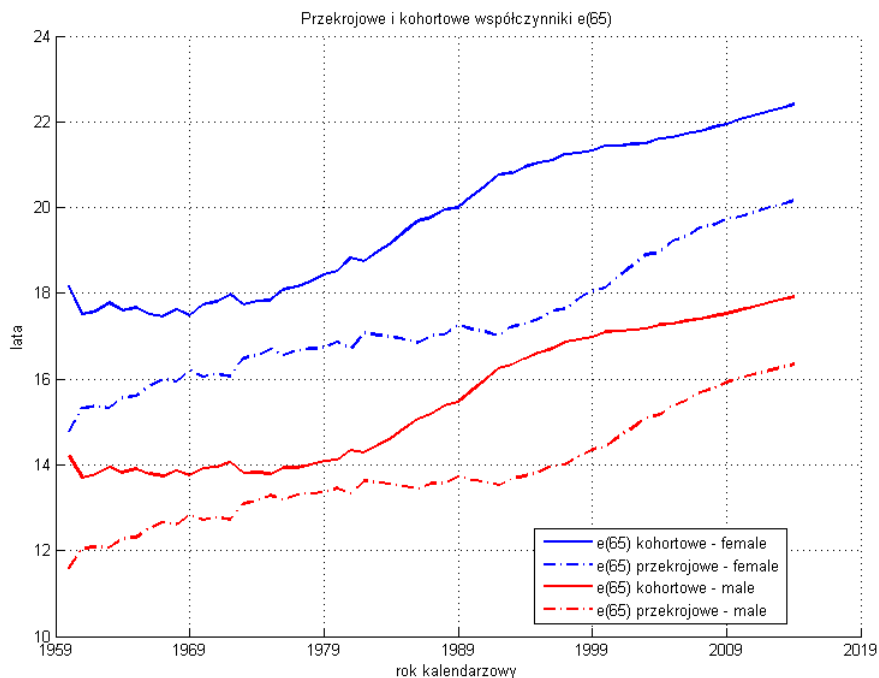


Na podstawie oszacowanego modelu można wyliczyć przeciętne dalsze trwanie życia osoby żyjącej w Polsce.

Rysunek 10. Przekrojowe dalsze trwanie życia noworodka i osoby w wieku 65 lat



Rysunek 11. Kohortowe dalsze trwania życia osoby w wieku 65 lat dla kobiet i mężczyzn z Polski



Obliczenia wskazują na wyraźne wydłużanie się długości życia w Polsce. Przeciętne przekrojowe dalsze trwanie życia dla obu płci w 1960 r. wynosiło 62–65 lat, 50 lat później zaś, w 2010 r., wzrosło do 74–79 lat. Przeprowadzone analizy przewidują, że w 2050 r. przeciętna długość życia noworodka w Polsce będzie wynosić 80 lat. Szczegółowa prognoza jest zamieszczona na rysunku 10. Ostatni γ_{t-x} wykres zawiera informacje o kohortowym przeciętnym dalszym trwaniu życia osoby w wieku 65 lat. Informacja taka może być szczególnie istotna w nawiązaniu do aktualnych problemów systemu emerytalnego w Polsce. Z wykresu wyraźnie wynika, jak bardzo podejście przekrojowe (stosowane przez GUS) nie doszacowuje przeciętnego dalszego trwania życia.

Powyższe badania mają charakter wstępny. Kolejnym krokiem będzie próba ekstrapolacji parametru γ opisującego wpływ efektu kohortowego oraz wyznaczenie przedziałów ufności dla prognozowanych wartości współczynników $\mu_{x,t}$.

Bibliografia

- [1] Bijak J., Więckowska B. (2008), *Prognozowanie przeciętnego dalszego trwania życia na podstawie modelu Lee-Cartera – wybrane zagadnienia*, w: *Statystyka aktuarialna – teoria i praktyka*, red. W. Ostasiewicz, WUE we Wrocławiu, Wrocław, s. 9–27.

- [2] Bowers N. i in. (1997), *Actuarial Mathematics*, Society of Actuaries, USA Schaumbury.
- [3] Brouhns N., Denuit M., Vermunt J. (2002), *A Poisson log-bilinear regression approach to the construction of projected lifetables*, „Insurance: Mathematics and Economics”, vol. 31, s. 373–393.
- [4] Imhoff E. van (1990), *The exponential multidimensional demographic projection model*, „Mathematical Population Studies”, vol. 2, s. 171–182.
- [5] Keilman N., Pham D. (2006), *Prediction intervals for Lee-Carter-based mortality forecasts*, referat na Europejską Konferencję Ludnościową w Liverpoolu, czerwiec 2006.
- [6] Koissi M.C., Shapiro A., Högnäs G. (2006), *Evaluating and extending the Lee-Carter model for mortality forecasting: Bootstrap confidence interval*, „Insurance: Mathematics and Economics”, vol. 38, s. 1–20.
- [7] Lee R., Carter L. (1992), *Modeling and Forecasting U.S. Mortality*, „Journal of the American Statistical Association”, vol. 87, s. 659–671.
- [8] Li N., Lee R., Tuljapurkar S. (2004), *Using the Lee-Carter method to forecast mortality for populations with limited data*, „International Statistical Review”, vol. 72, s. 19–36.
- [9] Plat R. (2009), *On stochastic mortality modeling*, „Insurance: Mathematics and Economics”, vol. 45, s. 393–404.
- [10] Renshaw A.E., Haberman S. (2006), *A cohort-based extension to the Lee-Carter model for mortality reduction factors*, „Insurance: Mathematics and Economics”, vol. 38, s. 556–570.

Stochastic modeling the mortality in Poland

Abstract

In this work I'll consider different ways of modeling mortality in Poland. I'll present two approaches (popular in recent years): Lee-Carter model and Renshaw-Haberman model. With these models, we can capture the age-period effects and cohort effect of an aging population. The advantage is also the possibility of extrapolating the estimated components. This allows to create projections of Polish mortality. At the time of speech I'll make a comparison of these models for Polish data and I'll examine life expectancy. The study will make for both sexes.

Autor:

Kamil Jodź, Instytut Nauk Ekonomicznych i Społecznych, Uniwersytet Przyrodniczy we Wrocławiu, ul. C.K. Norwida 25, 50-375 Wrocław,
e-mail: kamil.jodz@up.wroc.pl