

Zuzanna Słobodzian
nr albumu: 412204

Kraków, 10.06.2024



Akademia Górniczo-Hutnicza im. Stanisława Staszica w Krakowie

Sprawozdanie z projektu
zaliczeniowego
Data Mining:
dane CO2 Emissions

1. Wstęp

Celem projektu jest wykonanie dwóch modeli regresyjnych na danych dotyczących zmiany wskaźnika emisji CO₂ w różnych państwach świata na przestrzeni lat 1960–2014. Wybrane modele to kNN regresyjny oraz sieci neuronowe MLP, modele przewidują wartość wskaźnika w roku 2014, jest to ostatni rok, dla którego są dostępne dane. Projekt został wykonany w języku R, przy użyciu bibliotek: ggplot2, caret, neuralnet.

2. Opis danych

Oryginalnie dane mieszczą się w dwóch arkuszach, pierwszy zawiera nazwy i kody państw, nazwę i kod wskaźnika oraz wartości wskaźnika z lat 1960–2014. Drugi zawiera nazwy i kody państw, informacje o grupie dochodowej i regionie oraz dodatkowe uwagi.

Przed przystąpieniem do analizy, informacje z obu tabelek zostały połączone w jedną, tak aby uzyskać komplet informacji o nazwie państwa, przypisanym regionie, grupie dochodowej i wysokości wskaźnika na przestrzeni lat (rys. 2.1).

	A	B	C	D	E	F	G	H	I	J	K
4											
5	Country Name	Country Code	Indicator Name	Indicator Code	Region	Income Group	1960	1961	1962	1963	1964
6	Aruba	ABW	CO2 emissions (metric tons per capita)	EN.ATM.CO2E.PC	Latin America & Caribbean	High income					
7	Afghanistan	AFG	CO2 emissions (metric tons per capita)	EN.ATM.CO2E.PC	South Asia	Low income	0.04606	0.053604	0.073765	0.074233	0.086292
8	Angola	AGO	CO2 emissions (metric tons per capita)	EN.ATM.CO2E.PC	Sub-Saharan Africa	Lower middle income	0.097472	0.079038	0.201289	0.192535	0.201003

rys. 2.1

Kolumna „Country Name” zawiera nazwy państw oraz zbiorczych regionów jak np. „*East Asia & Pacific (excluding high income)*”, a kolumna „Country Code” trzy literowe unikalne id dla państw i zbiorczych regionów. Wartości „Indicator Name” i „Indicator Code” dla każdego wiersza są takie same: „*CO2 emissions (metric tons per capita)*” oraz „*EN.ATM.CO2E.PC*” – cały arkusz danych zawiera informacje o tylko jednym wskaźniku. W kolumnie „Region” jest wyróżnionych siedem grup: „*South Asia*”, „*Sub-Saharan Africa*”, „*Europe & Central Asia*”, „*Middle East & North Africa*”, „*Latin America & Caribbean*”, „*East Asia & Pacific*”, „*North America*”. Kolumna „Group Income” dzieli dane na cztery grupy: „*Low income*”, „*Lower middle income*”, „*Upper middle income*”, „*High income*”. Wartości w kolumnach zatytułowanych numerami kolejnych lat od 1960 do 2014 zawierają dane numeryczne oznaczające wysokość wskaźnika emisji CO₂, liczonego jako tona CO₂ na osobę.

3. Preprocessing i analiza danych

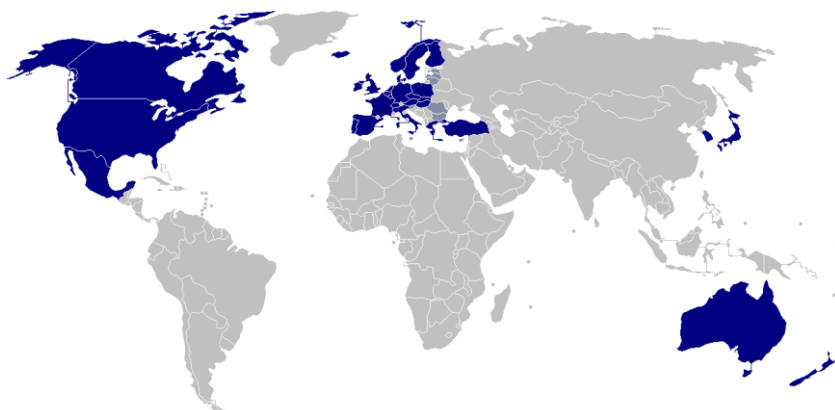
Po wczytaniu danych usunięto trzy górne puste wiersze, nadano odpowiednie nagłówki i usunięto wiersze z brakami danych.

Zdecydowano się na usunięcie całych wierszy zawierających braki danych, ponieważ układ występowania wartości NA uniemożliwiał zastąpienie ich konkretnymi wartościami. W przypadku kolumn z wartościami wskaźnika były to albo wiersze całkowicie puste, albo takie w których dłuższe ciągi braków wartości występowały obok siebie, zwykle w początkowych latach zbierania danych (rys. 3.1),.

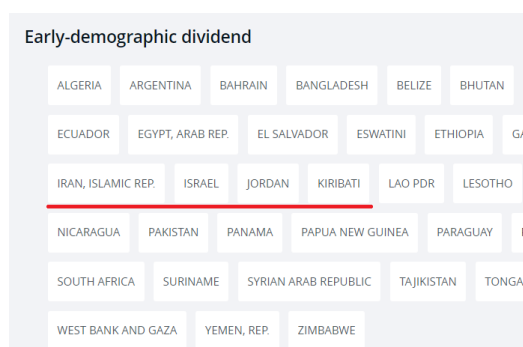
8.94799	0.54843	0.84433	0.74307	0.67747	0.72902	0.70383	0.84901	1.00036		1.0347	1.06737	1.13807	1.1889	1.08116	0.82388	1.28106	1.36748	1.35107	1.25263	1.06181	1.25673	1.04329	0.8325	0.81413	0.84205	0.66051	0.78086	0.87042	1.11727	0.91242	1.00475	0.99664	0.87701	0.98354	1.04708
6.20788	6.84435	0.86296	7.8237	8.87477	7.30642	7.91138	8.02038	4.30287	8.81156	8.67734	8.66668	9.26231	8.23775	9.27005	8.78946	8.26164	9.59822	9.11381	8.18035	7.8214	7.52549	1.77782	7.08828	6.74945	8.56707	4.43052	8.71418	6.42087	6.67189	6.87189	6.08925	8.89823	8.89412	6.20895	
3.26895	3.48131	3.73055	3.67364	3.72342	4.49027	4.79073	5.40744	4.16732	5.58135	5.80288	6.50565	7.42028	6.43139	7.8954	6.62971	8.39276	11.1454	8.09125	10.0105	10.7454	10.10109	11.0525	11.3455	10.9549	10.8072	11.2037	12.2603	14.8477	14.8477	14.8477	14.8477	14.8477	14.8477	14.8477	
0.17283	0.14208	0.24855	0.40629	0.76291	1.52748	1.80894	2.48153	2.52526	4.72261	4.9677	6.55678	9.58436	9.26282	8.03050	9.10291	10.8838	9.82434	10.9191	9.10427	9.24898	7.01238	7.43247	7.81155	6.24489	4.5876	4.8054	6.22892	4.72522	4.52709	4.90918	4.83836	4.14617	4.25737	4.28587	
2.2457	2.1316	2.2068	2.0602	2.2867	2.04	2.3477	2.9492	1.187	2.1119	2.51093	2.8887	2.61888	2.87945	2.88817	2.57531	2.31865	3.01791	2.80724	1.60342	0.7742	1.28862	1.88513	1.46676	1.34439	1.60371	2.38438	2.38101	4.90375	6.77865	7.98829	8.04828	8.44802	9.1228	9.29363	
0.26326	0.23588	0.21963	0.22113	0.18833	0.18658	0.17455	0.20131	0.30541	0.26055	0.26584	0.26479	0.30583	0.27937	0.24292	0.29671	0.26888	0.27188	0.26955	0.27458	0.26484	0.30735	0.26909	0.26135	0.23254	0.24372	0.24807	0.23495	0.26505	0.25182	0.25753	0.28004	0.30337	0.31464	0.32647	
2.26507	2.21306	1.83961	2.86052	2.29407	2.3834	2.34405	4.43851	2.79489	2.51093	2.8887	2.61888	2.87945	2.88817	2.57531	2.31865	3.01791	2.80724	1.60342	0.7742	1.28862	1.88513	1.46676	1.34439	1.60371	2.38438	2.38101	4.90375	6.77865	7.98829	8.04828	8.44802	9.1228	9.29363		
0.17885	0.18582	0.18375	0.17967	0.18201	0.18246	0.17965	0.18771	0.18983	0.19387	0.18571	0.19324	0.18324	0.20253	0.20509	0.20404	0.21227	0.21375	0.21572	0.20962	0.1982	0.19807	0.19544	0.19072	0.18466	0.18484	0.18111	0.16632	0.16186	0.1558	0.15204	0.15575	0.15389	0.15774		
0.06950	0.0734	0.07435	0.07318	0.06885	0.07027	0.06818	0.10939	0.10658	0.11053	0.11069	0.12726	0.12332	0.13001	0.13458	0.206	0.2454	0.25067	0.2699	0.2532	0.24891	0.23965	0.24583	0.23542	0.21948	0.22345	0.21748	0.20392	0.19189	0.18924	0.1834	0.18547	0.18837	0.18572		
0.06403	0.07489	0.08111	0.10115	0.09964	0.07078	0.09094	0.08482	0.0789	0.10086	0.08902	0.14546	0.13353	0.14603	0.12158	0.13951	0.13022	0.13882	0.05315	0.17508	0.17175	0.16744	0.184	0.1918	0.19894	0.20596	0.22883	0.2337	0.17029	0.16971	0.16887	0.16843	0.16465	0.16124		
0.08379	0.09628	0.08456	0.10575	0.09011	0.08748	0.09735	0.14389	0.11964	0.11964	0.2188	0.30079	0.10439	0.19055	0.24	0.24719	0.25272	0.28498	0.2523	0.22979	0.25799	0.25236	0.19823	0.23862	0.18391	0.21454	0.25378	0.26432	0.28335	0.15464	0.13457	0.14644	0.21038	0.06775		
1.1032	1.44957	1.55376	1.89125	2.16521	2.24841	2.55367	2.80383	2.74833	3.1496	3.55235	4.17251	4.04138	4.2971	4.67028	5.04915	5.08882	5.73208	5.33364	5.21329	5.20897	5.64902	5.76556	5.85349	6.54863	6.77393	7.35254	7.28452	7.17684	7.26026	7.23801	7.41079	7.45888	7.58385		
0.27773	0.15713	0.2259	0.22263	0.38921	0.38116	0.42519	0.50348	0.46602	0.50814	0.58749	0.57447	0.78986	0.51585	0.52335	0.5296	0.57938	0.58284	0.5358	0.64778	0.66966	0.64947	0.63329	0.63286	0.65912	0.72786	0.67226	1.06965	1.10448	1.14854	1.20727	1.23946	1.43969	1.45965		
7.54354	9.06578	8.38729	8.79332	9.58758	9.97704	12.585	9.0541	8.21914	9.44441	11.4641	11.3752	11.2693	10.2025	4.40896	14.6235	7.97488	11.5333	11.2484	10.3539	10.0387	9.868	10.7853	9.58187	6.10024	4.33804	9.87049	8.55412	9.95884	8.84474	8.68675	9.03484	9.00321	9.2465		
0.0647	0.1014	0.10411	0.09383	0.18435	0.18443	0.18135	0.14233	0.18058	0.25426	0.25028	0.24141	0.48943	0.24727	0.26286	0.26448	0.26653	0.26247	0.2534	0.27236	0.26854	0.27027	0.26384	0.26384	0.26384	0.26384	0.26384	0.26384	0.26384	0.26384	0.26384	0.26384	0.26384	0.26384		
1.11891	0.98953	1.01758	1.02386	1.77688	1.95054	1.81	1.96504	2.24202	1.81048	2.18378	2.47174	2.10722	2.48847	2.32436	2.4914	2.67124	1.90045	2.26378	2.30352	1.79772	1.80059	1.8117	1.84136	1.9071	1.7282	1.87678	1.95844	1.9244	1.48776	1.90501	1.38782	1.82784	1.92191		
7.86174	8.19831	8.51033	8.76201	9.868	9.37469	9.7853	10.3036	11.095	11.2022	11.6454	12.1412	11.7348	11.2476	11.7656	11.5886	12.1275	12.348	11.9819	11.4129	10.7966	10.6101	11.0135	11.0632	11.0544	11.2482	11.5796	11.792	11.4617	11.5104	11.6449	11.7139	11.6944	11.9318		
1.11901	1.24238	1.29538	1.33818	1.4234	1.67724	1.61352	2.08817	1.12203	2.3844	2.20164	2.39188	2.43919	2.46882	2.78651	3.0483	3.37513	3.2224	3.39104	3.05588	3.74258	3.83175	4.13989	4.21787	4.61162	3.98631	4.23884	4.71612	4.97162	4.87601	5.76548	6.03234	5.18338	5.11033		
0.24941	0.32093	0.31684	0.32043	0.33408	0.33048	0.47623	0.44629	0.51158	0.47572	0.53311	0.50886	0.59943	0.52913	0.54703	0.58716	0.59277	0.54183	0.55928	0.4805	0.45103	0.48665	0.47836	0.4654	0.41907	0.47878	0.53558	0.56872	0.52219	0.5292	0.58826	0.52776	0.60897			
0.17238	0.13643	0.13408	0.1426	0.15043	0.16112	0.17855	0.18627	0.20726	0.21605	0.23651	0.2502	0.25688	0.27381	0.29277	0.29884	0.31138	0.31851	0.28647	0.31644	0.3092	0.31288	0.31826	0.31826	0.31826	0.31826	0.31826	0.31826	0.31826	0.31826	0.31826	0.31826	0.31826			
0.07568	0.06076	0.07182	0.07126	0.06888	0.06587	0.05821	0.06823	0.08177	0.08438	0.08846	0.08501	0.08971	0.08417	0.12881	0.11335	0.11688	0.11671	0.11911	0.11911	0.11805	0.11476	0.11632	0.11462	0.11621	0.11608	0.11616	0.11618	0.11597	0.11377	0.11311	0.08814	0.09918			
0.02034	0.17905	0.19024	0.04004	0.06004	0.17824	0.06206	0.06091	0.77149	0.68248	0.67703	7.00424	7.07174	7.12338	7.54843	7.81405	8.23535	8.02571	8.30053	0.88621	0.11461	0.33981	0.51538	0.68094	7.85828	7.86513	7.88992	7.29033	0.71365	6.41886	5.24445	5.878	5.693	5.84471		
1.11951	1.12578	1.15403	1.2114	1.50817	1.2251	1.30467	1.43227	1.60234	1.60477	1.70115	1.7805	1.78488	1.88223	1.87288	2.06844	2.38407	2.22127	2.20009	2.14113	2.22239	2.25757	2.33086	2.43814	2.47304	2.56032	2.63323	2.66357	2.94841	2.64896	2.65644	2.67292	2.63828	2.67488		
0.9677	0.8772	1.08407	1.06514	1.1035	1.06781	1.1384	1.22703	1.58651	1.43768	1.48114	1.54148	1.58973	1.63902	1.70805	1.78182	1.87588	1.84753	1.81933	1.83558	1.8887	2.00515	2.10805	2.1749	2.2373	2.24689	2.18463	2.22518	2.26044	2.29331	2.17773	2.22303	2.25320	2.25020		
0.1244	0.24482	0.24641	0.29025	0.3014	0.30701	0.33344	0.34406	0.37887	0.41578	0.43488	0.46529	0.49037	0.44773	0.46875	0.49482	0.443	0.50241	0.51773	0.53126	0.52239	0.53175	0.53943	0.54260	0.53642	0.57941	0.57078	0.60409	0.49009	0.30304	0.48947	0.46235	0.44628			
0.41382	0.46877	0.48092	0.55048	0.57581	0.57215	0.64237	0.68713	0.71525	0.82313	0.86808	0.89774	1.096	0.82059	0.96883	0.93138	1.00521	1.12893	1.10039	1.08579	1.04743	1.09654	1.11213	1.12482	1.03633	1.14089	0.97906	0.93514	0.93183	0.86474	0.96468	0.90901	0.85451			
0.24845	0.23968	0.22936	0.24814	0.22889	0.23192	0.25351	0.29687	0.31195	0.33962	0.33881	0.39547	0.40218	0.41281	0.46124	0.48063	0.60778	0.60615	0.64285	0.66341	0.68232	0.6841	0.6944	0.73477	0.72282	0.71841	0.75021	0.73481	0.82434	0.71887	1.07887	1.14323	1.14163			
0.1831	0.19374	0.14756	0.15935	0.16343	0.1689	0.18293	0.21382	0.20854	0.21744	0.23449	0.21791	0.21855	0.22551	0.22919	0.24815	0.26387	0.24787	0.25207	0.26278	0.26208	0.26643	0.27231	0.2806	0.26431	0.25779	0.26386	0.2568	0.24636	0.23883	0.24144	0.24789	0.24789			

rys. 3.1

W sytuacji gdzie braki danych występowały w kolumnach „Region” i „IncomeGroup” dotyczyły one zbiorczych regionów (np. „OECD members”, „South Asia (IDA & IBRD)”, „Least developed countries: UN classification”). Jako że zbiorcze regiony zawierają państwa położone w różnych częściach globu („OECD members” (rys. 3.2)) lub należące do różnych grup dochodowych („Early-demographic dividend” (rys. 3.3, rys 3.4)) nie było możliwe przypisanie im wartości w tych kolumnach.



rys. 3.2



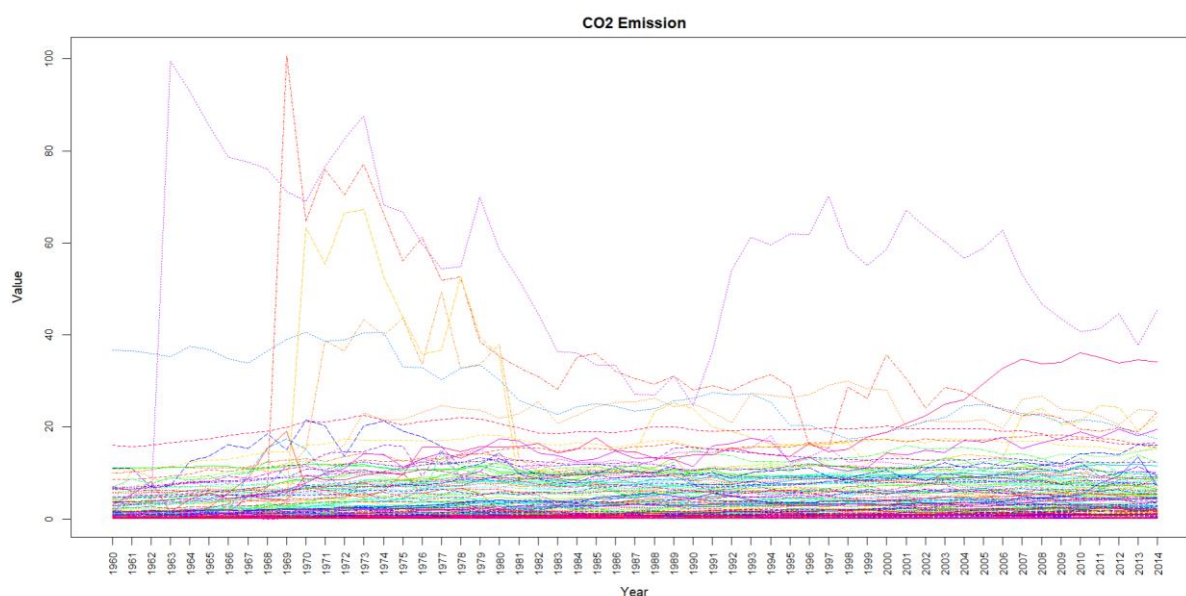
rys. 3.3

Ireland	IRL	CO2 emis: EN.ATM.Ci Europe & C	High income
<u>Iran, Islamic Rep.</u>	<u>IRN</u>	CO2 emis: EN.ATM.Ci Middle Ea	<u>Upper middle income</u>
Iraq	IRQ	CO2 emis: EN.ATM.Ci Middle Ea	Upper middle income
Iceland	ISL	CO2 emis: EN.ATM.Ci Europe & C	High income
<u>Israel</u>	<u>ISR</u>	CO2 emis: EN.ATM.Ci Middle Ea	<u>High income</u>
Italy	ITA	CO2 emis: EN.ATM.Ci Europe & C	High income
Jamaica	JAM	CO2 emis: EN.ATM.Ci Latin Amer	Upper middle income
<u>Jordan</u>	<u>JOR</u>	CO2 emis: EN.ATM.Ci Middle Ea	<u>Lower middle income</u>
Japan	JPN	CO2 emis: EN.ATM.Ci East Asia &	High income
Kazakhstan	KAZ	CO2 emis: EN.ATM.Ci Europe & C	Upper middle income
Kenya	KEN	CO2 emis: EN.ATM.Ci Sub-Sahar	Lower middle income
Kyrgyz Republic	KGZ	CO2 emis: EN.ATM.Ci Europe & C	Lower middle income
Cambodia	KHM	CO2 emis: EN.ATM.Ci East Asia &	Lower middle income
<u>Kiribati</u>	<u>KIR</u>	CO2 emis: EN.ATM.Ci East Asia &	<u>Lower middle income</u>
St. Kitts and Nevis	KNA	CO2 emis: EN.ATM.Ci Latin Amer	High income

rys. 3.4

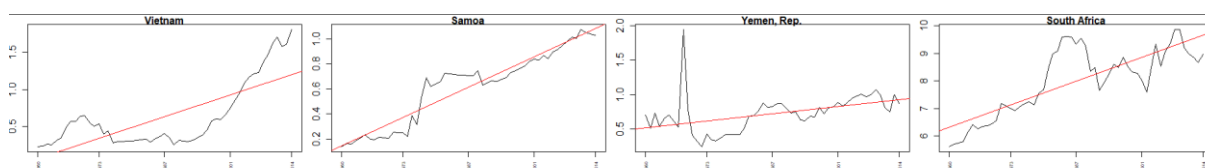
Kolumny „Indicator Name” i „Indicator Code” zostały usunięte, ponieważ nie przynoszą one żadnych informacji o danych – ich wartości dla każdego wiersza są takie same.

Następnie wykonano wykres porównawczy dla zmian wysokości emisji CO₂ na przestrzeni lat dla wszystkich państw. Można zauważyć cztery linie: fioletową, czerwoną, żółtą i niebieską wybijającą się ponad większość wartości.

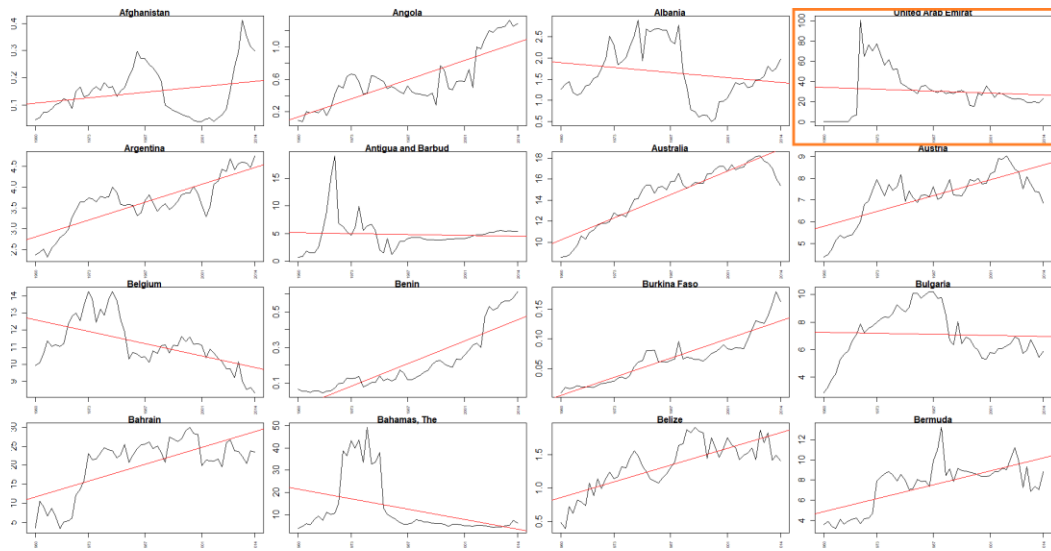


rys. 3.5

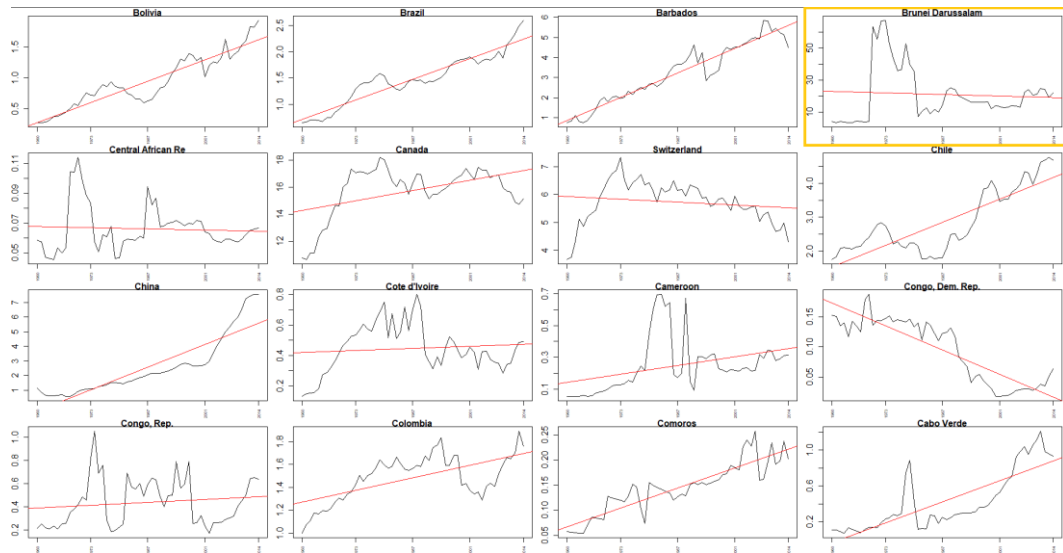
Po wykonaniu osobnych wykresów liniowych i linii trendu (rys. 3.6–15) dla każdego z państw można zauważyć że szukane państwa to: linia fioletowa – Katar, linia pomarańczowa – Zjednoczone Emiraty Arabskie, linia żółta – Brunei, linia niebieska – Luksemburg. W wielu państwach występuje wyraźny trend wzrostu emisji CO₂ (np. Wietnam, Benin, Indonezja, Indie).



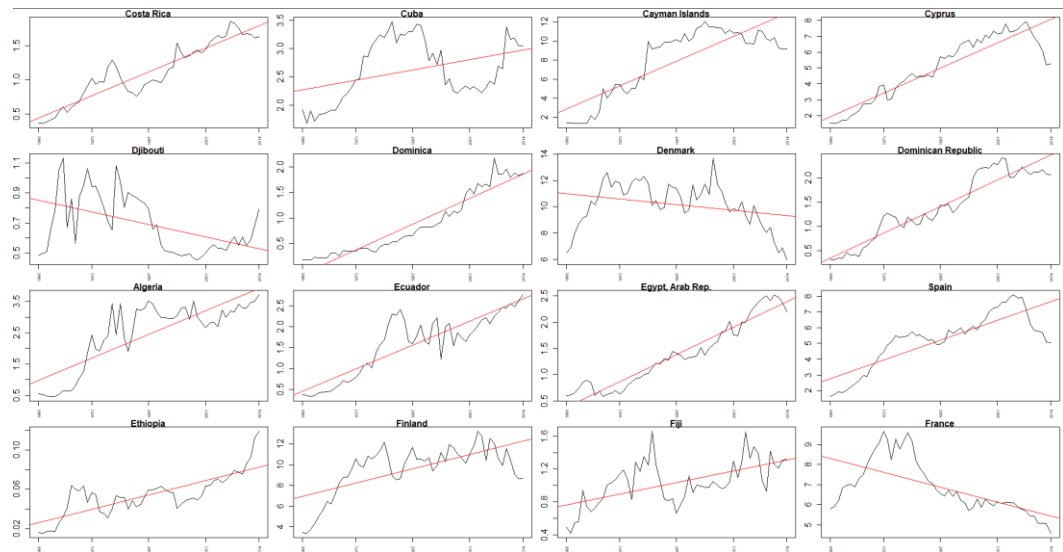
rys. 3.6



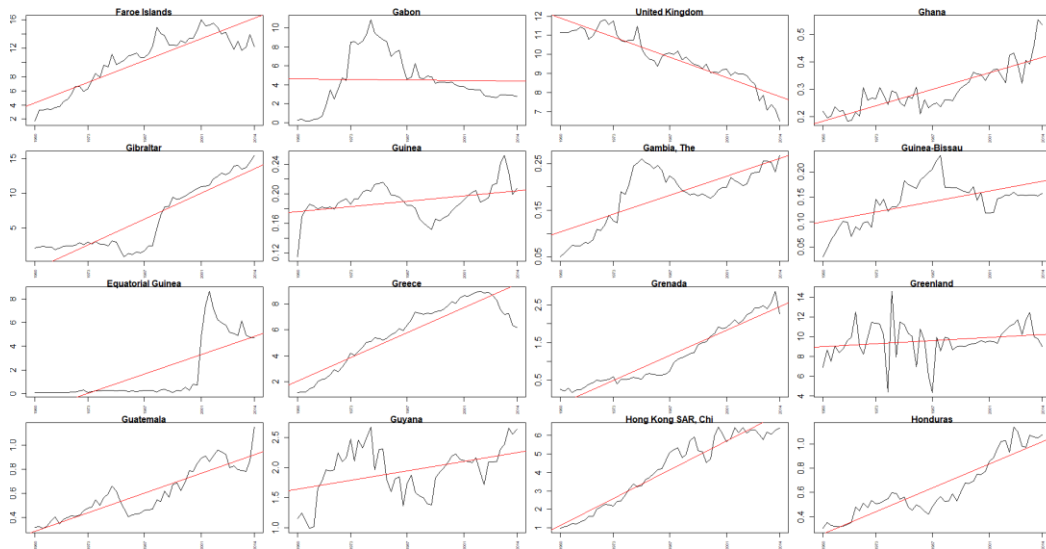
rys. 3.7



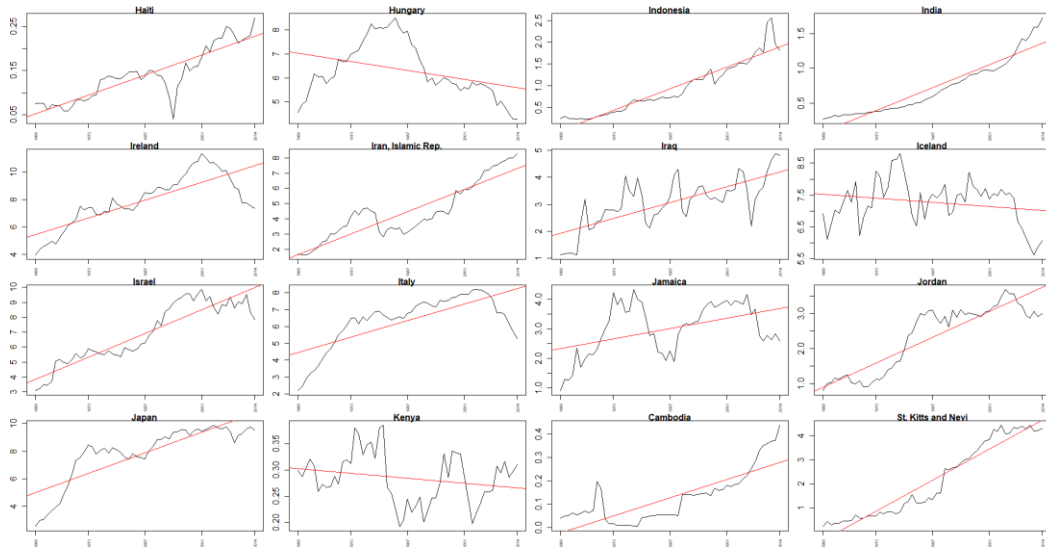
rys. 3.8



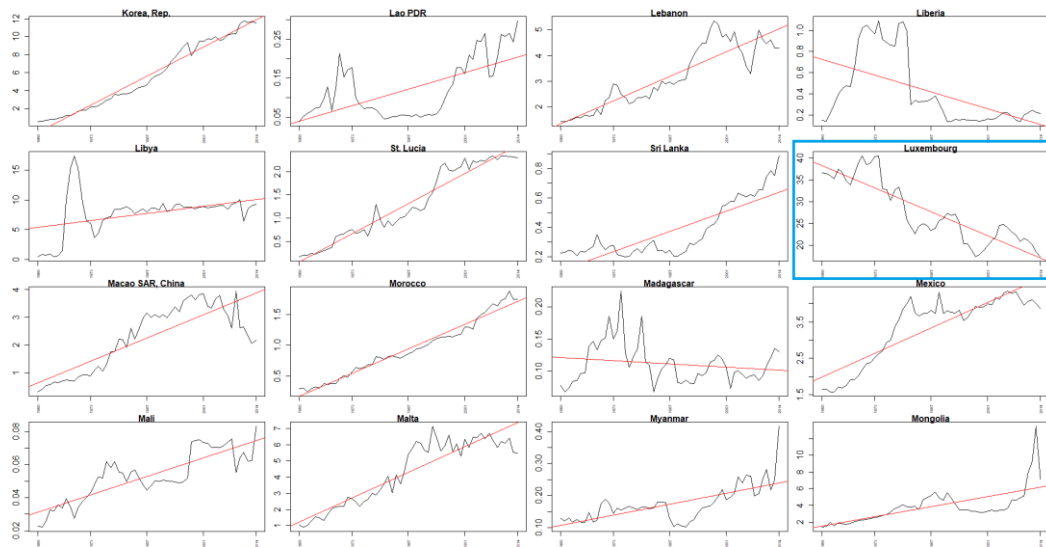
rys. 3.9



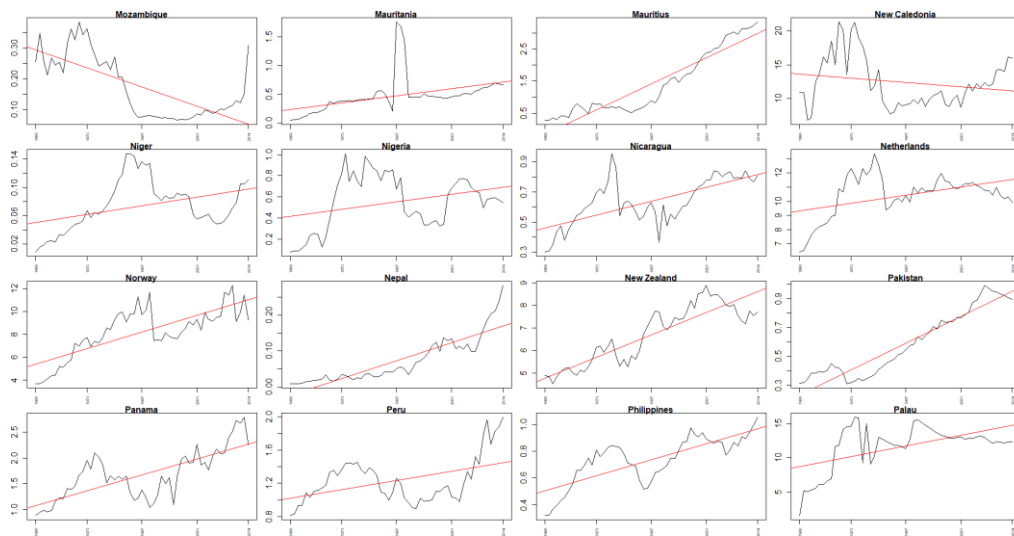
rys. 3.10



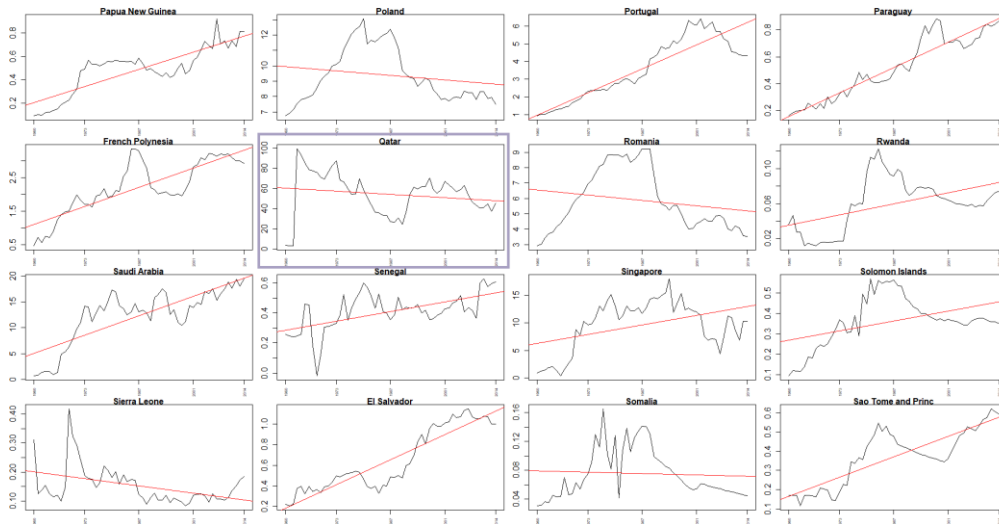
rys. 3.11



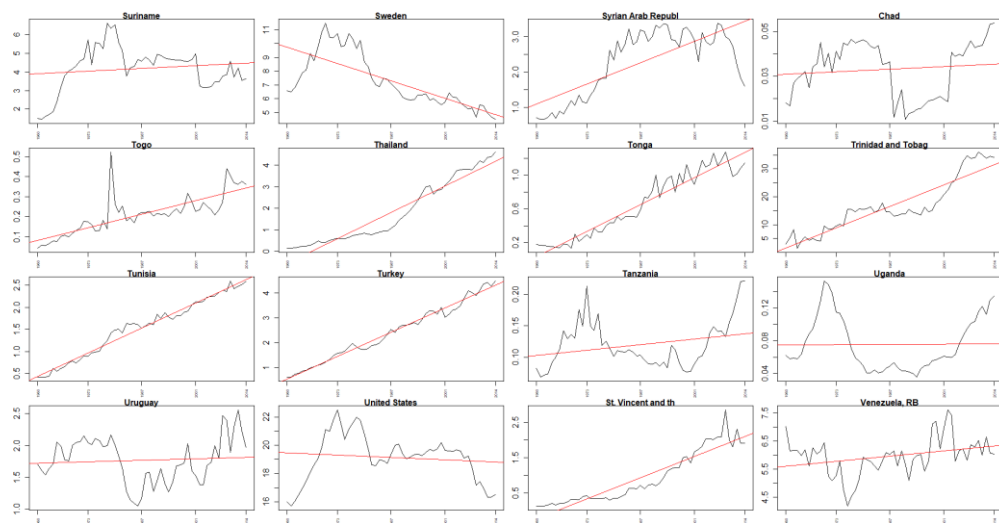
rys. 3.12



rys. 3.13



rys. 3.14

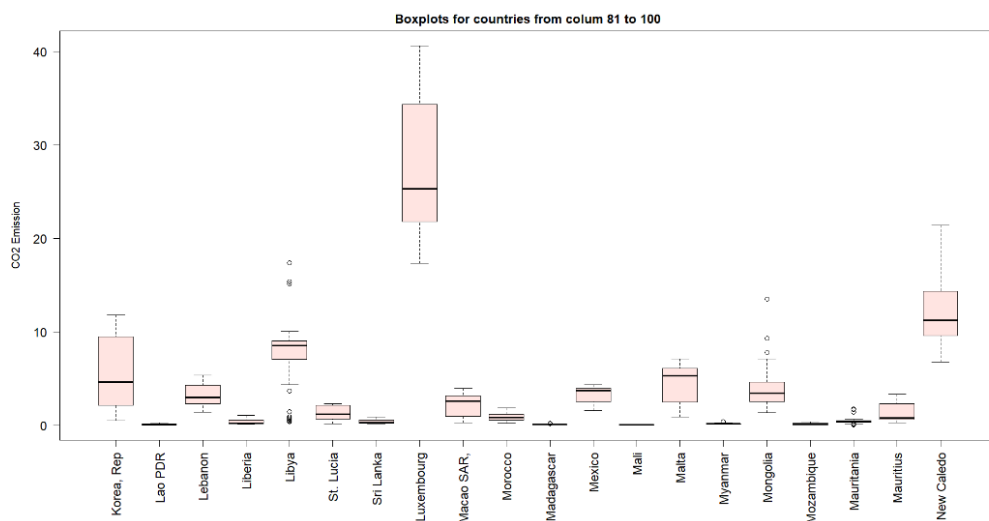


rys. 3.15

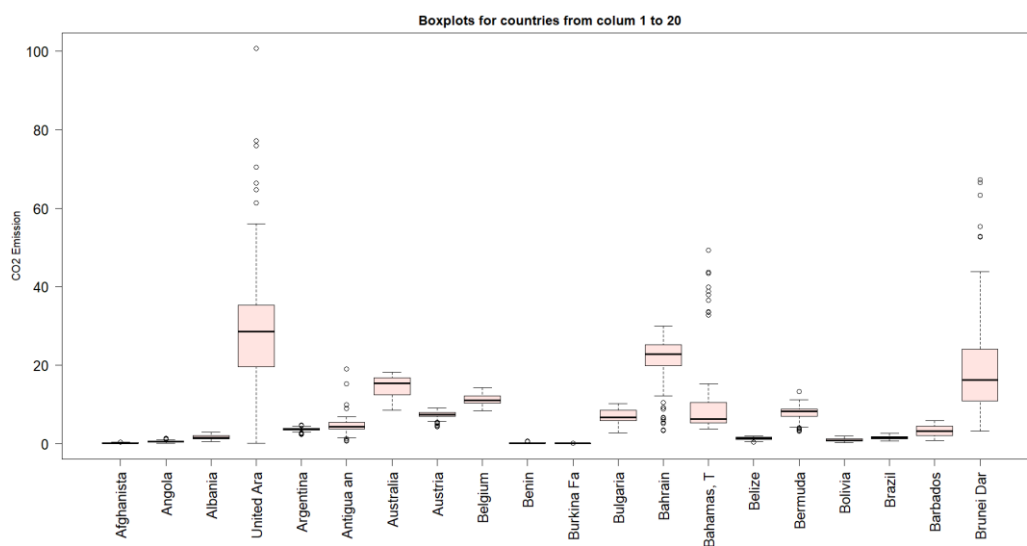
Na wykresach Kataru, Brunei i ZEA widać duży skok wartości w latach 60. Po tym skoku wartości stopniowo maleją do lat 80., nie jest więc to pojedyncza wartość ekstremalna, ale ciąg wyższych wartości trwających przez okres 20 lat. Taka sytuacja może mieć związek ze złożami ropy naftowej, które zaczęły odgrywać najważniejszą rolę w gospodarce tych państw.

Luksemburg natomiast, jest małym państwem o wysokich dochodach, co może wpływać na wysokie wartości wskaźnika CO₂, jednak stopniowo malejące na przestrzeni lat.

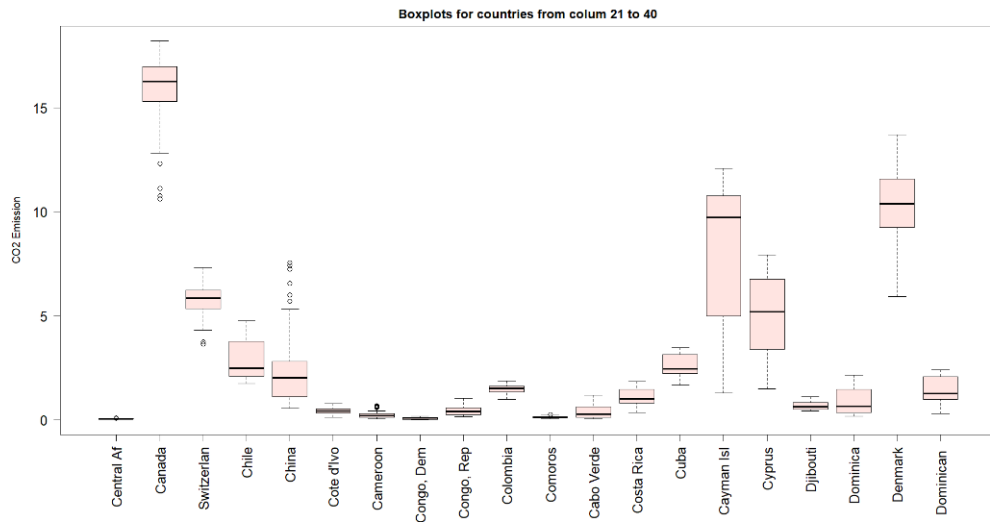
Na rys. od 16 do 24 zaprezentowano boxploty dla każdego z państw, większość z nich nie posiada obserwacji odstających, co oznacza brak poważnych skoków w wysokościach wskaźnika.



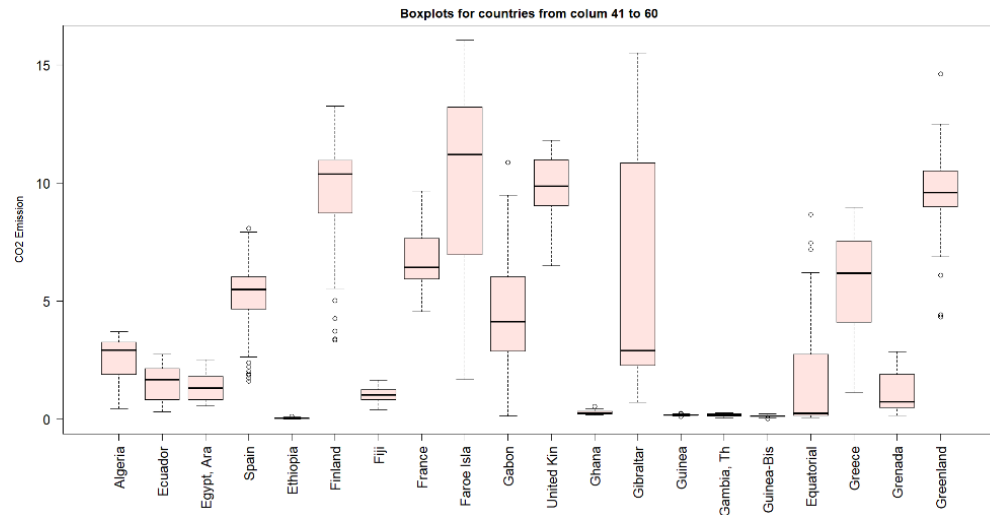
rys. 3.16



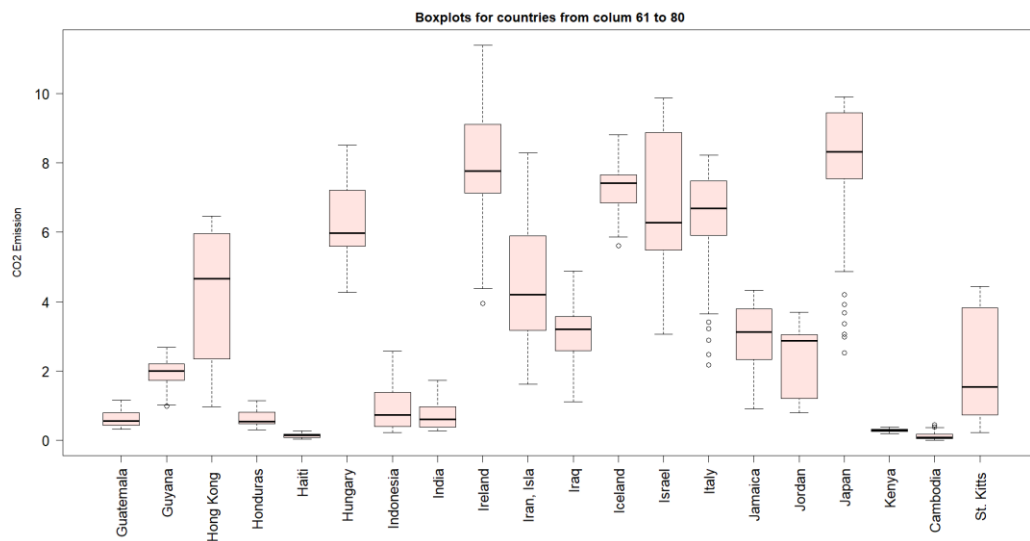
rys. 3.17



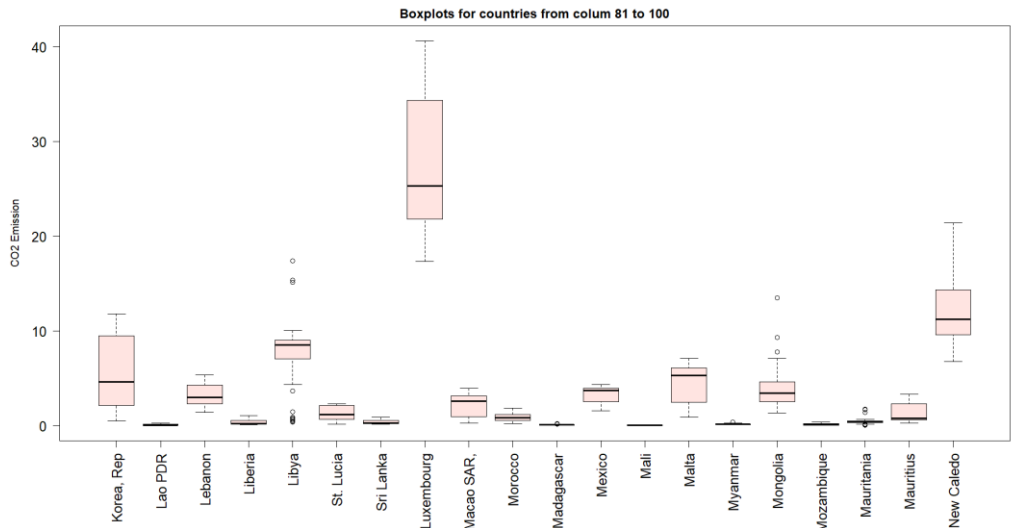
rys. 3.18



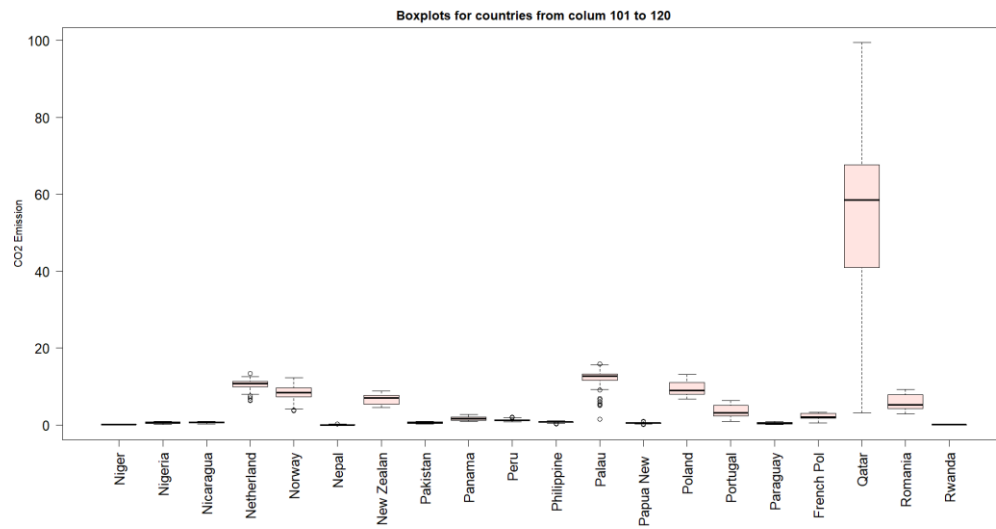
rys. 3.19



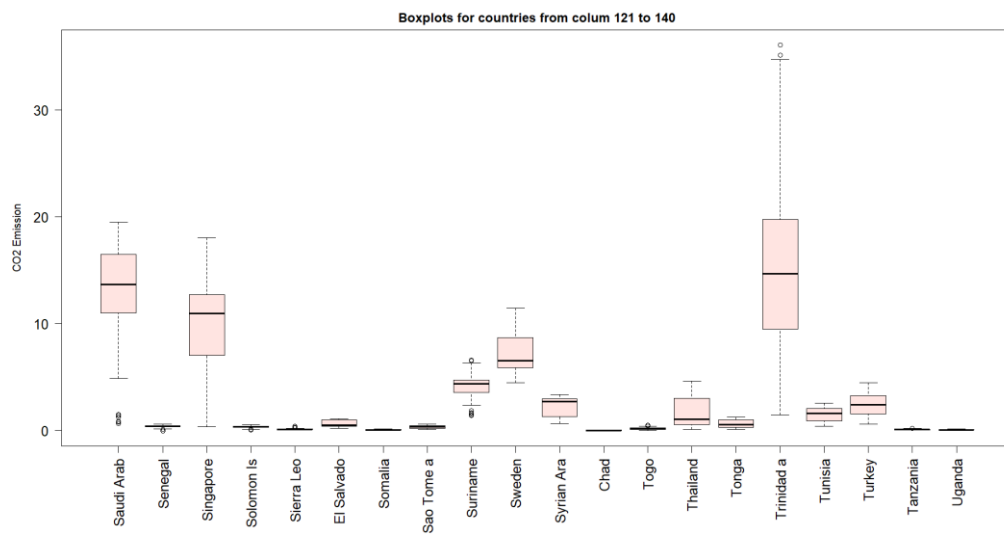
rys. 3.20



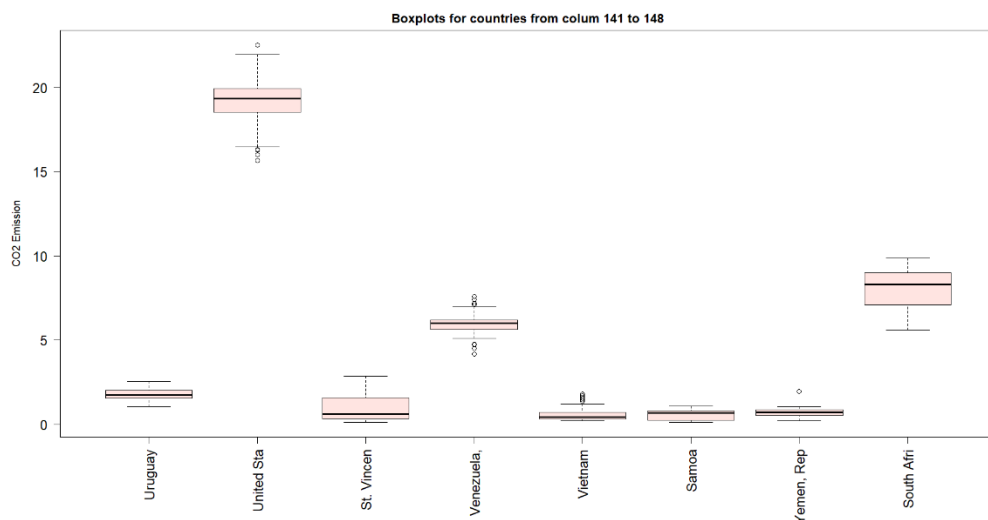
rys. 3.21



rys. 3.22

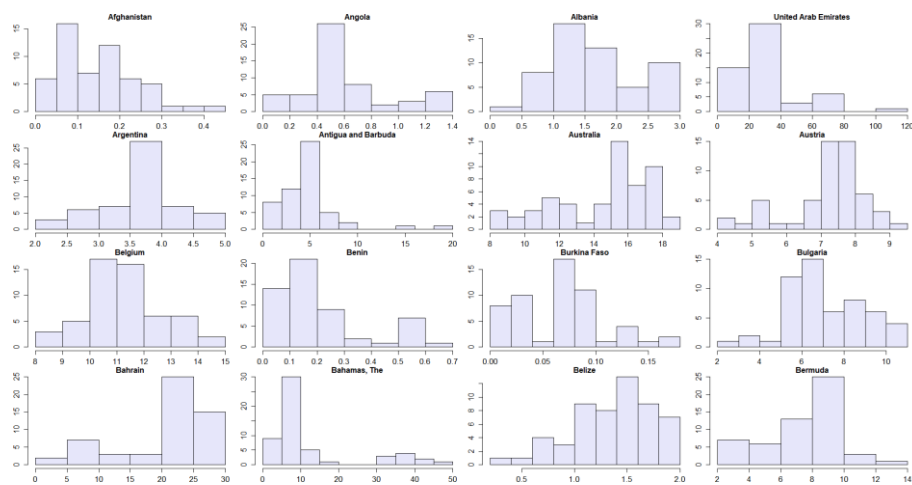


rys. 3.23

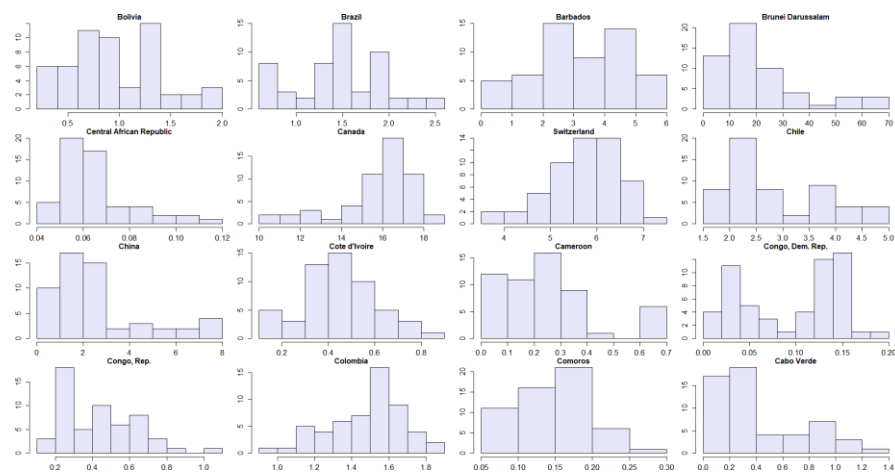


rys. 3.24

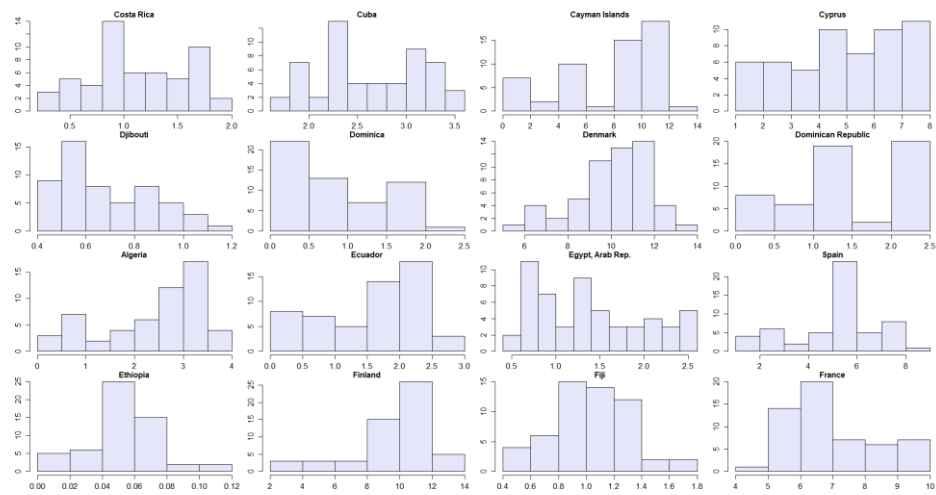
Zostały również przygotowane histogramy (rys. 3.25–34).



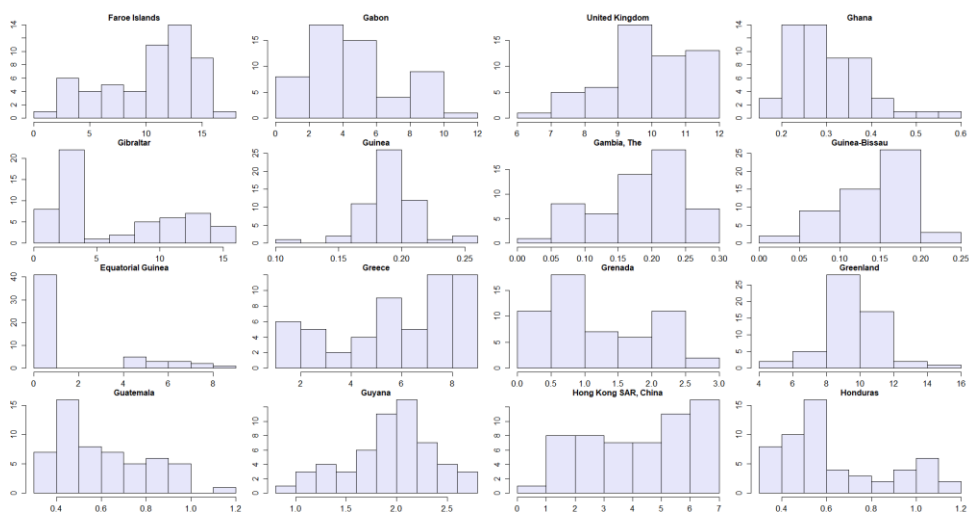
rys. 3.25



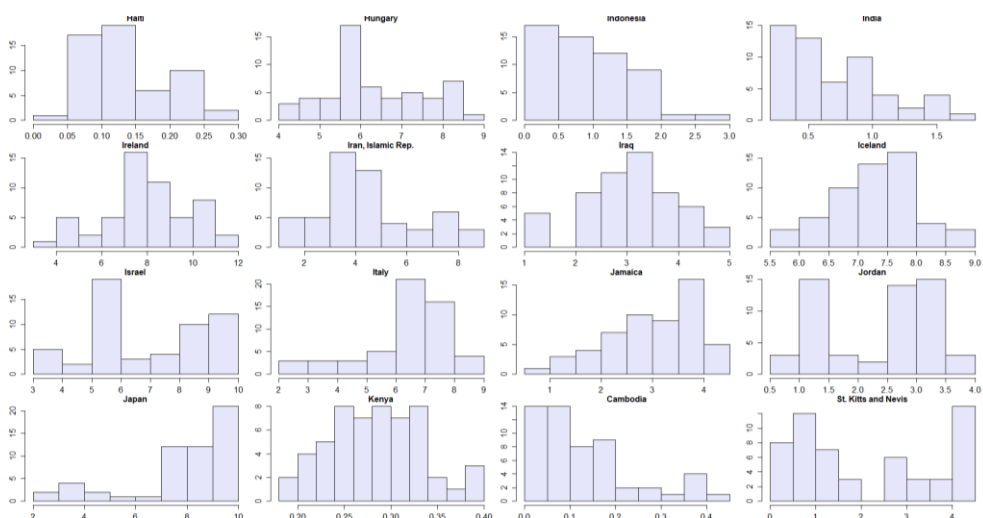
rys. 3.26



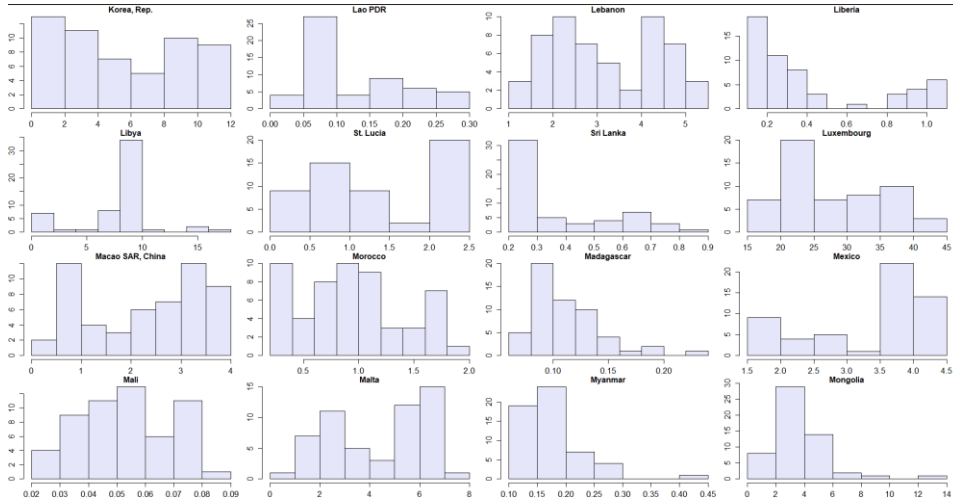
rys. 3.27



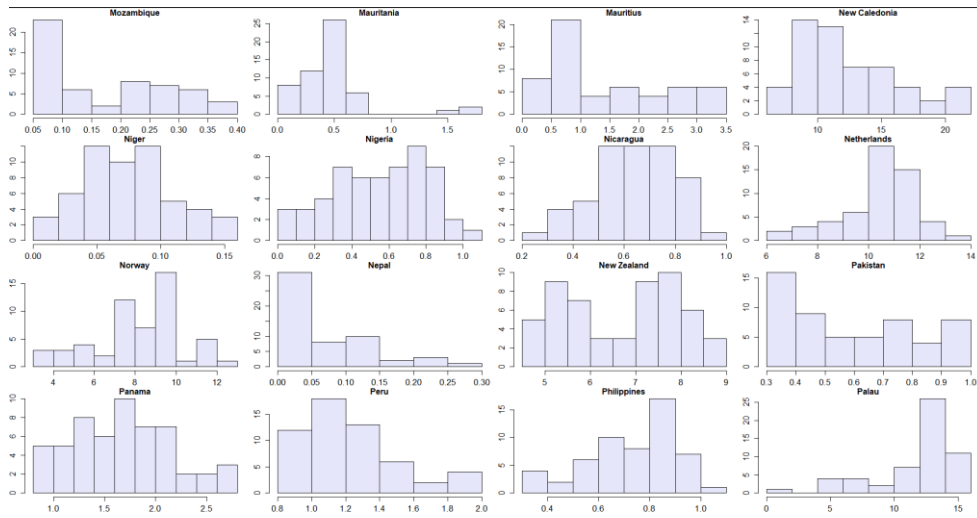
rys. 3.28



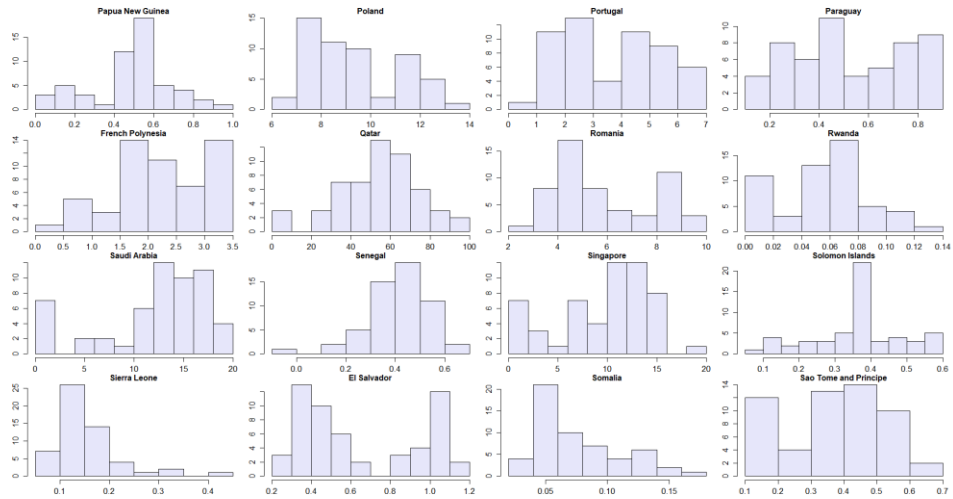
rys. 3.29



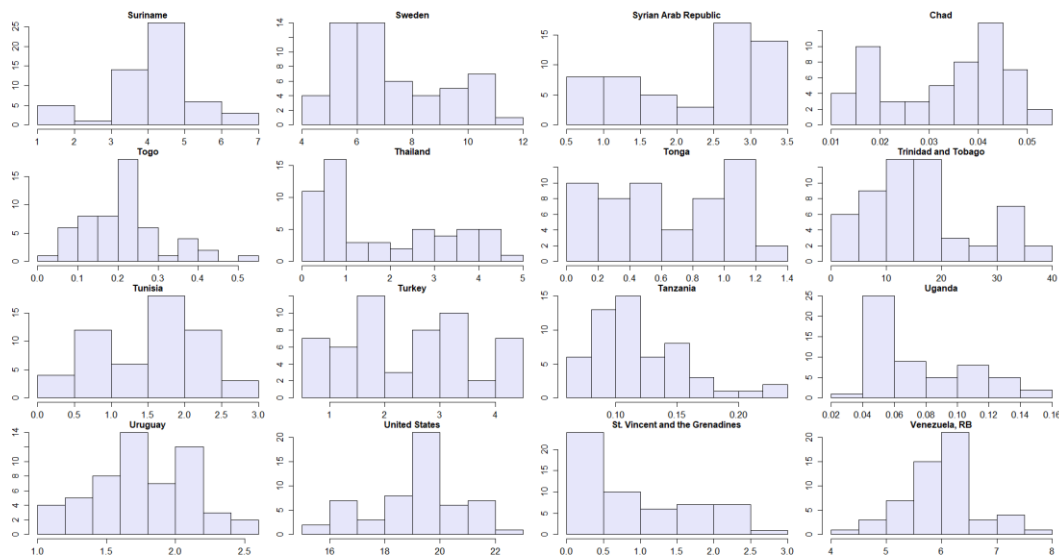
rys. 3.30



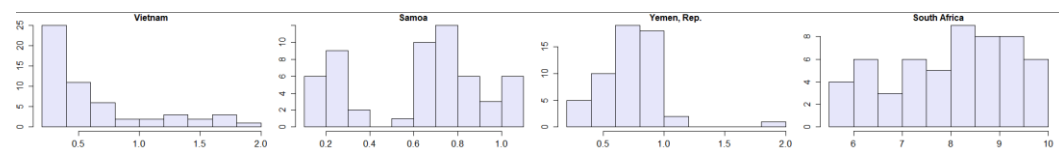
rys. 3.31



rys. 3.32



rys. 3.33



rys. 3.34

4. Model regresyjny Knn

Wylosowany model Knn, w odmianie regresyjnej, został wykorzystany w projekcie ze względu na możliwość przewidywania wartości wskaźnika w ostatnim roku szeregu czasowego. Model ten powinien sobie poradzić z nieliniowymi zmianami występującymi w przypadku większości państw, jest elastyczny i dzięki wykorzystywanej metodzie najbliższych sąsiadów może dostosowywać się do lokalnych zmian.

Aby uzyskać wartości numeryczne w całej ramce danych, dodano dwie nowe kolumny „NumRegion” i „NumIncomeGroup”. Kolumna „NumRegion” zawiera klasy od 1 do 7 odpowiadające wartościom opisowym w kolumnie Region, a kolumna „NumIncomeGroup” zawiera klasy od 1 do 4 odpowiadające wartościom opisowym w kolumnie „IncomeGroup”.

Po standaryzacji danych przy pomocy funkcji `scale()` podzielono zbiór na treningowy – 80%, testowy do walidacji – 20% i testowy do predykcji – 4 obserwacje. Następnie wykonano model Knn funkcją `knnreg()` z pakietu `caret` (rys. 4.1). Najefektywniejszą wartością parametru `k`, czyli liczby najbliższych sąsiadów uwzględnianych przy przewidywaniu wartości, była liczba 4.

```
model_knn<- knnreg(train_std_knn, train_y_knn, k = 4)
```

rys. 4.1

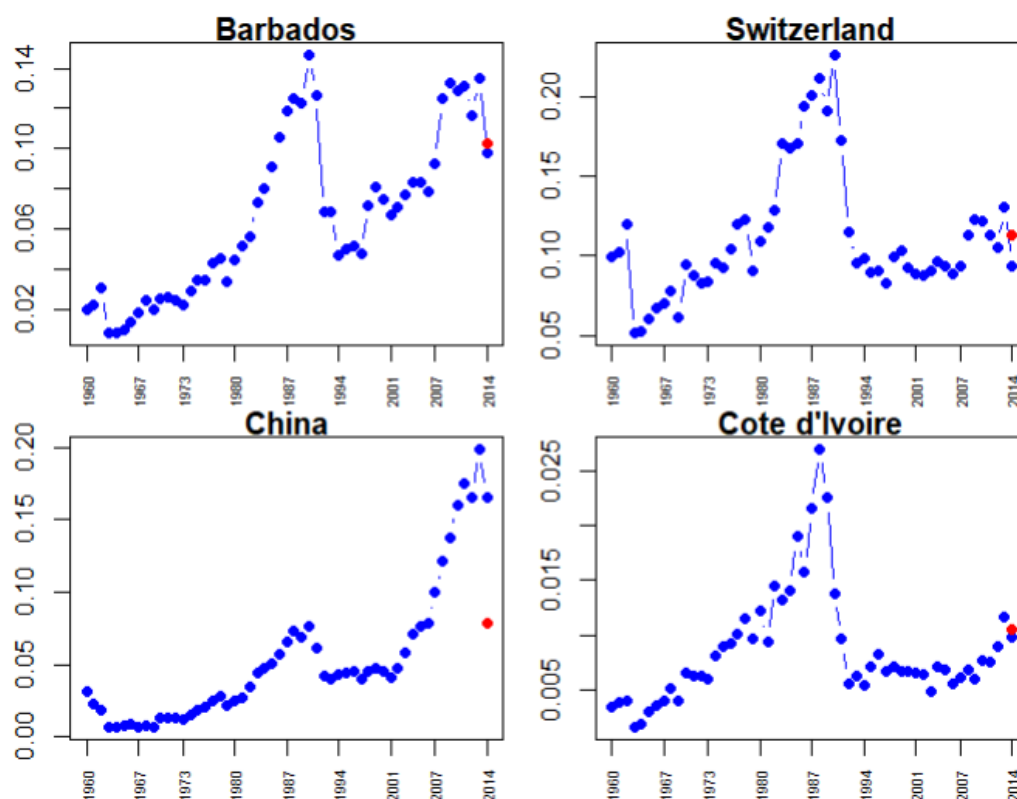
Najpierw wykonano model na większym zbiorze danych testowych do walidacji.

Na rys. 4.2 widać wartości rzeczywiste zestawione z przewidywanymi dla kolejnych wierszy zbioru testowego na danych zestandaryzowanych dla pierwszych czterech wartości.

	Barbados	Szwajcaria	Chiny	Wyb. Kości Słoniowej
real	0.0979	0.0940	0.1652	0.0098
predicted	0.1023	0.1126	0.0782	0.0105

rys. 4.2

Poniżej zaprezentowano wykres porównujący przewidywaną wartość dla roku 2014 (**czzerwony**) do wartości rzeczywistych na przestrzeni lat 1960–2014 (**niebieski**). Dla 3 na 4 wylosowanych państw, Barbadosu, Szwajcarii i Wybrzeża Kości Słoniowej, można optycznie stwierdzić, że przewidziane wartości były bliskie rzeczywistym. Większa różnica natomiast występowała w przypadku Chin. (rys. 4.3)



rys. 4.3

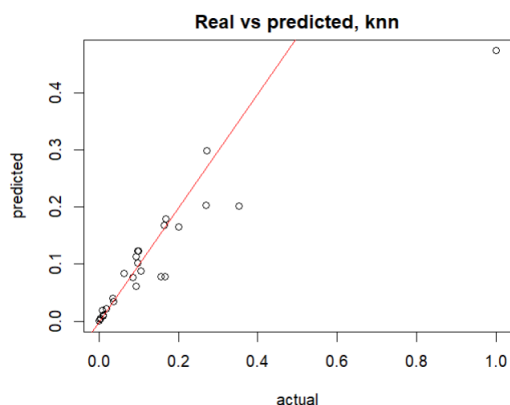
Obliczono również wskaźnik pozwalający ocenić jakość modelu i prognozy. Stosunkowo wysoka wartość współczynnika determinacji R^2 wskazuje na dobre dopasowanie modelu do danych. Wartości MSE i RMSE są bardzo niskie. Niepokojąca jest natomiast wysoka wartość MAPE, która sugeruje dużą średnią procentową różnicę między przewidywanymi, a rzeczywistymi wartościami. (rys. 4.4)

test: 20%	R^2	MSE	RMSE	MAPE
Knn	0.85	0.01	0.11	25.22%

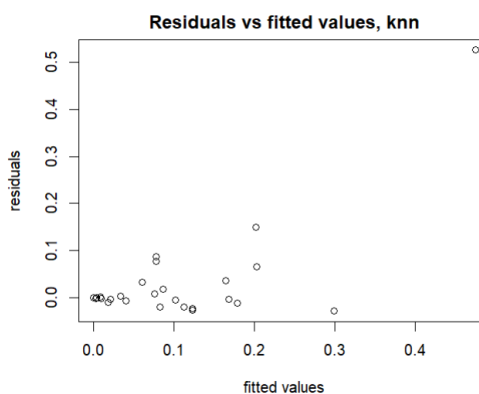
rys. 4.4

Następnie wykonano wykres porównujący rzeczywiste wartości z przewidywanymi (rys. 4.5) oraz wykresy dotyczące reszt: wykres rozrzutu reszt (rys. 4.6), histogram reszt (rys. 4.7) , porównanie rozkładu reszt z rozkładem normalnym (rys. 4.8) , ACF (rys. 4.9) i PACF (rys. 4.10) .

Na wykresie porównujący rzeczywiste wartości z przewidywanymi (rys. 4.5) wartości nie są bardzo oddalone od czerwonej linii, a na wykresie dotyczącym rozrzutu reszt (rys. 4.6) wartości są ułożone mniej więcej losowo. Oba wykresy zaburza jednak mocno odstająca pojedyncza wartość w prawym górnym rogu, która może świadczyć o nagłym skoku wskaźnika, którego model nie przewidział.

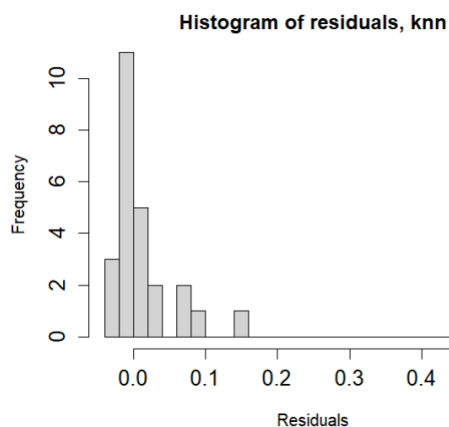


rys. 4.5

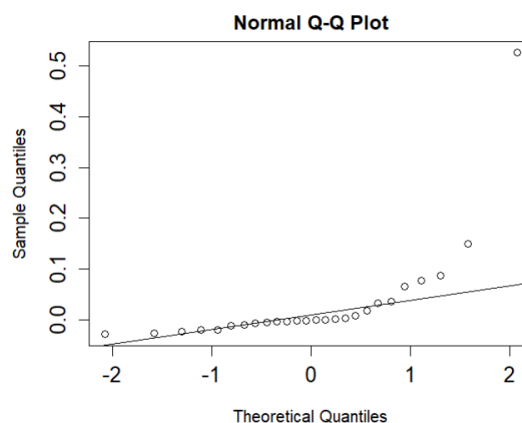


rys. 4.6

Histogram reszt (rys. 4.7) jest prawostronny i nie przypomina pożądanego kształtu rozkładu normalnego. Brak rozkładu normalnego potwierdza wykres porównania rozkładu reszt z rozkładem normalnym (rys. 4.8), gdzie punkty po prawej stronie wykresu są mocno odchylone od linii referencyjnej.

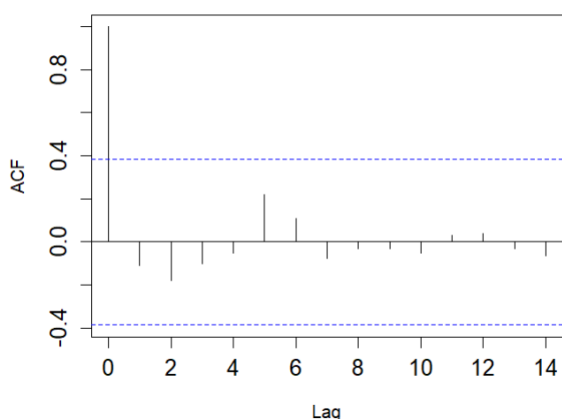


rys. 4.7

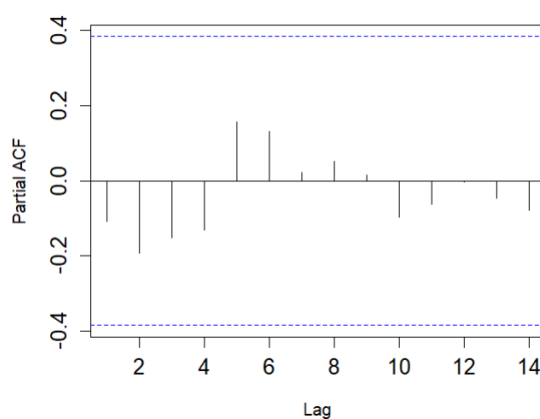


rys. 4.8

Wykresy ACF (rys. 4.9) i PACF (rys. 4.10) potwierdzają brak występowania trendu i sezonowości, pionowe czarne linie nie przebijają poziomych linii niebieskich (poza pierwszą linią na wykresie ACF oznaczającą korelację z samym sobą)



rys. 4.9



rys. 4.10

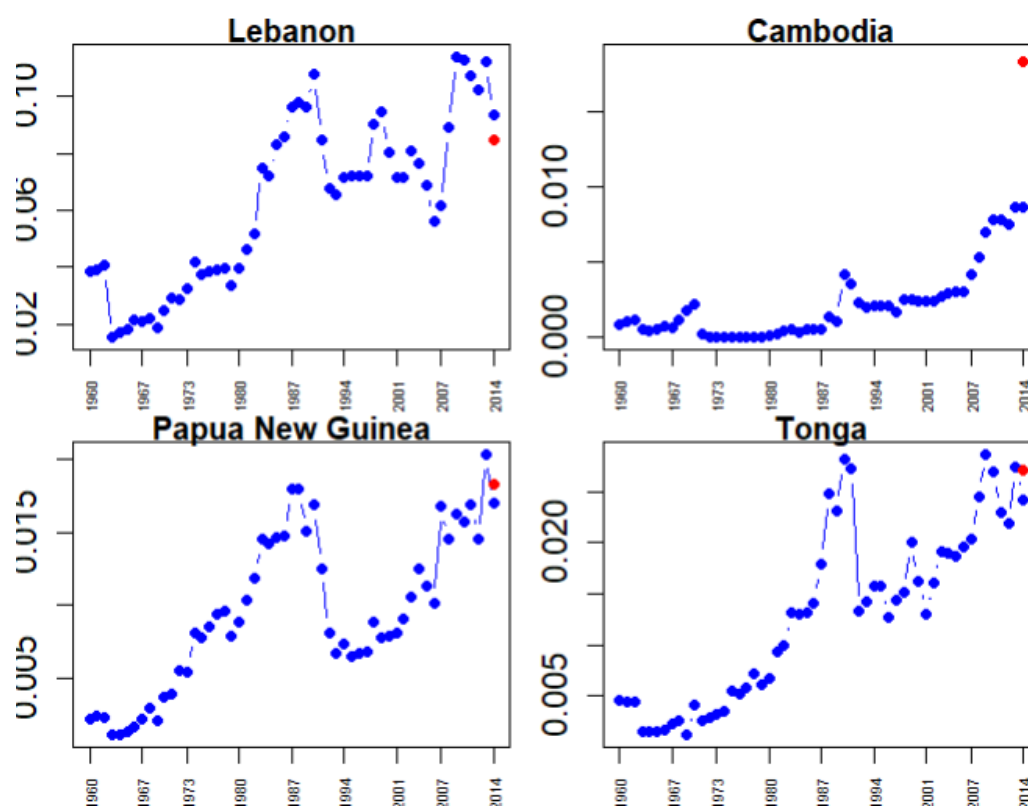
Następnie sprawdzono model na zbiorze testowym do predykcji zawierającym 4 obserwacje.

W poniższej tabelce przedstawiono wartości rzeczywiste zestawione z przewidywanymi dla wszystkich państw. (rys. 4.11)

	Liban	Kamboża	Papua Nowa Gwinea	Tonga
real	0.0936	0.0086	0.0169	0.0242
predicted	0.0846	0.0183	0.0183	0.0270

rys. 4.11

Przygotowano wykres porównujący zmianę w czasie danych rzeczywiste z przewidywaną wartością w roku 2014 (rys. 4.12). Bardzo dobrze przewidziane zostały wartości dla Libanu, Papui Nowej Gwinei i Tongy. Natomiast już dla Kambodży wartość przewidywana mocno odbiegała od rzeczywistej.



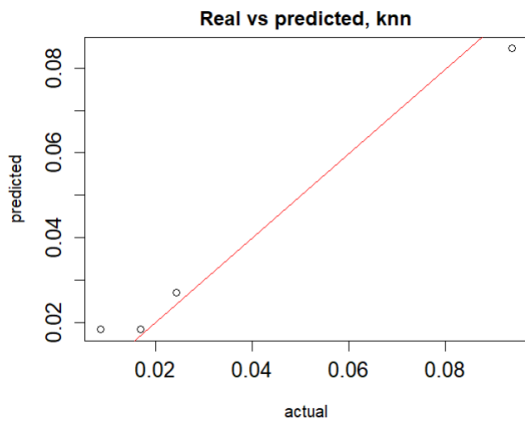
rys. 4.12

W poniższej tabelce (rys. 4.13) zostały przedstawione poprzednio liczone wskaźniki dla nowego podziału. Wartość MSE i za tym RMSE mocno spadły, natomiast współczynnik R2 bardzo wzrósł, co świadczy on o dobrym dopasowaniu modelu do danych. Niepokojący jest natomiast wzrost MAPE do aż 35%.

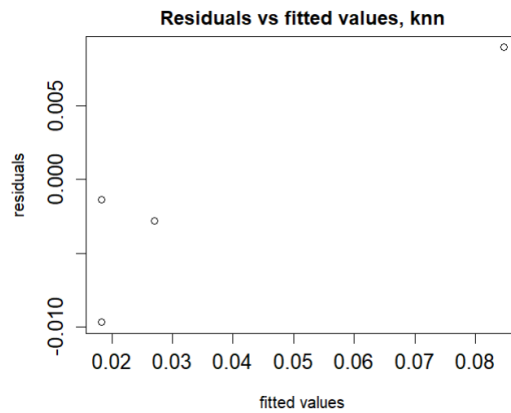
test: 4	R2	MSE	RMSE	MAPE
Knn	0.99	4.61E-05	0.0067	35.25%

rys. 4.13

Na wykresie porównującym rzeczywiste wartości z przewidywanymi (rys. 4.14) widać, że punkty układają się stosunkowo blisko przekątnej. Punkty na wykresie rozrzutu reszt (rys. 4.15) są rozrzucone losowo i nie układają się w żaden konkretny kształt.

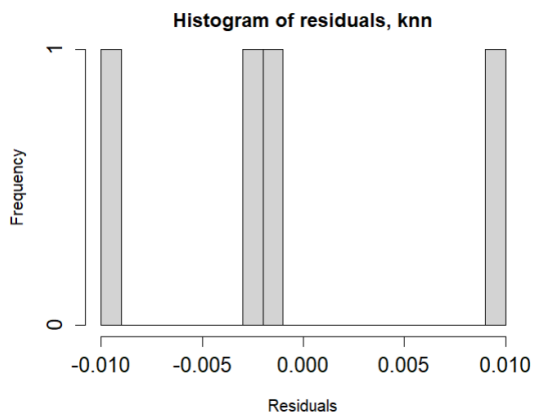


rys. 4.14

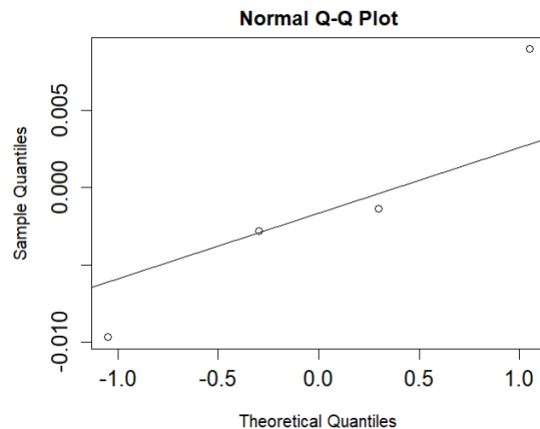


rys. 4.15

Histogram reszt (rys. 4.16) nie przypomina rozkładu normalnego jednak wpływ na to może mieć mała ilość obserwacji. Na wykresie porównania rozkładu reszt z rozkładem normalnym (rys. 4.17) odchylenia na krańcowych wartościach oznaczają większą ilość wartości ekstremalnych niż w przypadku występowania prawdziwego rozkładu normalnego.

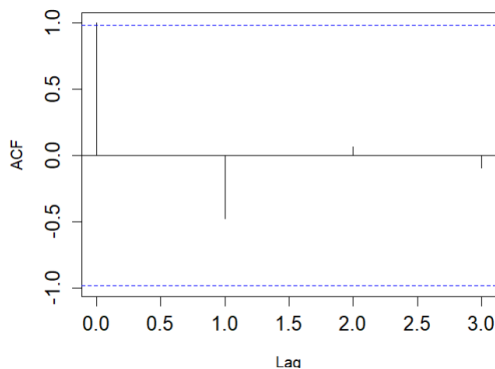


rys. 4.16

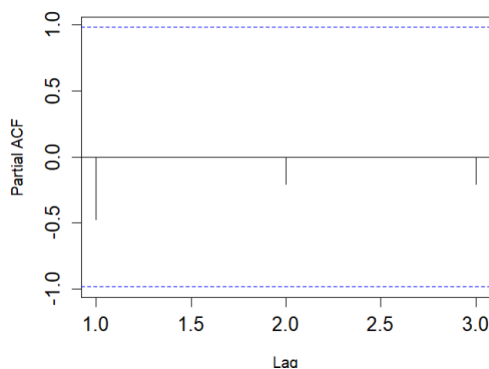


rys. 4.17

Ponownie wykresy ACF (rys. 4.18) i PACF (rys. 4.19) potwierdzają brak występowania trendu i sezonowości.



rys. 4.18



rys. 4.19

Po analizie danych stwierdzono, że odstający punkt reprezentuje Katar, który charakteryzuje się dużymi zmianami wskaźnika. Pomimo braku poprawnej predykcji w tym punkcie, model bardzo dobrze poradził sobie z predykcją dla państw z mniejszymi zmianami wartości, co potwierdzają wysokie wartości R², niskie MSE i wizualnie wykresy z nałożonymi wartościami rzeczywistymi i przewidywanymi (rys. 4.3, rys. 4.12).

5. Model regresyjny sieci neuronowych MLP

Do wykonania drugiego modelu wybrano sieci neuronowe MLP z uwagi na możliwość wykorzystania warstw ukrytych, które mogą przetwarzać złożone dane. Model jest również elastyczny i może dostosowywać się do zmieniających się warunków.

Ponownie podzielono zestandaryzowane dane na trzy zbiory, treningowy, walidacyjny i do predykcji. Do wykonania modelu użyto funkcji `neuralnet` z dwoma ukrytymi warstwami: pierwszą z czterema neuronami i drugą z trzema. (rys. 5.1)

```
model_mlp <- neuralnet(year2014 ~ . ,
                        data = data_train_mlp, hidden = c(4, 3),
                        linear.output = TRUE)
```

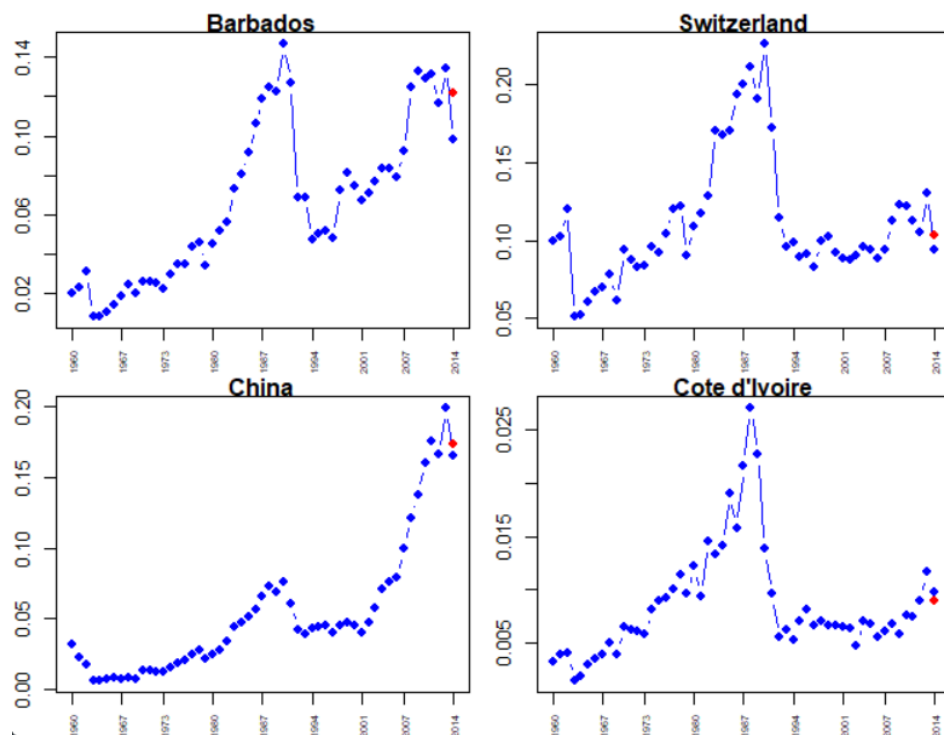
rys. 5.1

Na rys. 5.2 przedstawiono wartości przewidywane i rzeczywiste dla pierwszych czterech państw ze zbioru testowego do walidacji.

	Barbados	Szwajcaria	Chiny	Wyb. Kości Słoniowej
real	0.0979	0.0940	0.1652	0.0098
predicted	0.1217	0.1038	0.1734	0.0089

rys. 5.2

Następnie przygotowano wykres, dzięki któremu można porównać wizualnie jak zostały przewidziane wartości dla roku 2014. W każdym przypadku predykowane wartości są bliskie rzeczywistym (rys. 5.1). Porównując rys. 5.3 do zestawu wykresów na rys. 4.3 wykonanych przy pomocy algorytmu Knn można stwierdzić, że wartości zostały lepiej przyporządkowane przez algorytm MLP.



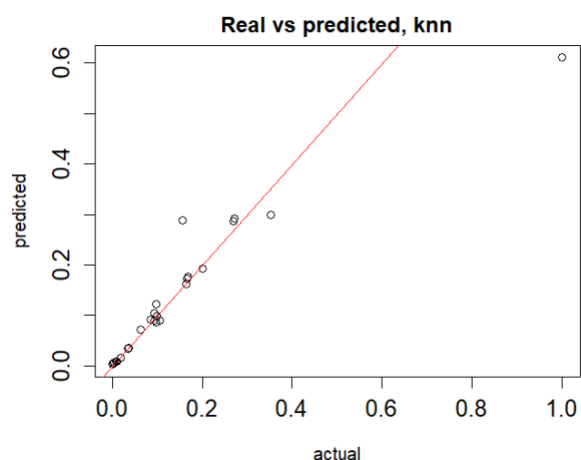
rys. 5.3

Na podstawie tabelki na rys. 5.4 można stwierdzić, że model jest stosunkowo dobrze dopasowany do danych (R^2 równe 0.88). MSE i RMSE są zadowalająco niskie, lecz współczynnik MAPE bardzo wysoki.

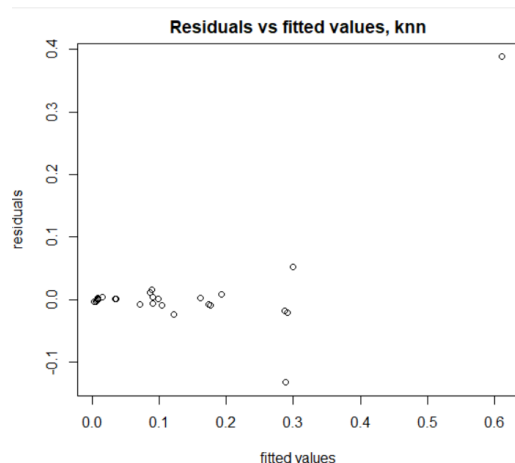
test: 20%	R2	MSE	RMSE	MAPE
Knn	0.88	0.006	0.081	33.93%

rys. 5.4

Ponownie wartości na wykresie porównującym przewidywane i rzeczywiste (rys. 5.5) wartości leżą blisko czerwonej linii, z wyjątkiem jednej odstającej wartości. Na wykresie rozrzutu reszt (rys. 5.6) powtarza się sytuacja z tą samą jedną obserwacją, która mocno odbiega od innych. Reszta wartości jest rozrzucona losowo wokół poziomu zera na osi reszt.

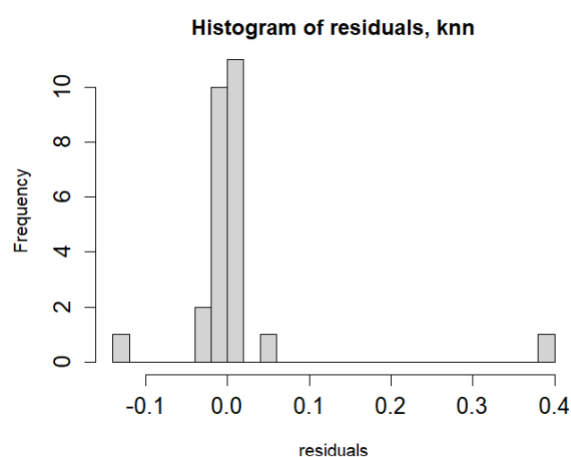


rys. 5.5

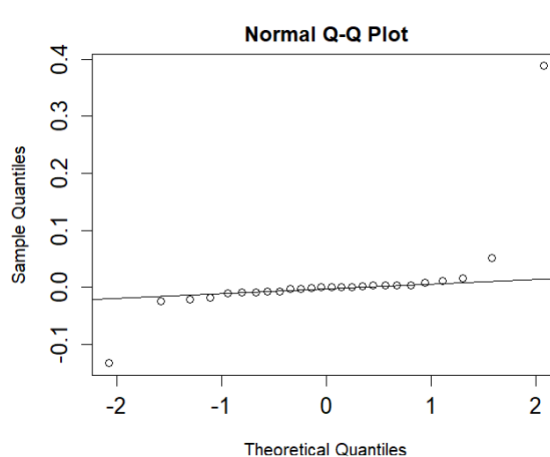


rys. 5.6

Histogram reszt (rys. 5.7) zaburza jedna odstająca wartość, poza tym może on w przybliżeniu przypominać rozkład normalny. Potwierdza to wykres porównania rozkładu reszt z rozkładem normalnym, gdzie poza krańcową wartością, większość obserwacji leży na linii referencyjnej (rys. 5.8).

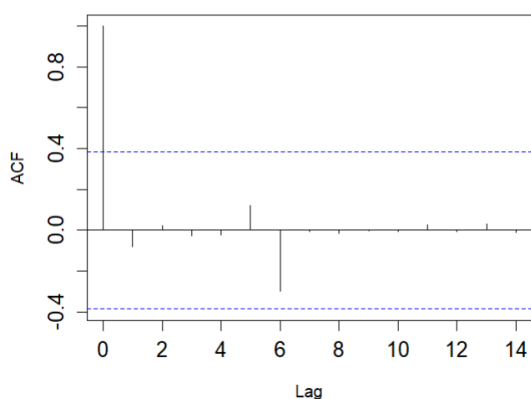


rys. 5.7

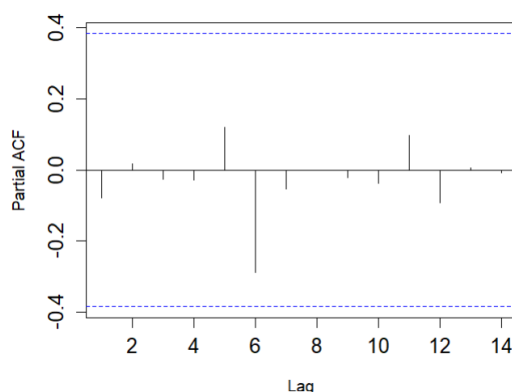


rys. 5.8

Wygląd wykresów ACF (rys. 5.9) i PACF (rys. 5.10) świadczy o braku sezonowości i trendu.



rys. 5.9



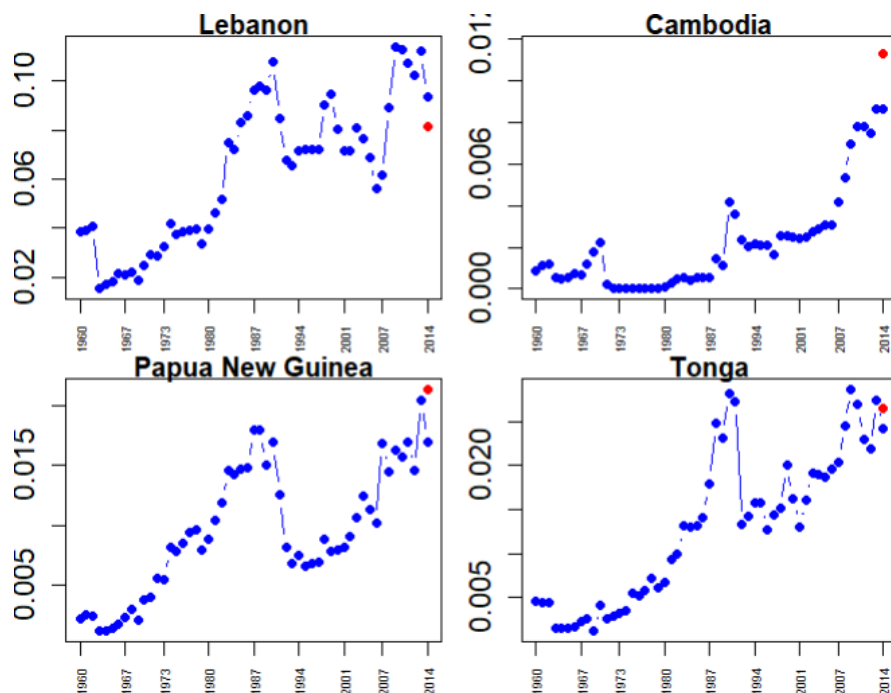
rys. 5.10

Następnie sprawdzono model na zbiorze przeznaczonym do predykcji zawierającym 4 obserwacje. Na rys. 3.11 porównano wartość predykowaną przez model do rzeczywistych wartości dla każdego państwa.

	Liban	Kamboża	Papua Nowa Gwinea	Tonga
real	0.0936	0.0086	0.0169	0.0242
predicted	0.0813	0.0112	0.0213	0.0265

rys. 5.11

Ponownie przygotowano wykres na którym wizualnie można stwierdzić jak dobrze wartości zostały przewidziane przez model. W przypadku Libanu i Tongy wartości zostały przewidziane stosunkowo poprawnie, dużo gorzej jednak oszacowano wartości dla Kambodży i Papui Nowej Gwinei. (rys. 5.12)



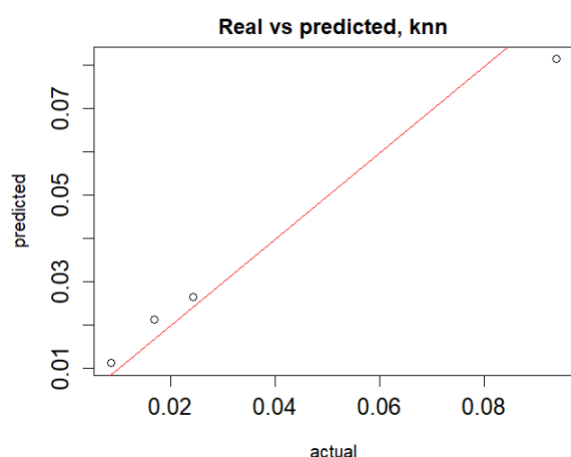
rys. 5.12

Wskaźnik R^2 ponownie jest bardzo wysoki, a MSE utrzymuje się na niskim poziomie. Warto zwrócić uwagę na MAPE, które spadło poniżej 19%. (rys. 5.13)

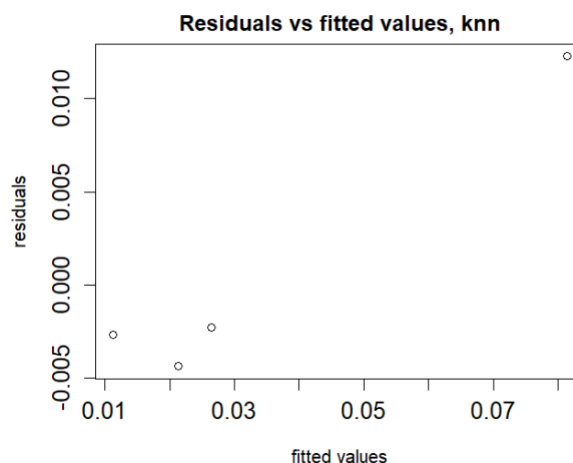
test: 4	R2	MSE	RMSE	MAPE
Knn	0.99	4.56E-05	0.0067	19.00%

rys. 5.13

Wszystkie wartości na wykresie porównującym rzeczywiste i przewidywane (rys. 5.14) wartości leżą blisko linii, a na wykresie rozrzutu reszt obserwacje są rozłożone losowo (rys. 5.15).

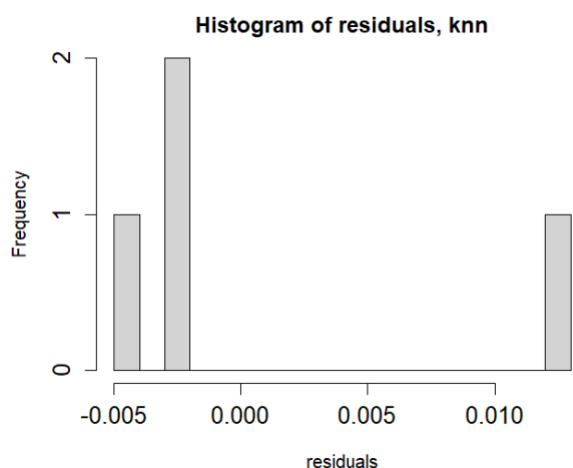


rys. 5.14

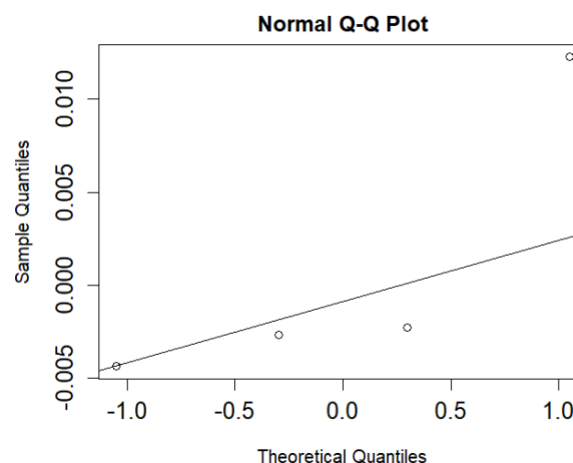


rys. 5.15

Histogram reszt (rys. 5.16) nie ma pożądanego kształtu rozkładu normalnego przez bardzo duże odsunięcie w prawo jednej wartości, obserwację tą potwierdza wykres porównania rozkładu reszt z rozkładem normalnym (rys. 5.17), gdzie wartości leżą w dość dużym oddaleniu od linii referencyjnej.

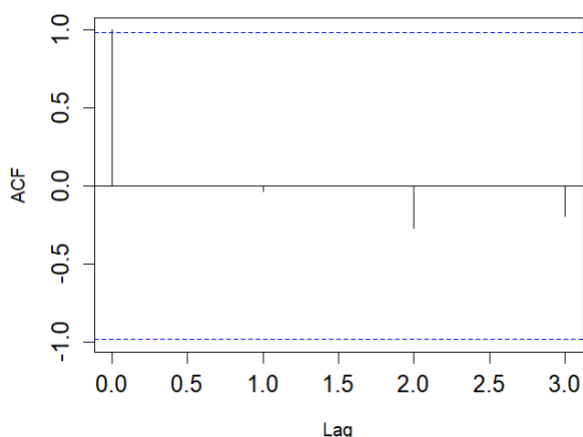


rys. 5.16

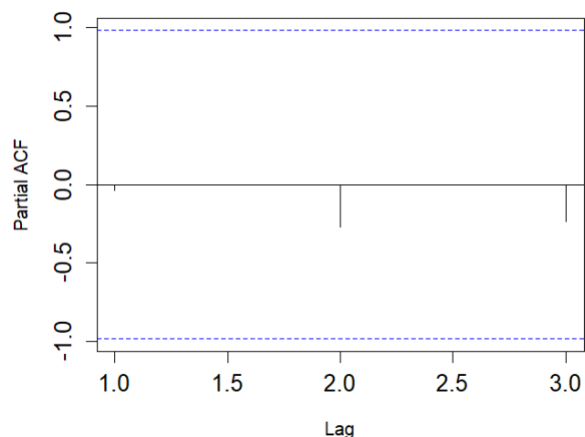


rys. 5.17

Tak jak w każdym poprzednim przypadku wykresy ACF (rys. 5.18) i PACF (rys. 5.19) informują o braku trendu i sezonowości.



rys. 5.18



rys. 5.19

6. Porównanie modeli

W celu porównania modeli wykonanych przy wykorzystaniu algorytmów Knn regresyjnego i sieci neuronowych MLP zestawiono dane z przetestowania modeli na zbiorze do predykcji zawierającym cztery obserwacje.

model	R2	MSE	RMSE	MAPE
Knn	0.99	4.61E-05	0.0067	35.25%
MLP	0.99	4.56E-05	0.0067	19.00%

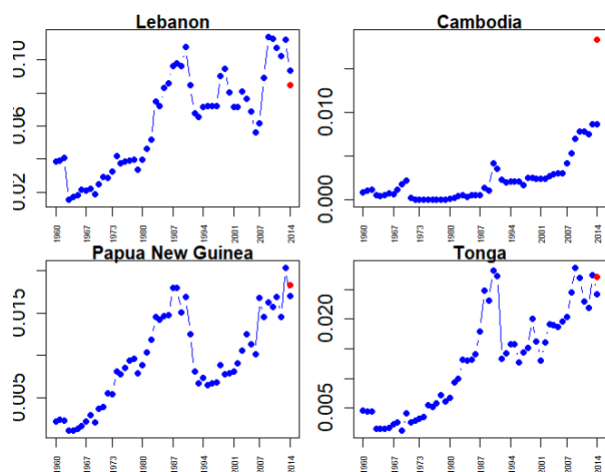
rys. 6.1

Wskaźnik R2 był identyczny dla obu modeli i osiągnął niezwykle wysoką wartość 0.99, co oznacza bardzo dobre dopasowanie modelu do danych. Również nie widać znaczącej różnicy między wskaźnikami MSE (i RMSE), w obu przypadkach były one niskie. Największą różnicę między modelami można zaobserwować porównując wartość MAPE, zdecydowanie niższe było ono dla modelu wykonanego przez MLP, aż o 16 punktów procentowych.

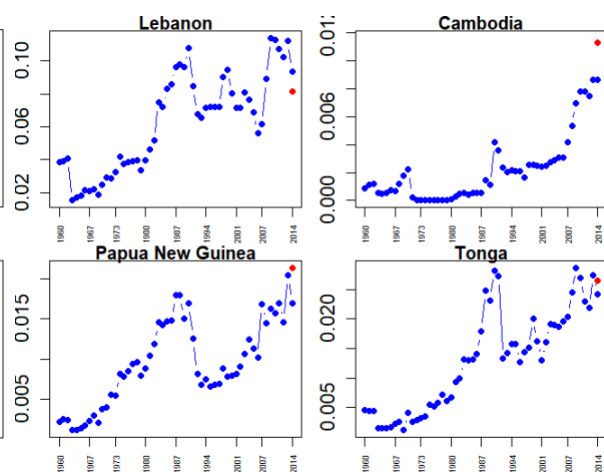
	Liban	Kamboża	Papua Nowa Gwinea	Tonga
real	0.0936	0.0086	0.0169	0.0242
Knn	0.0846	0.0183	0.0183	0.0270
MLP	0.0813	0.0112	0.0213	0.0265
różnica bzw. Knn	0.0090	0.0097	0.0014	0.0028
różnica bzw. MLP	0.0123	0.0026	0.0044	0.0023
zgodnie z kierunkiem Knn	tak	tak	tak	tak
zgodnie z kierunkiem MLP	tak	tak	nie	tak

rys. 6.2

Na rys. 6.2 porównano wartości przewidywane przez oba modele do wartości rzeczywistych, następnie obliczono wartości bezwzględne ich różnic. Model Knn okazał się być lepszy w pierwszym i trzecim przypadku, a model MLP w drugim i czwartym przypadku – nie można na tej podstawie wyciągnąć wniosków o większej skuteczności któregoś z modeli. Jednak kierunek zmian (tzn. wzrost lub spadek wartości w 2014 względem 2013) przewidział poprawnie w każdym przypadku Knn, a MLP pomylił się raz. Zmiana wartości rzeczywistych na przestrzeni lat oraz wartość przewidywana w 2014 zostały przedstawione wizualnie na rys. 6.3 i rys. 6.4.



rys. 6.3



rys. 6.4

7. Wniosek

Oba modele wykonane przez algorytmy Knn regresyjny i sieci neuronowych MLP osiągnęły podobną skuteczność w przewidywaniu wartości wskaźnika emisji CO₂ w zależności od państwa w 2014 roku.

Jako nieco lepszy można wskazać algorytm MLP, przemawiają za tym: niższa wartość wskaźnika MAPE i lepszy układ histogramu reszt oraz wykresu porównania rozkładu reszt z rozkładem normalnym.