

# Przewidywanie liczby ludności w Polsce za pomocą modeli regresji liniowej na podstawie zmiennych demograficznych

Krzysztof Kulka  
272667@student.pwr.edu.pl  
MSiD Lab Wtorek 9.15 NP

18 maja 2024

# Spis treści

<b>1</b>	<b>Wstęp</b>	<b>3</b>
<b>2</b>	<b>Zbiór danych i jego analiza</b>	<b>3</b>
2.1	Opis zbioru danych . . . . .	3
2.2	Przykładowe dane . . . . .	3
2.3	Obróbka danych . . . . .	4
2.3.1	Pozyskanie danych . . . . .	4
2.4	Dane po obróbce . . . . .	4
2.5	Analiza danych . . . . .	5
2.6	Populacja w Polsce . . . . .	6
2.7	Imigracja do Polski . . . . .	7
2.8	Współczynnik dzietności . . . . .	8
2.9	Oczekiwana długość życia . . . . .	9
2.10	Urbanizacja . . . . .	10
2.11	Wskaźnik zmiany populacji . . . . .	11
2.12	Korelacja . . . . .	12
2.13	Eliminacja danych odstających . . . . .	12
<b>3</b>	<b>Dobór metryk oceny</b>	<b>13</b>
3.1	Mean Squared Error . . . . .	13
3.2	Mean Absolute Percentage Error . . . . .	13
3.3	R2 Score . . . . .	13
3.4	Analiza metryk . . . . .	13
<b>4</b>	<b>Analiza modeli</b>	<b>13</b>
4.1	Dla podziału 80/20 . . . . .	13
4.1.1	Regresja liniowa . . . . .	14
4.1.2	Regresja Ridge . . . . .	15
4.1.3	Regresja Lasso . . . . .	16
4.1.4	Regresja Elastic Net . . . . .	17
4.1.5	Regresja Bayesian Ridge . . . . .	18
4.1.6	Podsumowanie dla podziału 80/20 . . . . .	18
4.2	Dla podziału 60/40 . . . . .	19
4.2.1	Regresja liniowa . . . . .	19
4.2.2	Regresja Ridge . . . . .	20
4.2.3	Regresja Lasso . . . . .	20
4.2.4	Regresja Elastic Net . . . . .	21
4.2.5	Regresja Bayesian Ridge . . . . .	22
4.2.6	Podsumowanie dla podziału 60/40 . . . . .	22
<b>5</b>	<b>Wnioski</b>	<b>22</b>

# 1 Wstęp

Problemem projektu jest analiza możliwości przewidywania liczby ludności w Polsce na podstawie zmiennych demograficznych, przy użyciu modeli regresji liniowej. Przedstawiona analiza ma na celu rozstrzygnięcie czy model regresji liniowej jest odpowiedni do przewidywania liczby ludności w Polsce, oraz jakie modele sprawdzają się do tego najlepiej. Analizie zostaną poddane następujące czynniki:

- Historyczna liczba ludności
- Imigracja do kraju
- Wskaźnik dzietności
- Oczekiwana długość życia
- Urbanizacja
- Wskaźnik zmiany populacji na przestrzeni ostatnich 5 lat

Zbadane zaś zostaną następujące modele regresji liniowej:

- Regresja liniowa
- Regresja typu Ridge
- Regresja Lasso
- Regresja Elastic Net
- Regresja Bayesian Ridge

## 2 Zbiór danych i jego analiza

### 2.1 Opis zbioru danych

Zbiór danych zawiera informacje na temat historycznej liczby ludności, imigracji do kraju, wskaźniku dzietności i oczekiwanej długości życia w momencie urodzenia na przestrzeni lat 1960-2023.

Dane zostały pobrane z serwisu internetowego World Bank[1] Dane dotyczą około 260 krajów. Dodatkowo informacje na temat urbanizacji zostały pobrane z serwisu internetowego Zintegrowana Platforma Edukacyjna Ministerstwa Edukacji Narodowej[2], a wskaźnik zmiany populacji na przestrzeni ostatnich 5 lat został obliczony na podstawie danych historycznych.

### 2.2 Przykładowe dane

"Country Name"	"Country Code"	"Indicator Name"	"Indicator Code"	"1960"	"1961"	...
"Aruba"	"ABW"	"Fertility rate, total (births per woman)"	"SP.DYN.TFRT.IN"	"4.82"	"4.655"	...
"Africa Eastern and Southern"	"AFE"	"Fertility rate, total (births per woman)"	"SP.DYN.TFRT.IN"	"6.72412501084242"	"6.74275210020318"	...
"Afghanistan"	"AFG"	"Fertility rate, total (births per woman)"	"SP.DYN.TFRT.IN"	"7.282"	"7.284"	...
"Africa Western and Central"	"AFW"	"Fertility rate, total (births per woman)"	"SP.DYN.TFRT.IN"	"6.45844789624312"	"6.47151755185967"	...
...						

Tabela 1: Wycinek danych ze zbioru danych nt. dzietności na kobiety

"Country Name"	"Country Code"	"Indicator Name"	"Indicator Code"	"1960"	"1961"	...
"Aruba"	"ABW"	"Population, total"	"SP.POP.TOTL"	"54608"	"55811"	...
"Africa Eastern and Southern"	"AFE"	"Population, total"	"SP.POP.TOTL"	"130692579"	"134169237"	...
"Afghanistan"	"AFG"	"Population, total"	"SP.POP.TOTL"	"8622466"	"8790140"	...
"Africa Western and Central"	"AFW"	"Population, total"	"SP.POP.TOTL"	"97256290"	"99314028"	...
...						

Tabela 2: Wycinek danych ze zbioru danych nt. liczby ludności

## 2.3 Obróbka danych

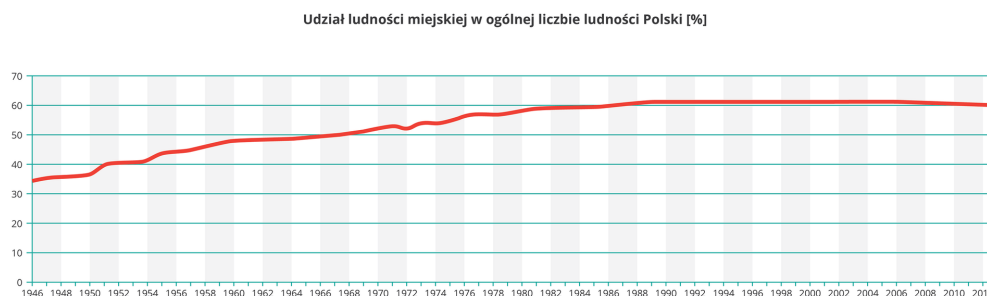
### 2.3.1 Pozyskanie danych

Najważniejszym krokiem w obróbce danych było wyizolowanie danych dotyczących Polski, oraz usunięcie kolumn, które nie były istotne dla analizy takich jak kod kraju, nazwa wskaźnika i kod wskaźnika.

Dodatkowo, z racji tego że dane dotyczące imigracji były rejestrowane jedynie co pięć lat, skorzystano z interpolacji liniowej, aby uzupełnić brakujące dane.

Największe wyzwanie pojawiło się z danymi dotyczącymi urbanizacji. Jedyne z zaufanego oficjalnego źródła były w formie obrazka .jpg. Aby więc uzyskać dane, i móc zamienić je w data frame, użyto następujących kroków:

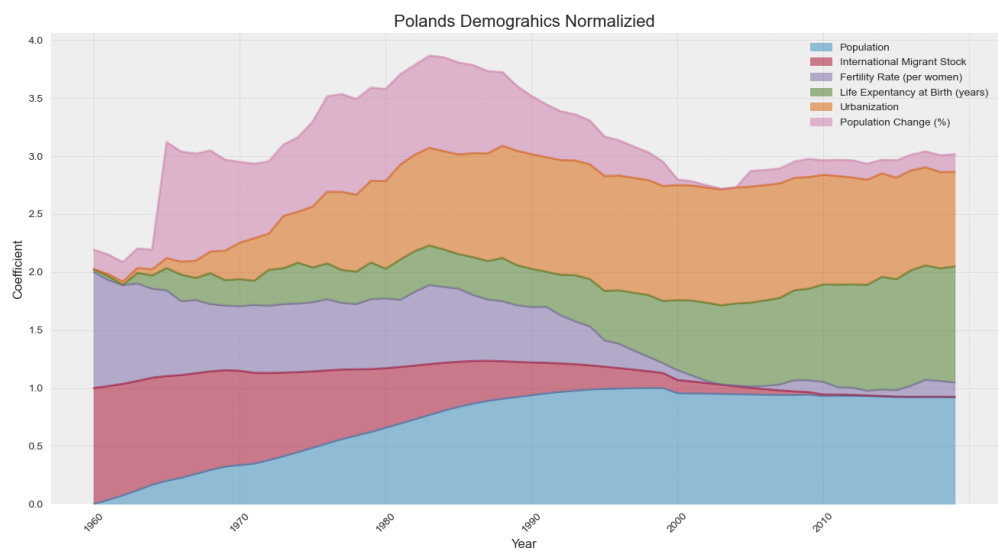
- Scraping obrazka za pomocą skryptu pythonowego `scrape_image.py`
- Następnie przekazanie zescrapowanego obrazka do zewnętrznego programu WebPlotDigitizer[3]
- Dalej, z racji niedoskonałości otrzymanego wyniku (m.in. potraktowanie lat jako liczb rzeczywistych, a nie całkowitych), dane zostały poprawione za pomocą kolejnego skryptu `python_fitter.py` napisanego w pythonie.
- Finalnie, ponieważ dane sięgały jedynie 2012 roku, dodano za pomocą skryptu `python_fitter.py` dane z lat 2013-2022 korzystając z witryny Gęografia24.pl[4].



Rysunek 1: Zescrapowany graf urbanizacji w Polsce

## 2.4 Dane po obróbce

Po wstępnej obróbce danych, tak prezentują się zmienne demograficzne po ich normalizacji:



Rysunek 2: Znormalizowane dane po obróbce

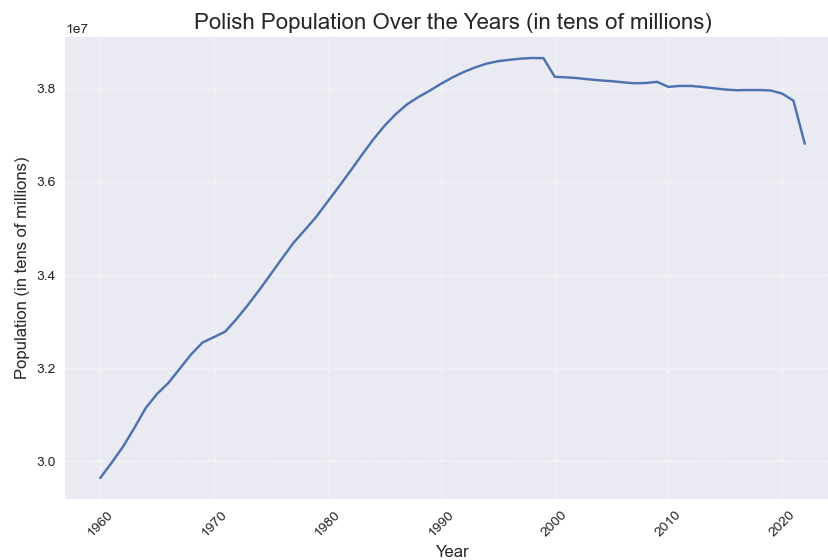
Year	Population	International Migrant Stock	Fertility Rate (per women)	Life Expentancy at Birth (years)	Urbanization	Population Change (%)
1987	37668045.0	1243055.6	2.15	70.8975609756098	60.20358735274085	3.967830291847352
1988	37824487.0	1204627.4	2.13	71.3317073170732	60.74798188743977	3.425258603566994
1989	37961529.0	1166199.2	2.08	71.0439024390244	60.99298968660394	2.865248104724527
1990	38110782.0	1127771.0	2.06	70.890243902439	61.02281597125872	2.443147706090709
1991	38246193.0	1095161.8	2.07	70.5878048780488	61.02661990321469	2.109332256232954
1992	38363667.0	1062552.6	1.95	71.090243902439	61.03232580114864	1.846716494046885

Tabela 3: Wycinek finalnych danych dla lat 1987-1992

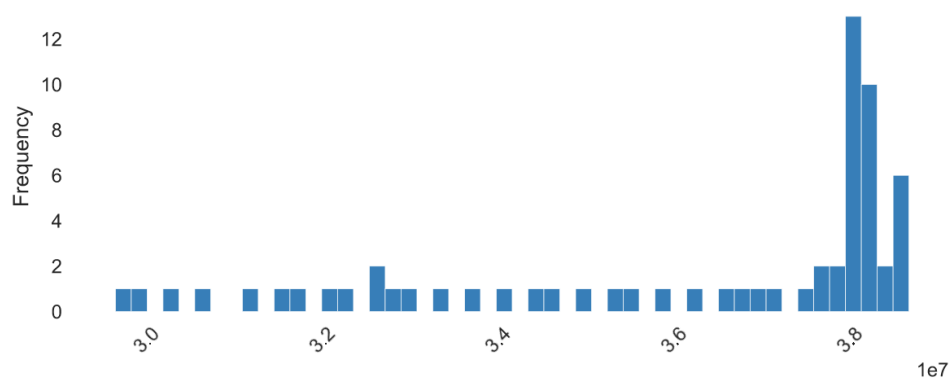
## 2.5 Analiza danych

Do analizy eksploracyjnej danych wykorzystano bibliotekę pandas-profiling[5]

## 2.6 Populacja w Polsce



Rysunek 3: Wizualizacja liczby ludności w Polsce na przestrzeni lat

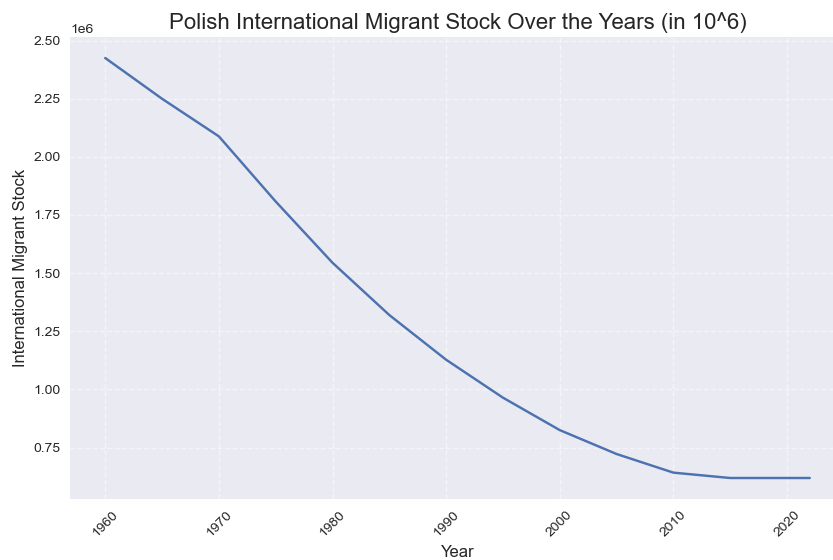


Rysunek 4: Histogram liczby ludności w Polsce

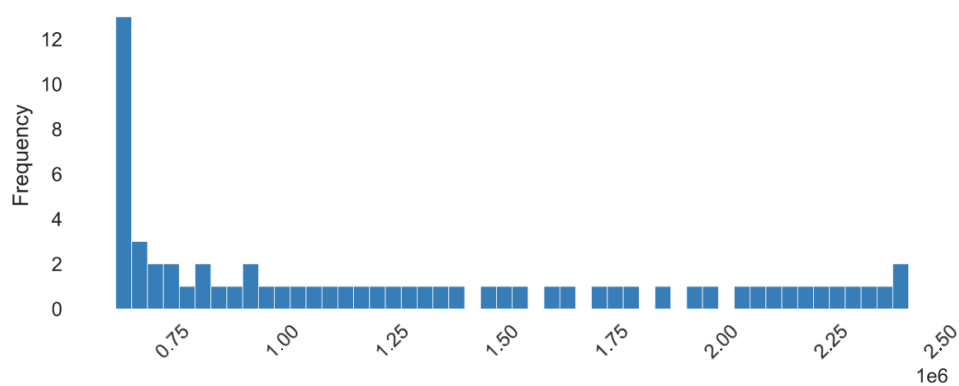
Minimum	Maximum	Mediana
29637450	38663481	37899070

Tabela 4: Statystyki liczby ludności w Polsce

## 2.7 Imigracja do Polski



Rysunek 5: Wizualizacja imigracji do Polski na przestrzeni lat

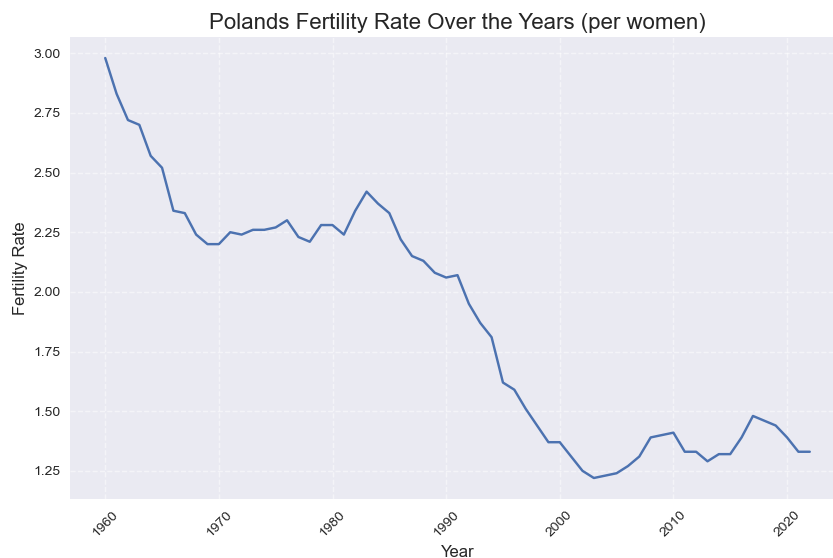


Rysunek 6: Histogram imigracji do Polski na przestrzeni lat

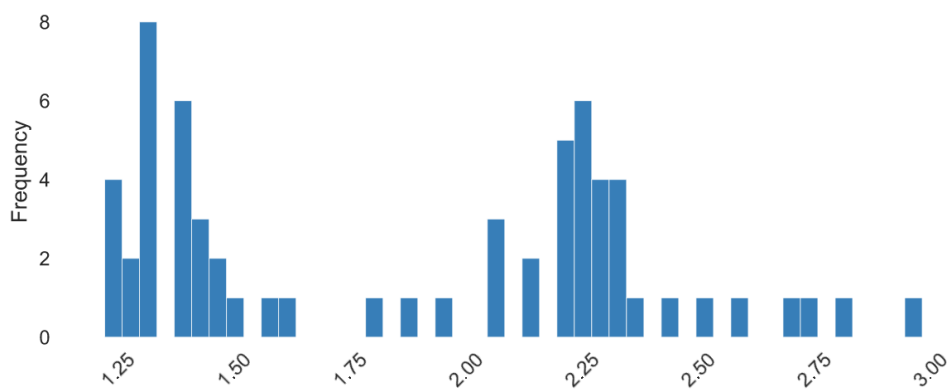
Minimum	Maximum	Mediana
619403	2424881	1095161

Tabela 5: Statystyki imigracji do Polski

## 2.8 Współczynnik dzietności



Rysunek 7: Wizualizacja współczynnika dzietności Polski na przestrzeni lat



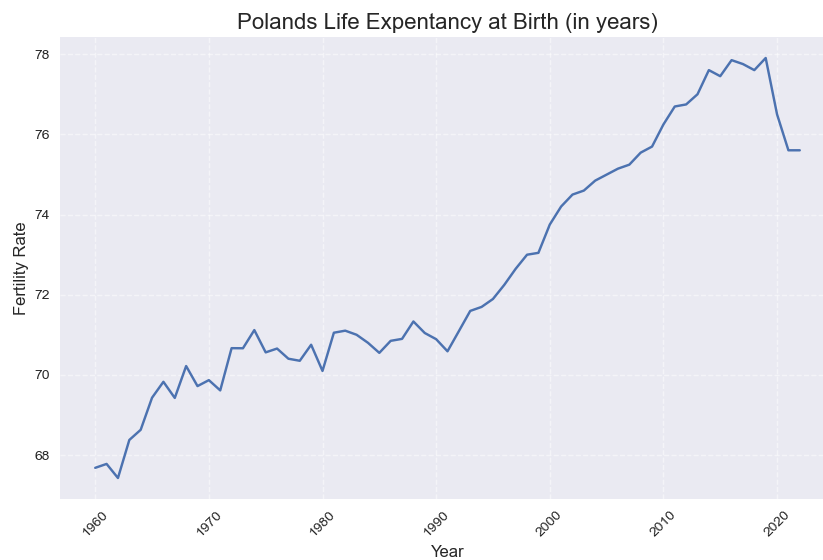
Rysunek 8: Histogram współczynnika dzietności na przestrzeni lat

Minimum	Maximum	Mediana
1.22	2.98	2.06

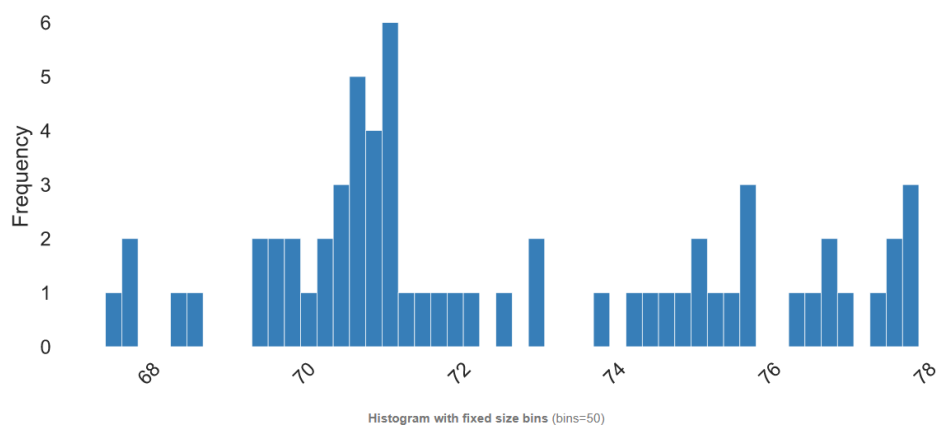
Tabela 6: Statystyki współczynnika dzietności



## 2.9 Oczekiwana długość życia



Rysunek 9: Wizualizacja oczekiwanej długości życia na przestrzeni lat

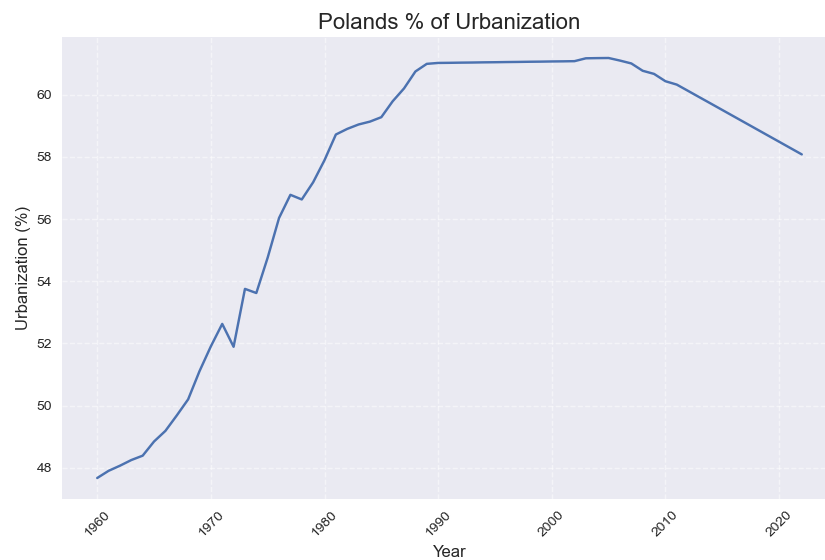


Rysunek 10: Histogram oczekiwanej długości życia na przestrzeni lat

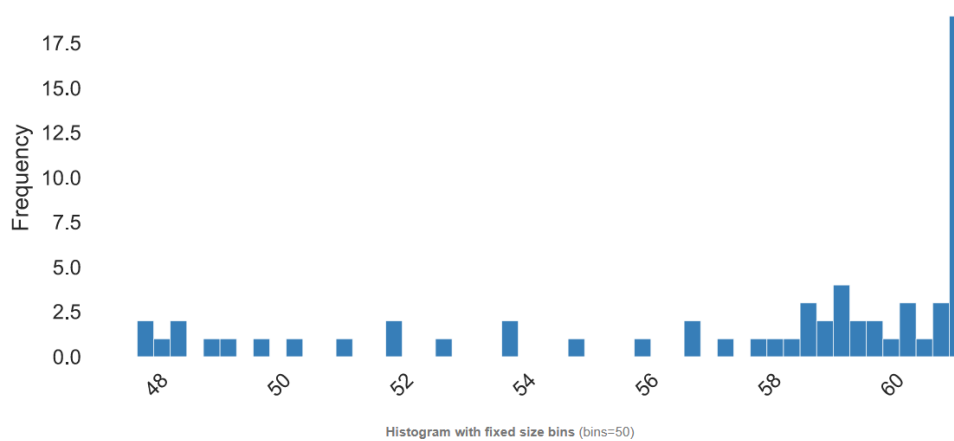
Minimum	Maximum	Mediana
67.42	77.9	71.1

Tabela 7: Statystyki oczekiwanej długości życia

## 2.10 Urbanizacja



Rysunek 11: Wizualizacja urbanizacji na przestrzeni lat



Rysunek 12: Wizualizacja urbanizacji na przestrzeni lat

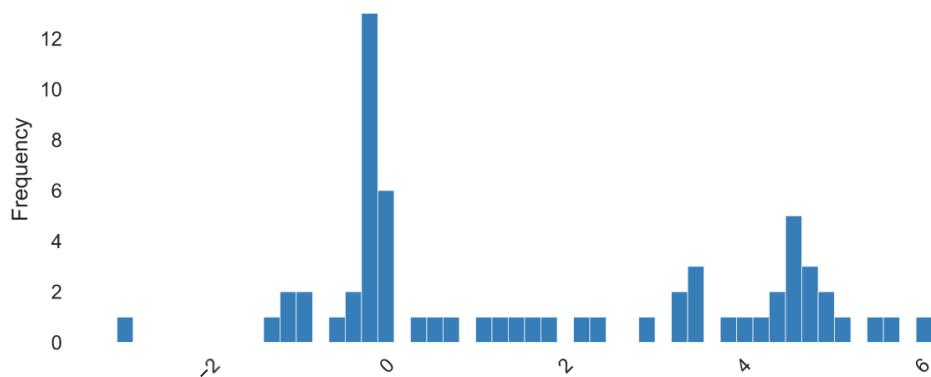
Minimum	Maximum	Mediana
47.66	61.18	59.27

Tabela 8: Statystyki urbanizacji

## 2.11 Wskaźnik zmiany populacji



Rysunek 13: Wizualizacja wskaźnika zmiany populacji na przestrzeni lat

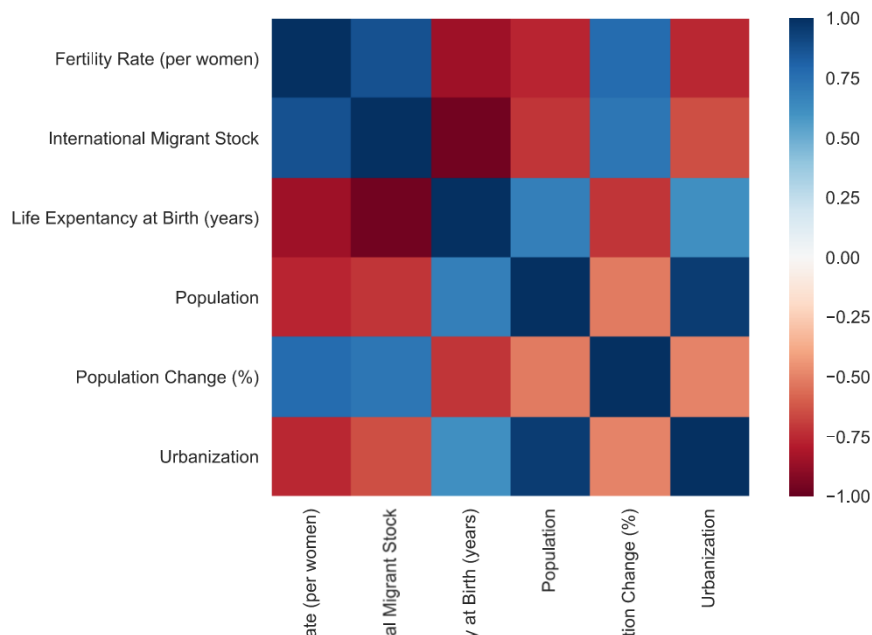


Rysunek 14: Histogram zmiany polskiej populacji na przestrzeni ostatnich 5 lat

Minimum	Maximum	Mediana
-3.03	6.09	0.98

Tabela 9: Statystyki wskaźnika zmiany populacji

## 2.12 Korelacja



Rysunek 15: Macierz korelacji

Z racji ręcznego doboru danych i ich selekcji, analiza wykazała bardzo dużą korelację (oraz anty-korelację) pomiędzy wybranymi współczynnikami. Pojawia się bardzo zaskakująca, przecząca logice korelacja, pomiędzy współczynnikiem dzietności a populacją wynosząca  $-0.763$ . Prawdopodobnie wynika ona z tego że przez większość badanego okresu, współczynnik był nadal na bardzo wysokim poziomie, więc mimo że malał, to populacja stale się zwiększała. Analogicznie zaskakuje negatywna korelacja migracji z populacją, co także dziwi, ponieważ zgodnie z intuicją, imigracja powinna zwiększać populację. Prawdopodobnie wynika to z faktu, że dane dotyczące imigracji są niewielkie, co powoduje że słabo, choć wciąż, przekładają się na polską populację. Jeśli zaś chodzi o spodziewane korelacje, należy szczególnie zwrócić uwagę:

- Oczekiwana długość życia z populacją: 0.683
- Urbanizacja z populacją: 0.949
- Urbanizacja z oczekiwaną długością życia: 0.611

Korelacje te dobrze wróżą dla modeli regresji, ponieważ są one na tyle silne, że powinny pozwolić na skuteczne przewidywanie liczby ludności w Polsce.

## 2.13 Eliminacja danych odstających

Podczas przeglądania grafów reprezentujących zebrane dane, nie trudno było zauważyć że dane po 2019 roku są znacznie odstające od reszty. Oczywiście jest że w 2020 roku, z racji pandemii, wiele wskaźników uległo zmianie, co sprawia że dane z tego roku są nieprzydatne do analizy. Z tego powodu, dane z 2020 roku wżwyz zostały usunięte.

### 3 Dobór metryk oceny

Zanim przystąpimy do analizy modeli, należy zdefiniować metryki, które pozwolą nam ocenić ich skuteczność. Dobór odpowiednich metryk jest kluczowy, ponieważ pozwala na obiektywną ocenę modeli, oraz porównanie ich ze sobą. W przypadku regresji, najczęściej stosowanymi metrykami są:

- Mean Squared Error (MSE) - średni błąd kwadratowy
- Mean Absolute Percentage Error (MAPE) - średni błąd procentowy bezwzględny
- R2 Score - współczynnik determinacji

#### 3.1 Mean Squared Error

MSE jest jedną z najczęściej stosowanych metryk w regresji. Oblicza ona średni błąd kwadratowy pomiędzy wartościami przewidywanymi przez model, a wartościami rzeczywistymi.

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (1)$$

#### 3.2 Mean Absolute Percentage Error

MAPE jest metryką, która mierzy średni błąd procentowy bezwzględny pomiędzy wartościami przewidywanymi przez model, a wartościami rzeczywistymi.

$$MAPE = \frac{1}{n} \sum_{i=1}^n \frac{|y_i - \hat{y}_i|}{y_i} \quad (2)$$

#### 3.3 R2 Score

R2 Score jest metryką, która mierzy jak dobrze model przewiduje dane w porównaniu do średniej wartości. Wartość R2 Score może przyjmować wartości od  $-\infty$  do 1, gdzie 1 oznacza idealne dopasowanie modelu.

$$R2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad (3)$$

#### 3.4 Analiza metryk

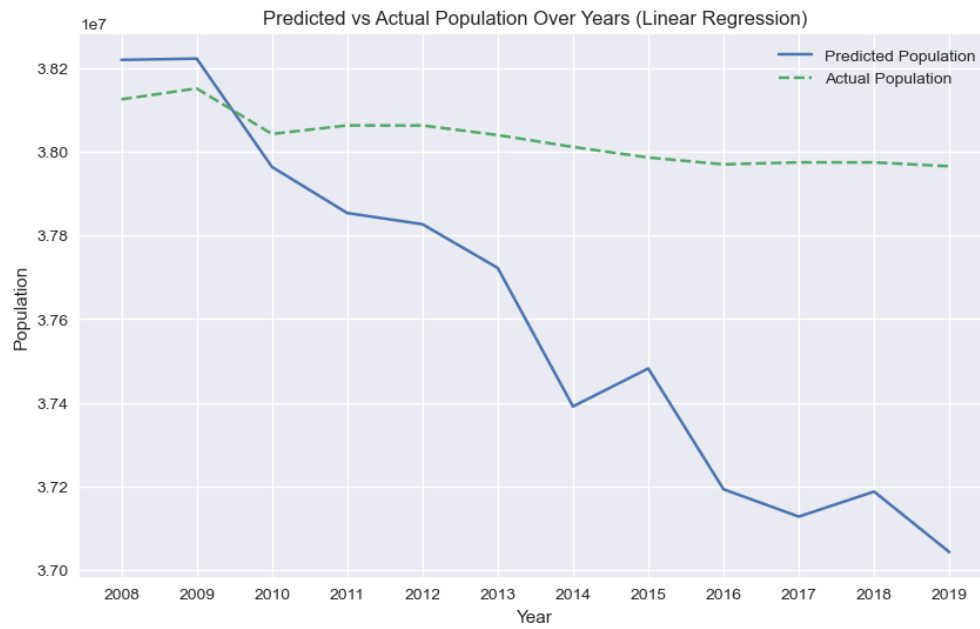
Metryki zostały wybrane tak, aby pokrywały one różne aspekty oceny modeli. MSE pozwala na ocenę jak dobrze model przewiduje wartości w skali, MAPE pozwala na ocenę jak dobrze model przewiduje wartości w procentach, a R2 Score pozwala na ocenę jak dobrze model przewiduje wartości w porównaniu do próby. Nie ulega jednak wątpliwości że najważniejszą metryką jest MAPE, która najlepiej nadaje się do predykcji dotyczących populacji[6], a więc jest najbardziej adekwatna do naszego problemu, i to przede wszystkim przez jej pryzmat będziemy oceniać modele.

### 4 Analiza modeli

#### 4.1 Dla podziału 80/20

Podział 80/20 oznacza to że model uczyć będzie się na danych z lat 1960-2007, a testowany będzie na danych z lat 2008-2019.

### 4.1.1 Regresja liniowa

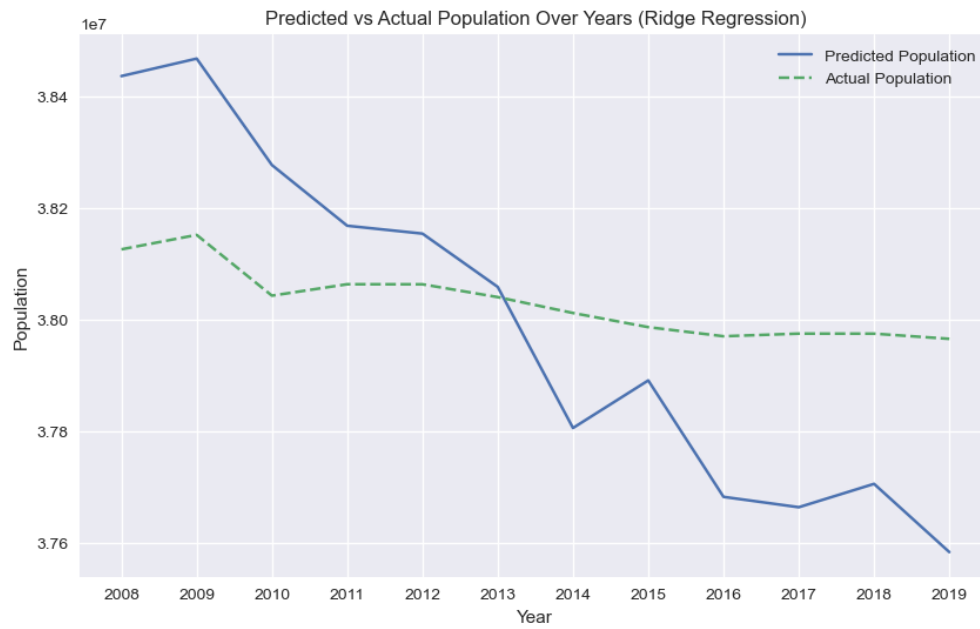


Rysunek 16: Regresja liniowa

Model	MSE	MAPE	R2 Score
Linear Regression	3.04e+11	0.012	-84.972

Tabela 10: Metryki dla regresji liniowej i podziału 80/20

### 4.1.2 Regresja Ridge

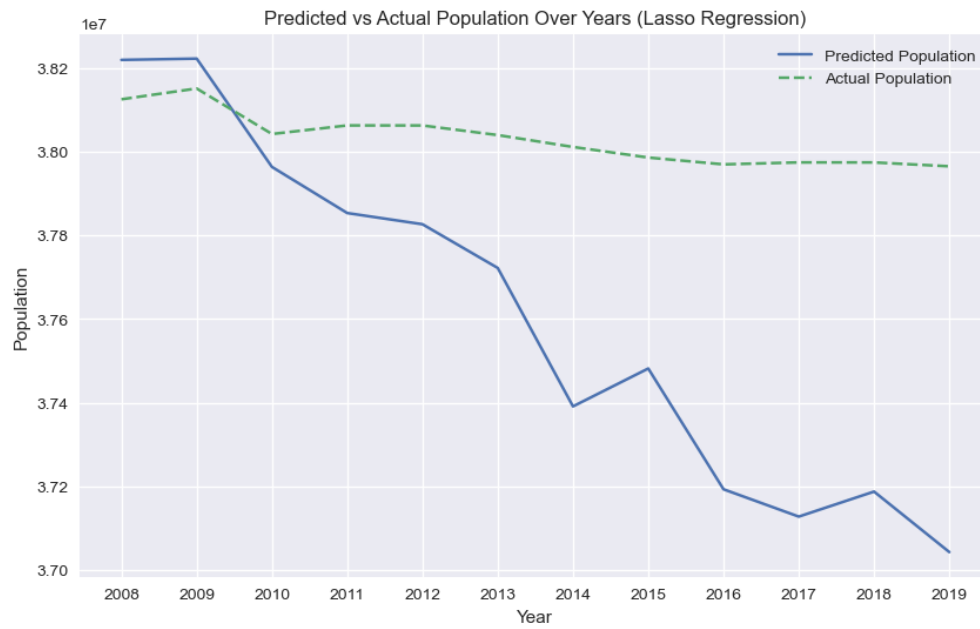


Rysunek 17: Regresja Ridge

Model	MSE	MAPE	R2 Score
Lasso Regression	6.00e+10	0.006	-15.926

Tabela 11: Metryki dla regresji Ridge i podziału 80/20

### 4.1.3 Regresja Lasso



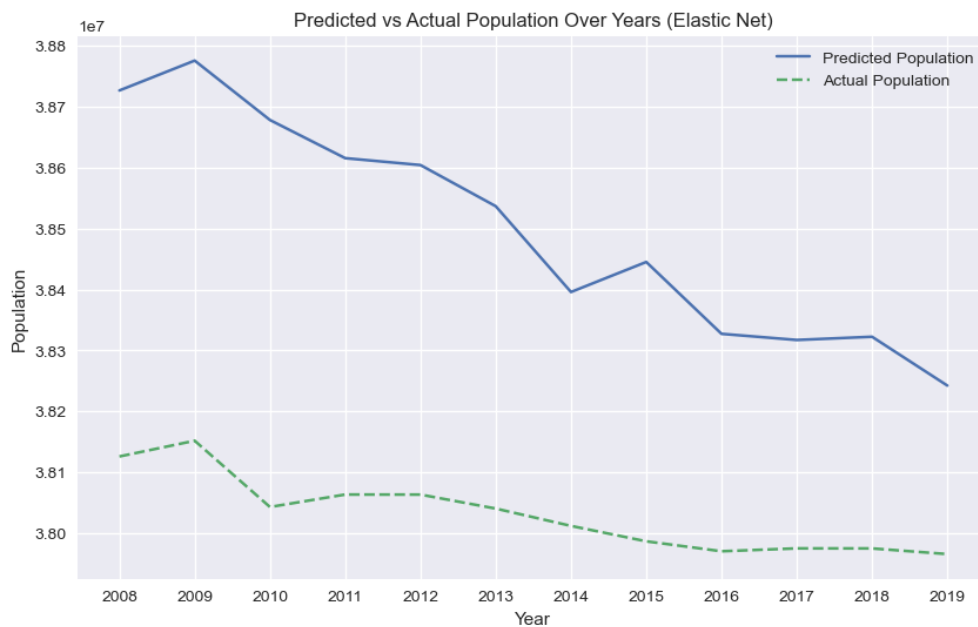
Rysunek 18: Regresja Lasso

Model	MSE	MAPE	R2 Score
Lasso Regression	3.04e+11	0.011	-84.952

Tabela 12: Metryki dla regresji Lasso i podziału 80/20



#### 4.1.4 Regresja Elastic Net

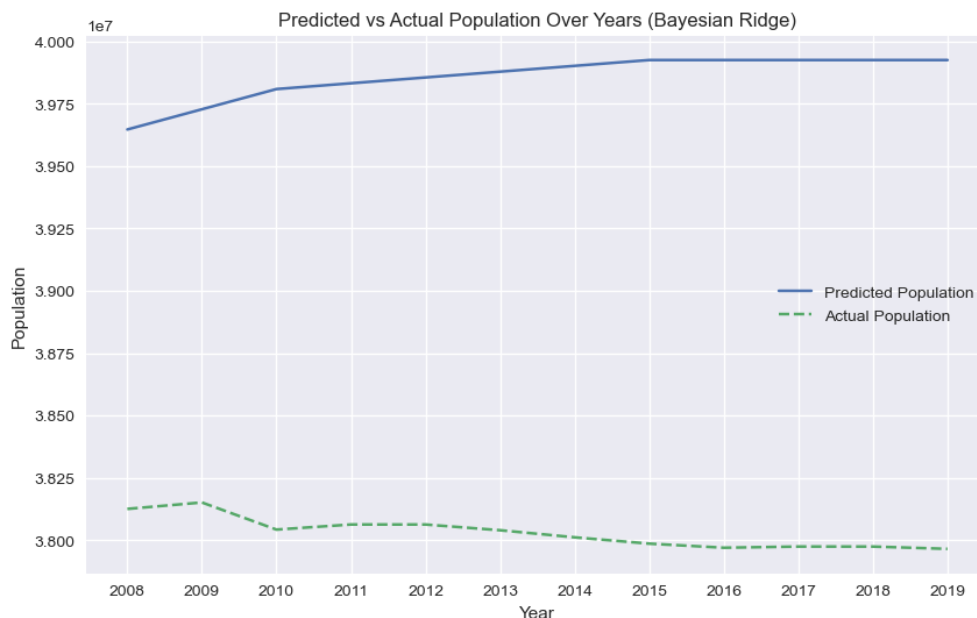


Rysunek 19: Regresja Elastic Net

Model	MSE	MAPE	R2 Score
Elastic Net	2.33e+11	0.012	-64.851

Tabela 13: Metryki dla regresji Elastic Net i podziału 80/20

### 4.1.5 Regresja Bayesian Ridge



Rysunek 20: Regresja Bayesian Ridge

Model	MSE	MAPE	R2 Score
Bayesian Ridge Regression	3.35e+12	0.048	-946.190

Tabela 14: Metryki dla regresji Bayesian Ridge i podziału 80/20

### 4.1.6 Podsumowanie dla podziału 80/20

Model	MSE	MAPE	R2 Score
Linear Regression	3.04e+11	0.012	-84.972
Ridge Regression	6.00e+10	0.006	-15.926
Lasso Regression	3.04e+11	0.011	-84.952
Elastic Net	2.33e+11	0.012	-64.851
Bayesian Ridge	3.35e+12	0.048	-946.190

Tabela 15: Wyniki dla podziału 80/20

Metryka MSE jest bardzo wysoka, co jest spowodowane tym że dane do przewidzenia są bardzo wysokie (rzędu  $10^7$ ), co sprawia że błędy są również bardzo wysokie. Dodatkowo, modele otrzymały niewielkie ilości danych szkoleniowych, co sprawia że są w pewnym stopniu niewystarczająco dopasowane do danych testowych. MAPE, która jest główną metryką, daje jednak nadzieje na to że model może być użyteczny, ponieważ wartość błędu procentowego dla Ridge wynosi zaledwie 0.006.

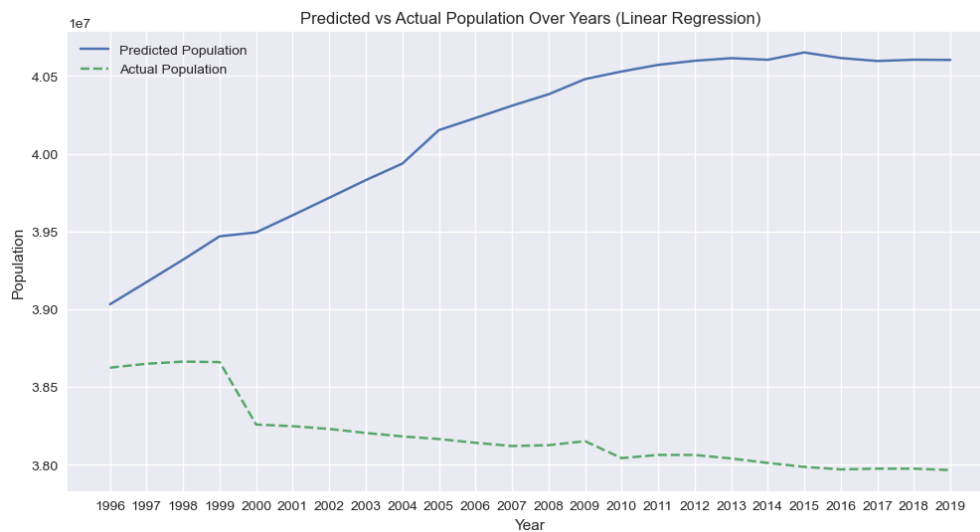
- Najlepszym modelem okazała się regresja Ridge, która osiągnęła najniższe wartości MSE, MAPE oraz najwyższe R2 Score.

- Najgorszym modelem okazała się regresja Bayesian Ridge, która osiągnęła najwyższe wartości MSE, MAPE oraz niespotykane niski R2 Score.
- Regresja Lasso oraz regresja liniowa osiągnęły bardzo zbliżone wyniki, co sprawia że nie można jednoznacznie stwierdzić która z nich jest lepsza.
- Regresja Elastic Net osiągnęła wyniki zbliżone, choć gorsze, od regresji Ridge, ale lepsze R2 i MSE niż Lasso i Liniowa.

## 4.2 Dla podziału 60/40

Podział 60/40 oznacza to że model uczyć będzie się na danych z lat 1960-1995, a testowany będzie na danych z lat 1996-2019.

### 4.2.1 Regresja liniowa

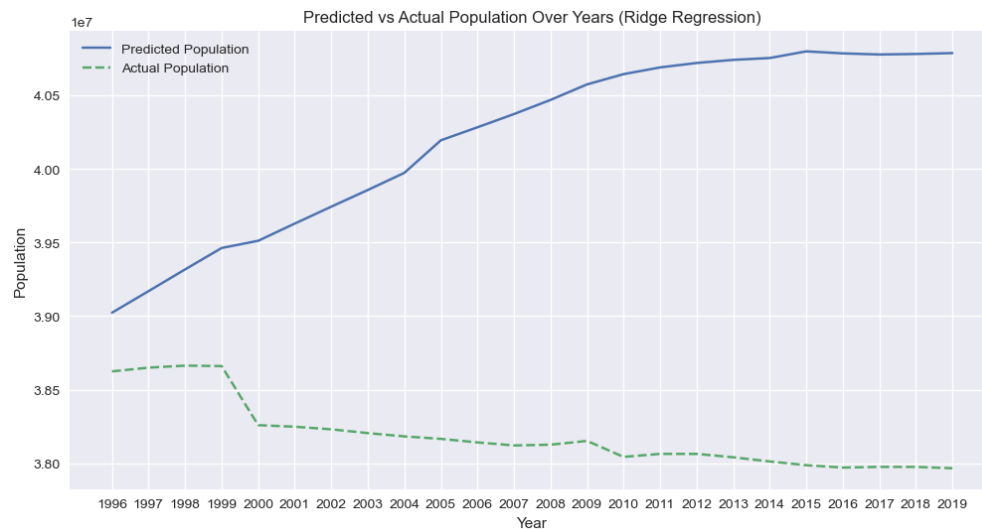


Rysunek 21: Regresja liniowa

Model	MSE	MAPE	R2 Score
Linear Regression	4.31e+12	0.051	-84.631

Tabela 16: Metryki dla regresji liniowej i podziału 60/40

## 4.2.2 Regresja Ridge

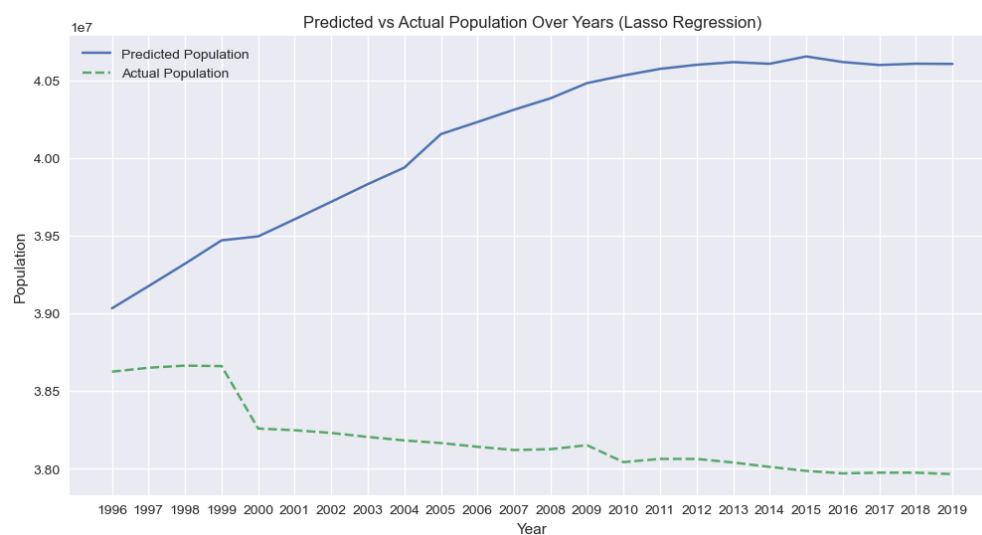


Rysunek 22: Regresja Ridge

Model	MSE	MAPE	R2 Score
Ridge Regression	4.73e+12	0.053	-92.887

Tabela 17: Metryki dla regresji Ridge i podziału 60/40

## 4.2.3 Regresja Lasso

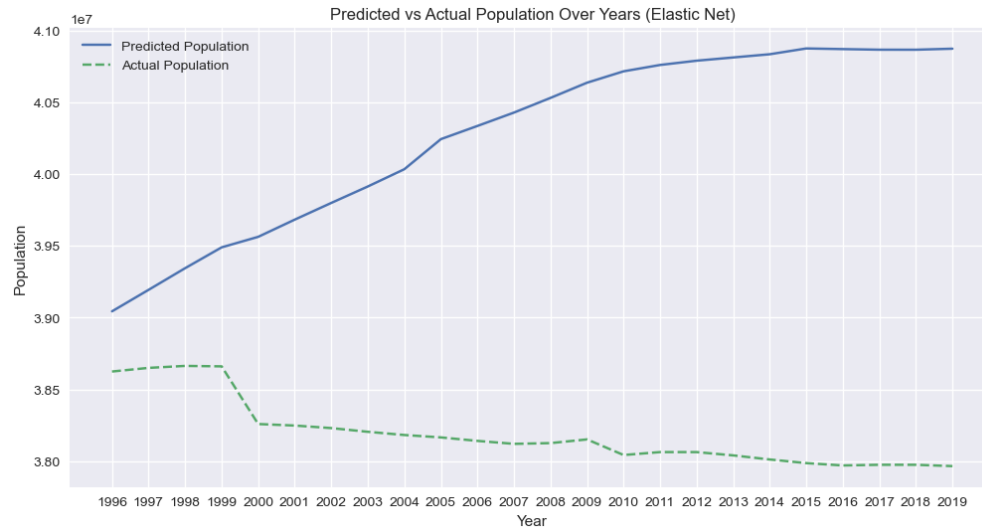


Rysunek 23: Regresja Lasso

Model	MSE	MAPE	R2 Score
Lasso Regression	4.31e+12	0.050	-84.638

Tabela 18: Metryki dla regresji Lasso i podziału 60/40

#### 4.2.4 Regresja Elastic Net

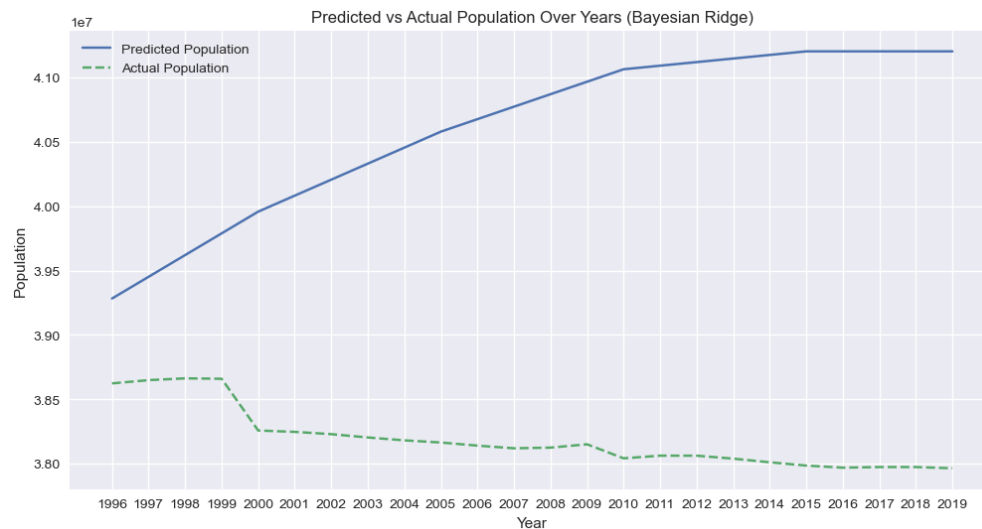


Rysunek 24: Regresja Elastic Net

Model	MSE	MAPE	R2 Score
Elastic Net	4.99e+12	0.054	-98.200

Tabela 19: Metryki dla regresji Elastic Net i podziału 60/40

## 4.2.5 Regresja Bayesian Ridge



Rysunek 25: Regresja Bayesian Ridge

Model	MSE	MAPE	R2 Score
Bayesian Ridge	6.55e+12	0.063	-129.007

Tabela 20: Metryki dla regresji Bayesian Ridge i podziału 60/40

## 4.2.6 Podsumowanie dla podziału 60/40

Model	MSE	MAPE	R2 Score
Linear Regression	4.31e+12	0.051	-84.631
Ridge Regression	4.73e+12	0.053	-92.887
Lasso Regression	4.31e+12	0.050	-84.638
Elastic Net	4.99e+12	0.054	-98.200
Bayesian Ridge	6.55e+12	0.063	-129.007

Tabela 21: Wyniki dla podziału 60/40

Jak widać, wyniki dla podziału 60/40 są znacznie gorsze niż dla podziału 80/20. Prawdopodobnie wynika to z faktu że model otrzymał mniej danych szkoleniowych, co powoduje niewystarczające dopasowanie do danych testowych. Split 80/20 okazał się zdecydowanie lepszy, co sprawia że jest on bardziej odpowiedni do analizy, i zapewnia dobry bilans między szkoleniem i testowaniem modelu.

## 5 Wnioski

Po przeanalizowaniu wielu różnych modeli regresji, na różnych podziałach danych, można stwierdzić że regresja Ridge jest najlepszym modelem do przewidywania liczby ludności w Polsce w zaprezentowanej sytuacji. Model ten osiągnął najniższe wartości MSE, MAPE oraz najwyższe R2 Score, co

sprawa że jest on najbardziej adekwatny do naszego problemu. Warto szczególnie zwrócić uwagę na wartość MAPE, która wynosi zaledwie 0.006, co oznacza że model przewiduje wartości zaledwie o 0.6% odchylające się od wartości rzeczywistych. Kategoryzuje to ten model w kategorię bardzo dobrych[7], co sprawia że jest on najbardziej odpowiedni do przewidywania liczby ludności w Polsce.

Dodatkowo, projekt ten udowadnia że zmienne demograficzne, takie jak dzietność, oczekiwana długość życia, urbanizacja, oraz imigracja, mają duży wpływ na liczebność populacji, do tego stopnia, że na ich podstawie można z dużą dokładnością przewidzieć liczebność populacji w przyszłości.

Oczywiście, jak pokazał przykład pandemii COVID-19, zdarzenia losowe mogą znacząco wpłynąć na liczebność populacji, co sprawia że model ten nie jest idealny, ale nadal jest on bardzo dobrym narzędziem do przewidywania liczby ludności w Polsce, i może być użyty do analizy trendów demograficznych w przyszłości, zwłaszcza w sytuacjach stabilnych, gdzie nie występują znaczące zmiany w liczbie ludności.

Ciekawą i wartą uwagi obserwacją, bez wątpienia jest fakt, że najlepsze modele regresji liniowej, przewidywały populację niższą niż rzeczywista. W istocie może być to wytłumaczone faktem, że dane imigracyjne są prawdopodobnie zaniżone, co przełożyło się na fakt że modele przewidywały zaniżoną liczbę mieszkańców Polski.

## Literatura

- [1] World Bank Data poland. <https://www.worldbank.org/pl/country/poland>.
- [2] Zintegrowana Platforma Edukacyjna Ministerstwa Edukacji Narodowej urbanizacja w polsce. <https://zpe.gov.pl/a/zroznicowanie-poziomu-urbanizacji-w-polsce/D19MUchJD>.
- [3] WebPlotDigitizer. <https://automeris.io/>.
- [4] Geografia24, Urbanizacja w Polsce i w Europie. <https://geografia24.pl/urbanizacja-w-polsce-i-w-europie/>.
- [5] Fabian Clemente. YData-Profiling. <https://github.com/ydataai/ydata-profiling>.
- [6] Aybeniz S. Aliyev Makrufa Sh. Hajirahimov. Development of a prediction model on demographic indicators based on machine learning methods: Azerbaijan example.
- [7] Research Gate, Interpretation of typical MAPE values. [https://www.researchgate.net/figure/nterpretation-of-typical-MAPE-values\\_tbl1\\_257812432](https://www.researchgate.net/figure/nterpretation-of-typical-MAPE-values_tbl1_257812432).