

Projekt z Szeregów Czasowych

Jan Moskal i Szymon Makulec

2025-01-18

Wstęp

Dane, które będziemy analizować, pochodzą ze strony Głównego Urzędu Statystycznego (<https://bdl.stat.gov.pl/bdl/dane/podgrup/temat>) znajdują się w grupie “Przeciętne ceny detaliczne towarów i usług konsumpcyjnych”, podgrupie “Ceny detaliczne wybranych towarów i usług konsumpcyjnych (dane miesięczne)” i dotyczą cen węgla kamiennego za toną. Dane o przeciętnych cenach obejmują notowania co miesiąc dla całej Polski. Projekt ma na celu analizę tego szeregu czasowego, aby zrozumieć zmiany cen węgla kamiennego w Polsce w latach 2006-2019 i stworzyć prognozy na przyszłość.

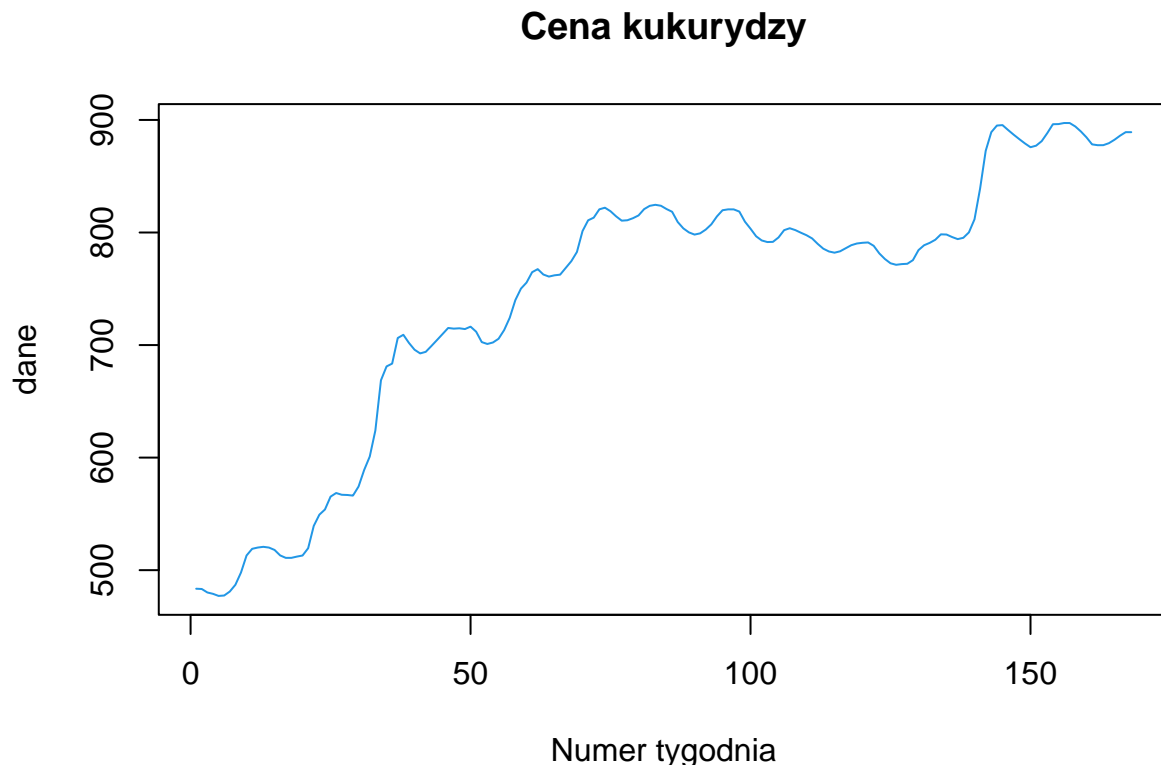
Wczytywanie danych

Zamieniamy wektor w macierz, aby ustawić dobrą kolejność danych (nawet mamy dane wypisane, w ten sposób, że jeden miesiąc dla czternastu lat i dopiero następny miesiąc, a chcemy żeby było chronologicznie)

Przekształcenie macierzy w wektor czytany kolumnowo (od góry do dołu). W ten sposób otrzymujemy dane w odpowiedniej kolejności.

Wstępna analiza szeregu

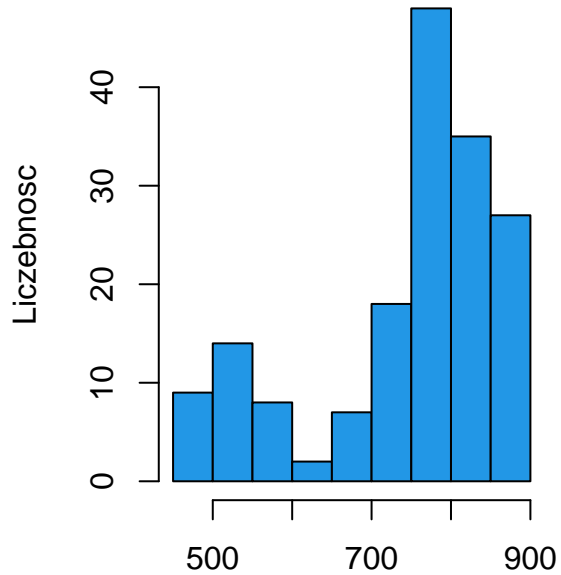
Wykres liniowy dla naszych danych w czasie t .



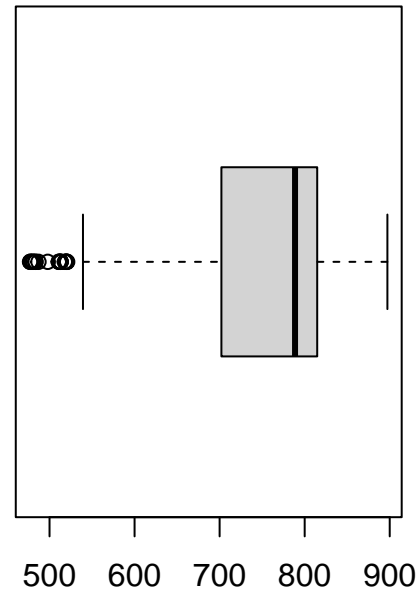
Jak widzimy z wykresu, cena dość szybko wzrosła do cen powyżej 700 zł. Widzimy również, że ogólny trend jest rosnący.

Robimy wykresy typu boxplot oraz histogram, żeby zobaczyć rozkład danych.

Histogram cen węgla kamiennego



Wykres ramka-wasy



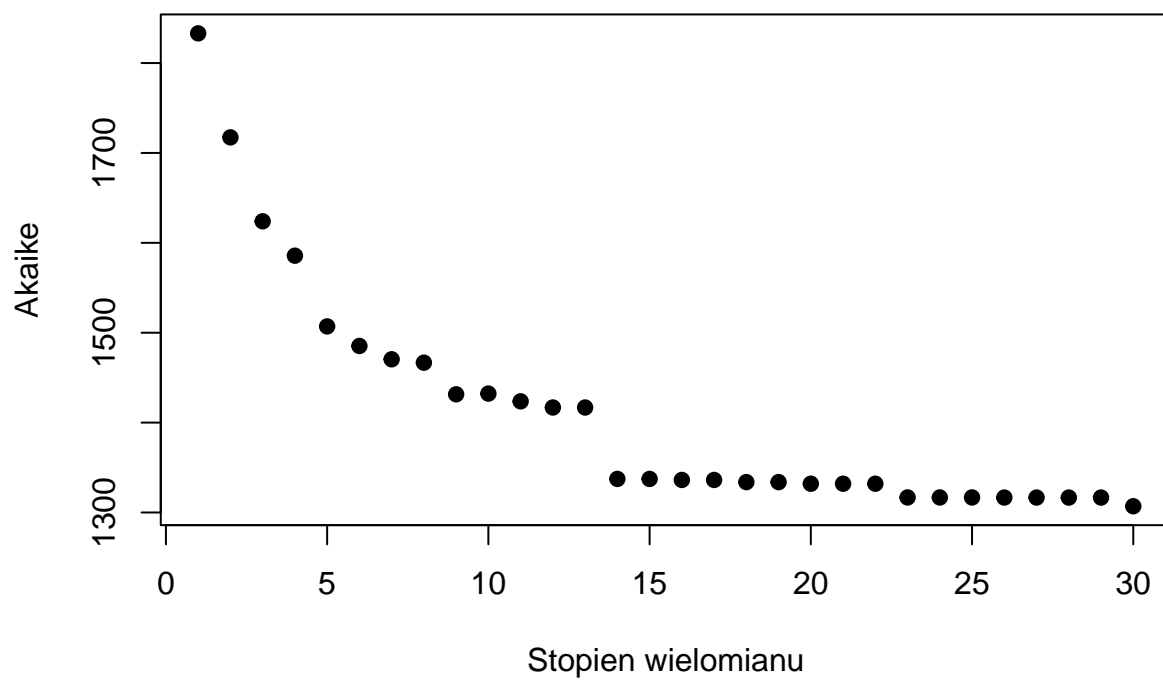
Możemy zauważyć, że rozkład jest lewostronnie asymetryczny, bierzemy się to z tego co już zauważyliśmy z wykresu liniowego czyli, że ceny od 700 zł za tonę zaczęły się już po 2 latach od pierwszej obserwacji z szeregu a pozostałe 12 lat oscylowało co do wartości od 700 do 900 zł za tonę. Z wykresu pudełkowego możemy zauważyć nawet dokładniej, że kwartył pierwszy wynosi około 700 a kwartył trzeci około 810 co w przełożeniu na nasz problem oznacza, że połowa obserwacji, czyli z 7 lat znajduje się na tym małym przedziale.

Podstawowe statystyki

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	477.1	702.2	788.7	744.1	814.5	897.3

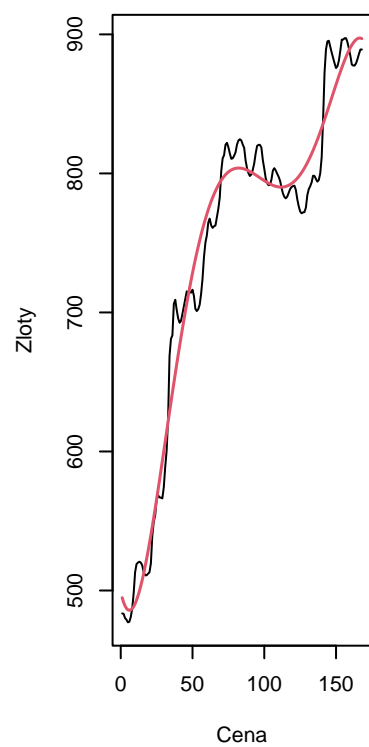
Szukamy najlepszego wielomianu opisującego nasz szereg

Kryterium AIC dla wielomianu stopnia i

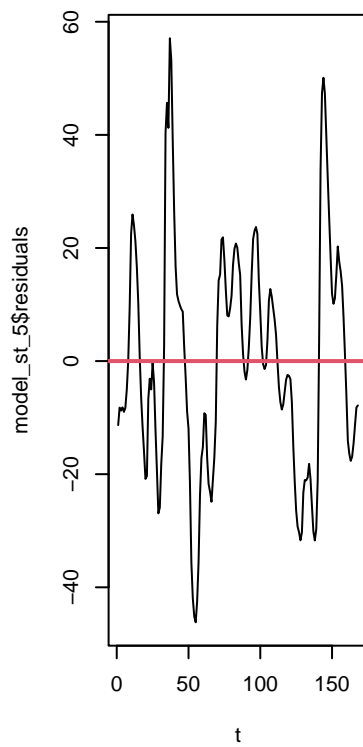


Z kryterium osuwiska wybieramy wielomian stopnia 5.

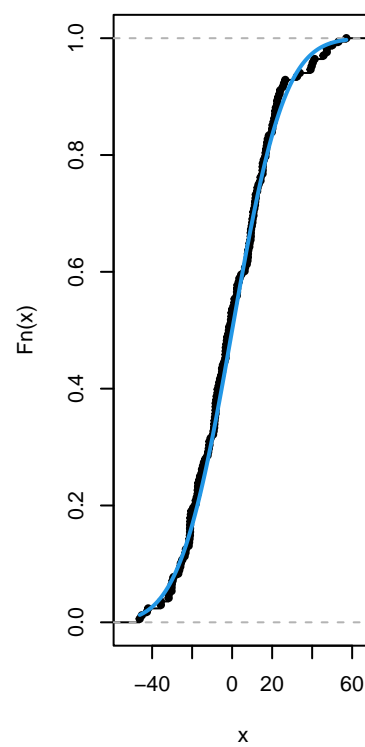
Dopasowanie wiel. st.: 5



Reszty



Dystrybuanta



Badanie reszt

Zbadamy reszty z modelu stopnia 5.

```
##
##  Runs Test
##
## data:  reszty
## statistic = -11.143, runs = 13, n1 = 84, n2 = 84, n = 168, p-value <
## 2.2e-16
## alternative hypothesis: nonrandomness

##
##  Asymptotic one-sample Kolmogorov-Smirnov test
##
## data:  reszty
## D = 0.48095, p-value < 2.2e-16
## alternative hypothesis: two-sided

##
##  Lilliefors (Kolmogorov-Smirnov) normality test
##
## data:  reszty
## D = 0.035656, p-value = 0.8673

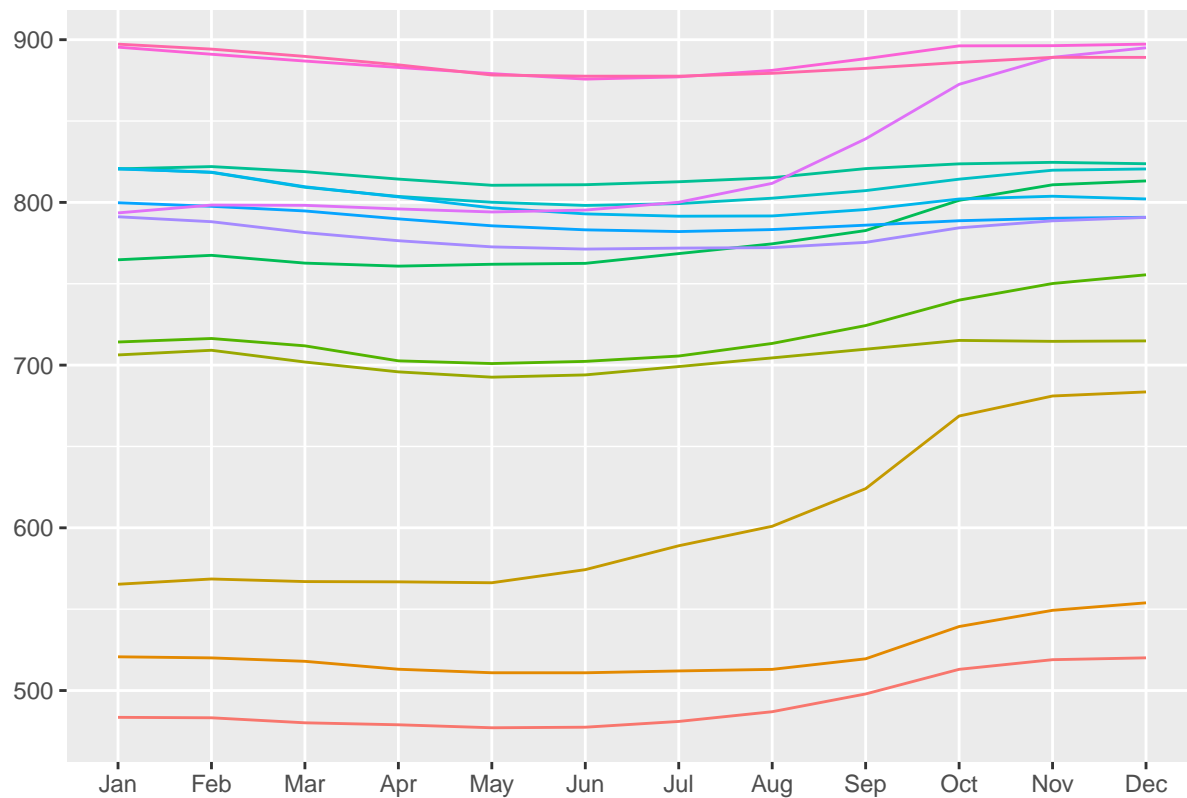
##
##  Shapiro-Wilk normality test
##
## data:  reszty
## W = 0.98779, p-value = 0.1533

##
##  Box-Ljung test
##
## data:  reszty
## X-squared = 390.21, df = 12, p-value < 2.2e-16
```

Odrzucamy hipotezy o losowości reszt, o ich średniej w zerze, o ich normalności oraz o ich nieskorelowaniu.

Analiza trendów fazowych

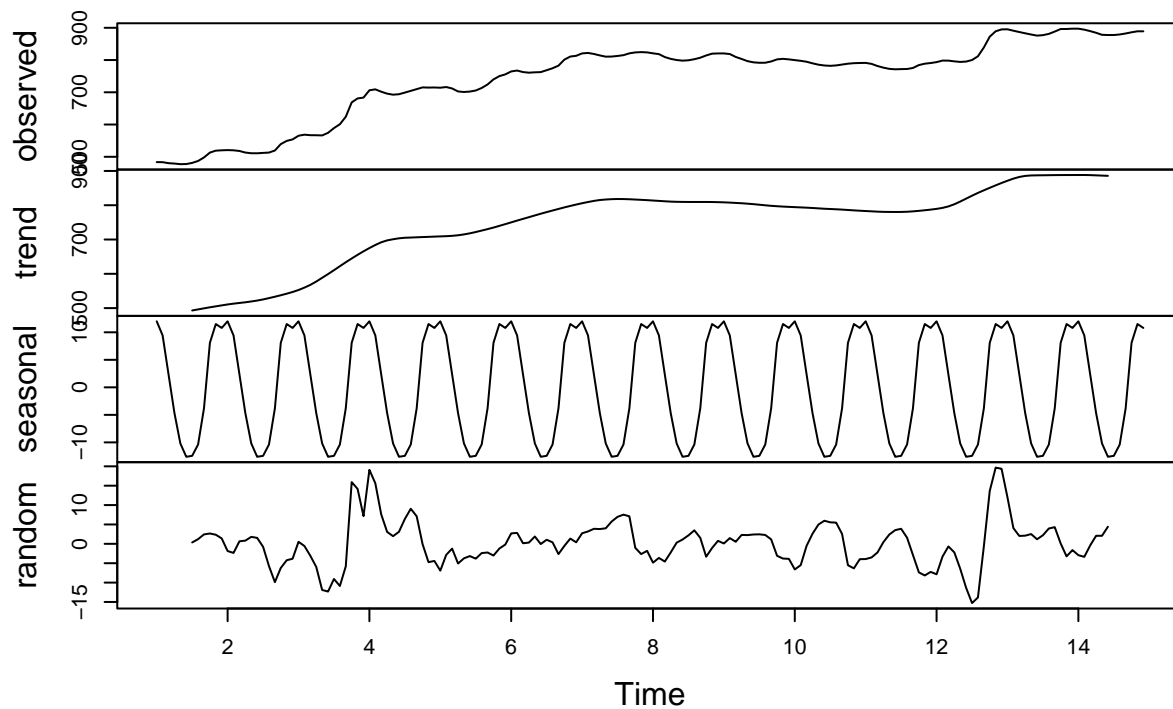
Wykres sezonowości dla lat 2006–2019



Z wykresu sezonowości widzimy powtarzający się trend wzrostu cen węgla kamiennego w okresie od sierpnia do listopada. W okresie od stycznia do maja zauważalny jest nieznaczny trend spadkowy cen.

Dekompozycja

Decomposition of additive time series



Po wykonaniu dekompozycji widzimy, że trend jest rosnący w dziedzinie. Widzimy również, że występuje sezonowość w częstotliwości 12 miesięcznej. W kwestii reszt wydają się one oscylować wokół zera, jednak przez ich nieregularność będzie się trzeba im lepiej przyjrzeć.

Stacjonarność szeregu

```
##
## Augmented Dickey-Fuller Test
##
## data: dane
## Dickey-Fuller = -1.7305, Lag order = 5, p-value = 0.6888
## alternative hypothesis: stationary

## Warning in kpss.test(dane): p-value smaller than printed p-value

##
## KPSS Test for Level Stationarity
##
## data: dane
## KPSS Level = 2.7326, Truncation lag parameter = 4, p-value = 0.01

##
## Phillips-Perron Unit Root Test
##
## data: dane
## Dickey-Fuller Z(alpha) = -4.5142, Truncation lag parameter = 4, p-value
## = 0.856
## alternative hypothesis: stationary
```

Z wszystkich testów wynika, że szereg jest niestacjonarny.

```
## Series: dane
## ARIMA(4,1,0) with drift
##
## Coefficients:
##          ar1      ar2      ar3      ar4      drift
##          0.9059  -0.4425  0.5447  -0.4685  2.4353
## s.e.    0.0676   0.0881  0.0872   0.0670  0.7722
##
## sigma^2 = 21.63: log likelihood = -491.99
## AIC=995.97   AICc=996.5   BIC=1014.68
```

Model ARIMA(4,1,0) wskazuje, że szereg czasowy wymagał różnicowania pierwszego rzędu, aby stać się stacjonarny. Współczynniki autoregresyjne sugerują zależność od czterech poprzednich wartości, a obecność dryfu oznacza trend wzrostowy. Nie mamy składnika średniej ruchomej. Teraz przejdziemy do zbadania reszt.

```
##
## Runs Test
##
## data:  reszty
## statistic = -0.31874, runs = 82, n1 = 75, n2 = 93, n = 168, p-value =
## 0.7499
## alternative hypothesis: nonrandomness
```

Odrzucamy hipotezę o losowości reszt.

```
##
## One Sample t-test
##
## data:  reszty
## t = -0.016825, df = 167, p-value = 0.9866
## alternative hypothesis: true mean is not equal to 0
## 95 percent confidence interval:
## -0.7036653  0.6917731
## sample estimates:
## mean of x
## -0.00594608
```

Nie odrzucamy hipotezy o średniej równej 0.

```
##
## Lilliefors (Kolmogorov-Smirnov) normality test
##
## data:  reszty
## D = 0.10715, p-value = 6.983e-05
```

```
##
## Shapiro-Wilk normality test
##
## data:  reszty
## W = 0.89614, p-value = 1.765e-09
```

```
##
## Box-Ljung test
##
## data:  reszty
## X-squared = 38.365, df = 12, p-value = 0.0001338
```

Z box.testu odrzucamy hipotezę o braku korelacji w resztach.