

Wykorzystanie poznanych metod dotyczących analizy zależności liniowej do wybranych danych rzeczywistych

Autorzy

Kacper Budnik, 262286
Maciej Karczewski, 262282



Politechnika Wrocławska

Wydział Matematyki 22 grudnia 2022

Spis treści

1	Wprowadzenie	2
2	Transformacja danych	2
2.1	Dodanie nowej kolumny	2
2.2	Usunięcie wartości odstających	2
3	Jednowymiarowa analiza	3
3.1	Objętość	3
3.2	Cena	4
4	Analiza zależności między ceną, a objętością diamentu	6
4.1	Estymacja punktowa	7
4.2	Estymacja przedziałowa	8
4.3	Predykcja danych	8
5	Analiza residuów	9
6	Wnioski autorów	12

1 Wprowadzenie

Diamenty są krystalicznym węglem, najtwardszym znanym materiałem na Ziemi. Są one tworzone w głębi Ziemi, w wysokich temperaturach i ciśnieniach, a ich wydobywanie wymaga specjalistycznej wiedzy i wielu lat doświadczenia. Diamenty są najcenniejszymi kamieniami szlachetnymi na świecie i są cenione zarówno przez osoby prywatne, jak i przez przemysł jubilerski. Błask i niezwykła trwałość diamentów sprawiają, że są one bardzo cenione, a ich cena zależy od wielu czynników, takich jak kolor, czystość, kształt i wielkość. Najcenniejsze diamenty są przezroczyste i bezbarwne, choć mogą też być innego koloru np. żółte, niebieskie, czerwone czy zielone.

Diamenty są często wykorzystywane do produkcji biżuterii. Są również używane w narzędziach do cięcia i szlifowania oraz w elektronice. W ostatnich latach diamenty zaczęły być również używane jako narzędzia do badań naukowych, ponieważ ich niezwykła trwałość i twardość sprawiają, że są one idealnymi materiałami do wielu zastosowań.

Dane o diamentach pozyskujemy z platformy Kaggle. Nasze dane zawierają 53940 rekordów i pierwsze 7 wierszy wygląda następująco:

	carat	cut	color	clarity	depth	table	price	x	y	z
1	0.23	Ideal	E	SI2	61.5	55	326	3.95	3.98	2.43
2	0.21	Premium	E	SI1	59.8	61	326	3.89	3.84	2.31
3	0.23	Good	E	VS1	56.9	65	327	4.05	4.07	2.31
4	0.29	Premium	I	VS2	62.4	58	334	4.2	4.23	2.63
5	0.31	Good	J	SI2	63.3	58	335	4.34	4.35	2.75
6	0.24	Very Good	J	VVS2	62.8	57	336	3.94	3.96	2.48
7	0.24	Very Good	I	VVS1	62.3	57	336	3.95	3.98	2.47

Tabela 1: Oryginalne dane.

Dla naszych danych sprawdzimy liniową zależność ceny ("price") od iloczynu wymiarów diamentu.

2 Transformacja danych

2.1 Dodanie nowej kolumny

Do naszych danych dodamy nową kolumną i nazwiemy ją "V". Nowa kolumna będzie iloczynem wymiarów diamentu (kolumny "x", "y" oraz "z"). Ten iloczyn, ponieważ rozprzyskane diamenty mają podobne kształty, możemy interpretować jako ich objętość, z dokładnością do przemnożonej stałej. Nową zmienną będziemy wyjaśniać cenę diamentu. W dalszej części raportu tą zmienną będziemy nazywać objętością lub iloczynem wymiarów.

2.2 Usunięcie wartości odstających

Przed przystąpieniem do analizy danych, oczyściliśmy je, poprzez usunięcie obserwacji odstających z wykorzystywanych danych. W celu klasyfikacji tych obserwacji, wykorzystaliśmy znane nam metody. Za obserwacje odstające uznaliśmy dane, które

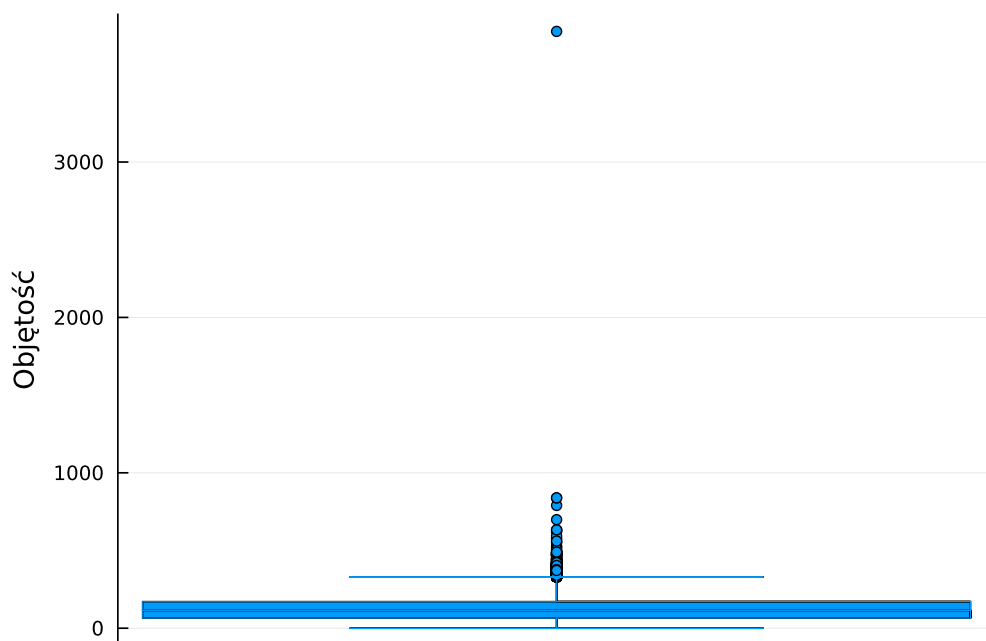
$$y \notin (Q_1 - IQR, Q_3 + IQR),$$

gdzie y to kandydat na obserwacje odstającą, Q_1 , Q_3 to odpowiednio pierwszy i trzeci kwantyl, a IQR jest rozstęp międzykwantylowym. Po tej transformacji odrzuciliśmy blisko 10% danych.

3 Jednowymiarowa analiza

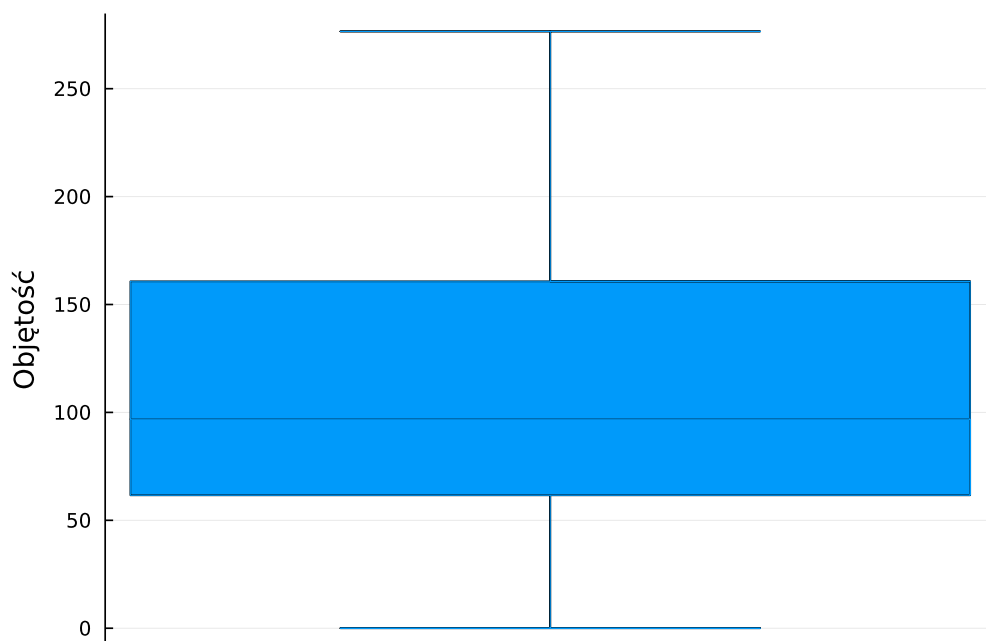
3.1 Objętość

Analizę danych zaczniemy od analizy objętości. Box plot objętości wygląda następująco.



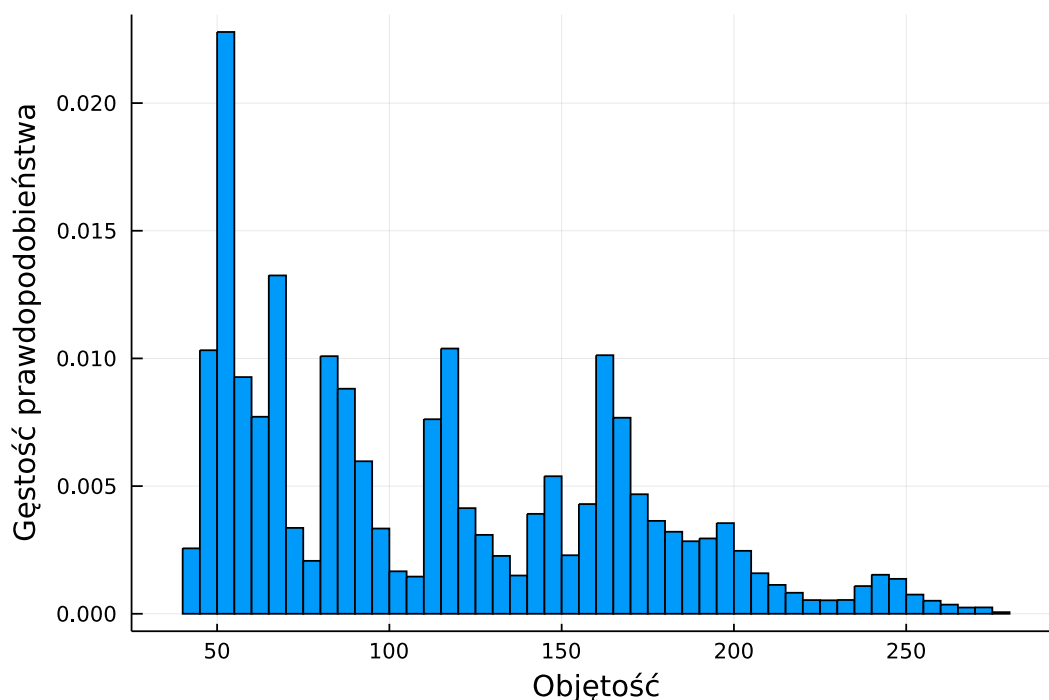
Rysunek 1: Boxplot objętości dla podstawowych danych.

Jak możemy zobaczyć na box plocie nasze dane mają dużo wartości odstających. Z tego powodu ciężko analizować wykres 1, a także stworzyć poprawny model, dlatego pozbędziemy się ich. Po usunięciu obserwacji odstających otrzymamy czytelniejszy wykres.



Rysunek 2: Box plot objętości dla oczyszczonych danych.

Tutaj nasze dane są o wiele bardziej czytelne. Z wykresu 2 możemy odczytać, że mediana naszych danych wynosi około 95 oraz, że mamy do czynienia z rozkładem prawostronnie skośnym. Widzimy też, że pierwszy kwantyl wynosi około 60 a trzeci około 160. Zobaczmy teraz jak wygląda rozkład objętości na histogramie.



Rysunek 3: Histogram objętości dla oczyszczonych danych.

Jak widzimy na histogramie gęstość nie zachowuje się monotonicznie, ciężko stwierdzić jaki rozkład ma objętość. Pomimo tego faktu możemy zauważyć, że diamenty dzielą się głównie na 4 grupy. Pierwszą, najliczniejszą, jest grupa o iloczynie wymiarów należących do przedziału $[45; 70]$, następną jest grupa o iloczynie należącym do $[80; 100]$. Pozostałymi grupami są iloczyny należące do przedziałów $[110; 120]$ oraz $[160; 170]$.

Zobaczmy teraz jak wyglądają statystyki opisowe dla naszych oczyszczonych danych.

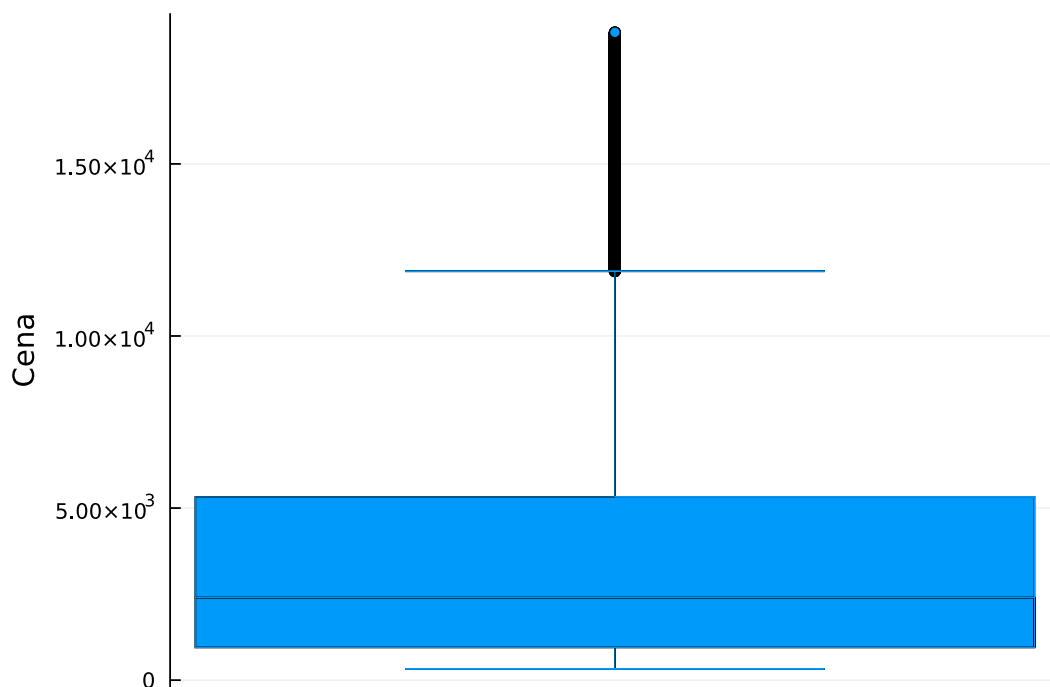
	Średnia	Mediana	Wariancja	Skośność	Kurtoza
V	113.48	100.87	3057.10	0.60	-0.63

Tabela 2: Podstawowe statystyki opisowe dla objętości.

Możemy zobaczyć, że średnia jest większa od mediany, oraz że skośność jest dodatnia co sugeruje nam, że rozkład objętości jest prawostronnie skośny, co się pokrywa za informacją zawartą na box plocie 2. Kurtoza jest ujemna to znaczy, że rozkład ma ogony węższe niż rozkład normalny czyli jest rozkładem platykurtycznym.

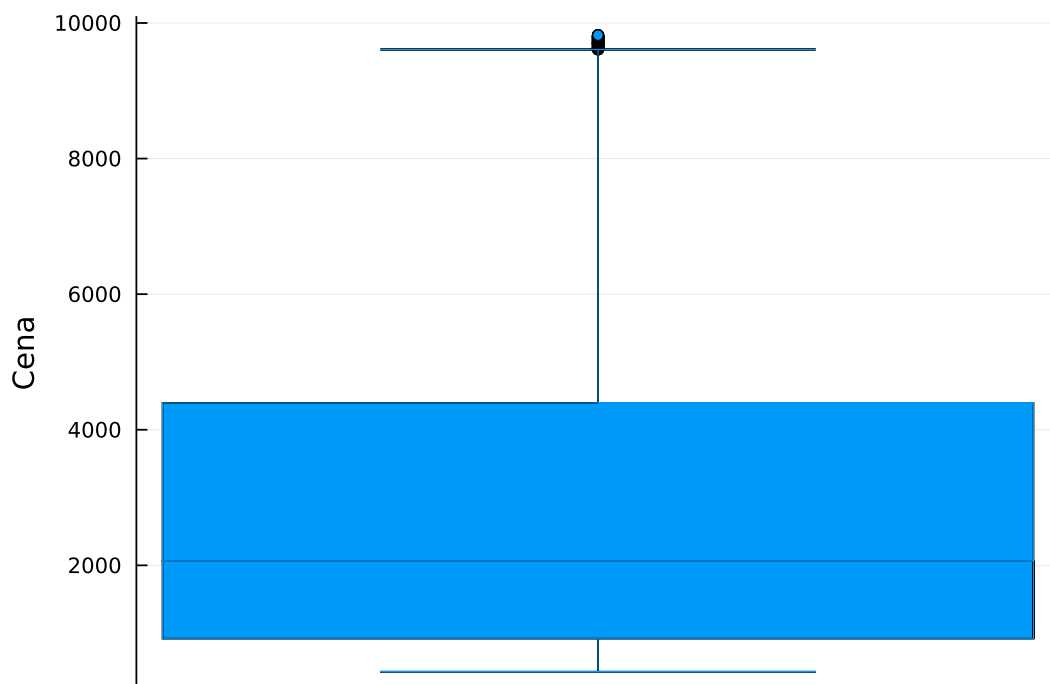
3.2 Cena

Sprawdziliśmy jak wygląda rozkład objętości to teraz sprawdzimy rozkład ceny diamentów. Box Plot dla ceny diamentów bez usunięcia wartości skrajnych wygląda następująco.



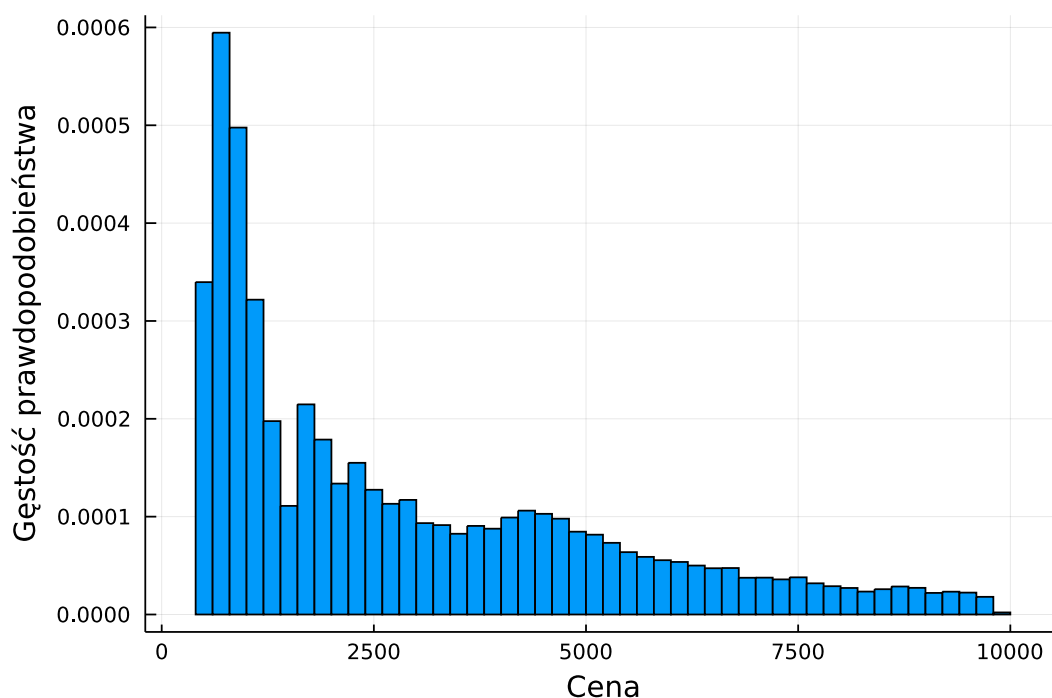
Rysunek 4: Boxplot ceny dla podstawowych danych.

Z box plotu możemy odczytać, że ponownie mamy dużo wartości odstających, które przeszkadzają w analizie danych. Właśnie z tego powodu usuniemy 1% najmniejszych cen oraz 10% największych cen. Po usunięciu wartości nasz Box plot wygląda następująco



Rysunek 5: Box plot ceny dla oczyszczonych danych.

Po usunięciu danych skrajnych wykres pudełkowy staje się czytelniejszy. Z wykresu możemy odczytać, że mediana cen jest w okolicy 2000, natomiast pierwszy w okolicach 1000, a trzeci w okolicach 4300. Możemy sądzić, że rozkład ceny jest prawostronnie skośny.



Rysunek 6: Histogram ceny diamentów.

Z histogramu możemy odczytać, że gęstość prawdopodobieństwa ceny jest prawie monotoniczna. Im większa cena tym na ogół, rzadszy diament. Możemy zobaczyć, że mamy najwięcej diamentów tanich nieprzekraczających 2500. Możemy także zobaczyć większą liczebność diamentów z ceną w okolicach 4000.

Zobaczmy teraz jak wyglądają statystyki opisowe dla naszych oczyszczonych danych

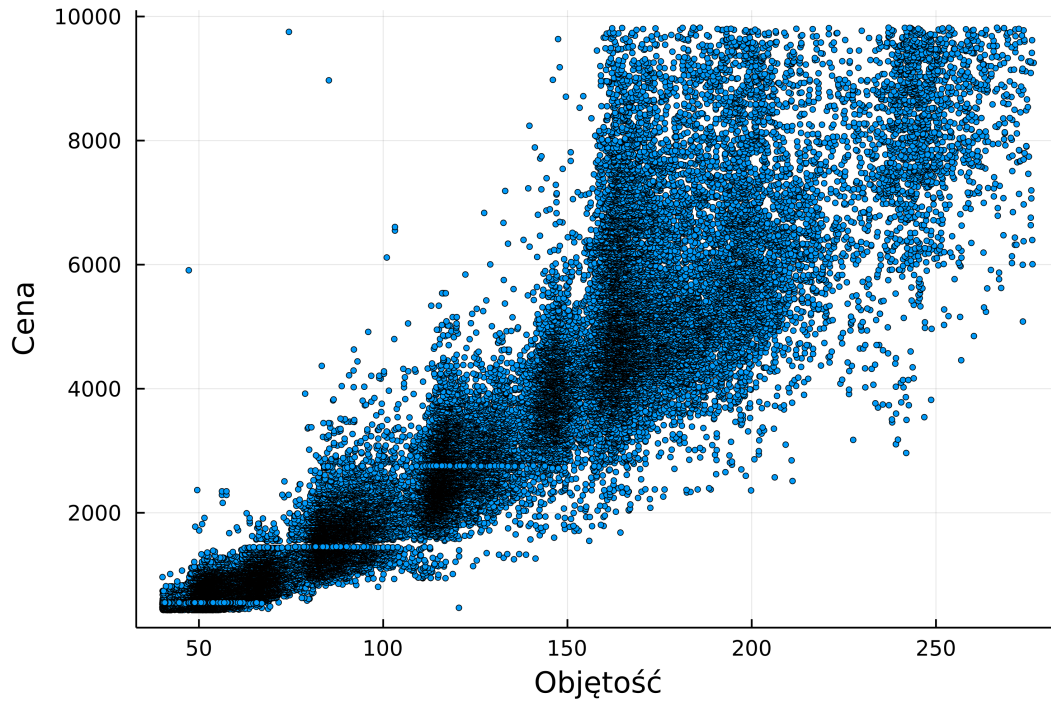
	Średnia	Mediana	Wariancja	Skośność	Kurtoza
Cena	2890.46	2064.0	$5.49 \cdot 10^6$	1.02	0.10

Tabela 3: Podstawowe statystyki opisowe dla objętości.

Możemy zobaczyć, że średnia jest większa od mediany, oraz że skośność jest dodatnia co sugeruje nam, że rozkład objętości jest prawostronnie skośny co się pokrywa za informacją zawartą na box plocie 5. Kurtoza jest dodatnia, ale bliska 0 to znaczy, że rozkład ma ogony grubsze niż rozkład normalny lub ma porównywalne do rozkładu normalnego. Zatem jest to rozkład leptokurtyczny lub mezokurtyczny. Wariancja ceny jest znacząco większa niż wariancja objętości w tabeli2.

4 Analiza zależności między ceną, a objętością diamentu

Po omówieniu danych oraz ich wstępnym przetworzeniu, możemy przystąpić do analizy. Celem naszym będzie przeanalizować zależność między ceną, a objętością oraz dopasowanie krzywej regresji. Zanim rozpoczniemy analizę, w celu określenia rodzaju zależności, przyjrzyjmy się wykresowi zależności między danymi.



Rysunek 7: Wykres zależności ceny od objętość.

Na wykresie możemy zauważyć silną zależność między naszymi danymi. Z powodu rozłożenia naszych danych, zależność ta może być liniowa. W celu określenia tej zależności obliczymy podstawowe statystyki

- $SST = \sum_{i=1}^n (y_i - \bar{y})^2 \approx 1.409e11$,
- $SSE = \sum_{i=1}^n (\hat{y}_i - y_i)^2 \approx 2.015e10$,
- $SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \approx 1.207e11$,

oraz, najbardziej nas interesujący, współczynnik korelacji Pearsona

$$\rho = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} \approx 0.93,$$

gdzie x_i , y_i to i -te obserwacje odpowiednio objętości oraz ceny, \bar{x} oznacza średnią z x , a n to rozmiar próby, oznaczeniami tymi będziemy posługiwali się w całym raporcie. Wartość ta jest blisko wartości 1, zatem nasze dane są silnie skorelowane dodatnie oraz liniowo. Dodatkowo wartości SST oraz SSR są relatywnie blisko siebie, a suma błędów SSE jest rzęd mniejsza od pozostałych. W tym przypadku regresja liniowa powinna być odpowiednim wyborem. Dlatego, model którym będziemy opisywać dane będzie miał postać

$$Y_i = \beta_0 + \beta_1 x_i + \xi_i, \quad (1)$$

gdzie ξ_i jest losowym błędem pomiarowym. Y_i jest zmienną losową, której realizacją jest zaobserwowana wartość y_i . Zgodnie z wytycznymi dotyczącymi naszego zadania, zakładamy, że wszystkie zminne losowe ξ_i są iid. o rozkładzie normalny ze średnią 0. Naszym celem będzie estymować wartość \hat{y}_i poprzez estymowanie realizacji zmiennych losowych $\hat{\beta}_0$ oraz $\hat{\beta}_1$.

4.1 Estymacja punktowa

By stworzyć estymator ceny \hat{y} skorzystamy z metody najmniejszych kwadratów. Do estymacji parametrów β_0 oraz β_1 , występujące we wzorze (1), wykorzystamy losowo wybrane 80% naszych danych. Pozostałe 20% posłuży nam w celu sprawdzenia poprawności modelu. Dane zostały wylosowane przy użyciu

podstawowych funkcji języka Julia z ustalonym ziarnem. Estymowane parametry, w zaobserwowanej realizacji $Y = (Y_1, Y_2, \dots, Y_n)$, mają wartość

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \approx 39.089 \quad \text{oraz} \quad \hat{\beta}_0 = \bar{y} - \beta_1 \bar{x} \approx -1548.34.$$

Zatem estymowana wartość ceny \hat{y} modelu (1) będzie miała postać

$$\hat{y}_i = 38.089 \cdot x_i - 1548.34. \quad (2)$$

4.2 Estymacja przedziałowa

W tym przypadku będziemy szukać przedziału, w którym będzie należał nasz Y_i z dużym prawdopodobieństwem. Będziemy chcieli by prawdopodobieństwo to wynosiło $1 - \alpha$. Znany jest fakt, że w modelu (1) zmienna \hat{Y}_i ma poniższe parametry

$$\mathbb{E}\hat{Y}_i = \beta_0 + \beta_1 x_i \quad \text{oraz} \quad \text{Var}\hat{Y}_i = \text{Var}(\xi_1) \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right).$$

Jeśli $\text{Var}(\xi_1)$ nie jest znana w jej miejsce wstawiamy jej estymator

$$s^2 = \frac{1}{n-2} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

otrzymując estymator $\text{Var}\hat{Y}_i$. W tym przypadku zmienna losowa

$$\frac{\hat{Y}_i - \mathbb{E}\hat{Y}_i}{\sqrt{s^2 \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)}}$$

ma rozkład t-studenta z $n-2$ stopniami swobody. Przez $z_{\alpha/2}$ będziemy oznaczać $1-\alpha/2$ kwantyl z właśnie tego rozkładu. Dla estymowanej wartości \hat{Y}_0 , dla znanej objętości, równej x_0 będzie zachodziła równość

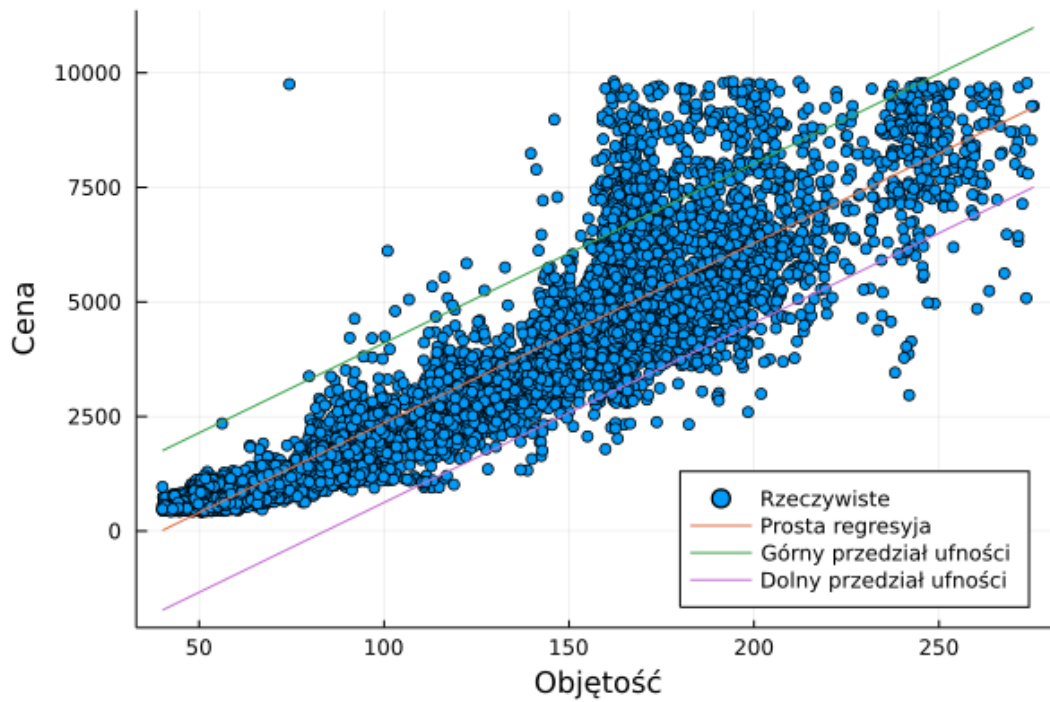
$$\mathbb{P} \left(\hat{Y}_0 - z_{\alpha/2} s \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}} \leq Y_0 \leq \hat{Y}_0 + z_{\alpha/2} s \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}} \right) = 1 - \alpha.$$

Zatem naszym przedziałem ufności ceny dla objętości x_0 , w zaobserwowanej próbie, będzie przedział

$$\left(\hat{y}_0 - z_{\alpha/2} s \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}, \hat{y}_0 + z_{\alpha/2} s \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}} \right). \quad (3)$$

4.3 Predykcja danych

Tak jak wspominaliśmy, 20% danych posłuży nam do sprawdzenia modelu. Dla tych danych przyjrzymy się wartościom rzeczywistym, estymowanym punktowo oraz przedziałowo.



Rysunek 8: Dane testowe.

Analizując wykres możemy zauważyć, że dane są w okolicach prostej regresyjnej (2) dla małych objętości. Wraz ze wzrostem iloczynu wymiarów rośnie też odległość danych od prostej. Do podobnych wniosków możemy dojść analizując przedziały ufności (3) dla naszych danych. Fakt ten może sugerować nam, że objętość nie wpływa liniowo na cenę diamentu lub, że założenia dotyczące modelu nie okazały się poprawne. Na poprawność regresji może wpływać również fakt, że diamenty mają jeszcze inne cechy, których nie bierzemy pod uwagę.

5 Analiza residuów

Podczas tworzenia modelu regresji liniowej oraz dalszych obliczeń zakładaliśmy następujące warunki

1. $\mathbb{E}\xi_i = 0 \forall i$,
2. $Var\xi_i = \sigma^2 < \infty \quad \forall i$,
3. ξ_i mają rozkład normalny,
4. $\xi_i \perp \xi_j$ dla $i \neq j$.

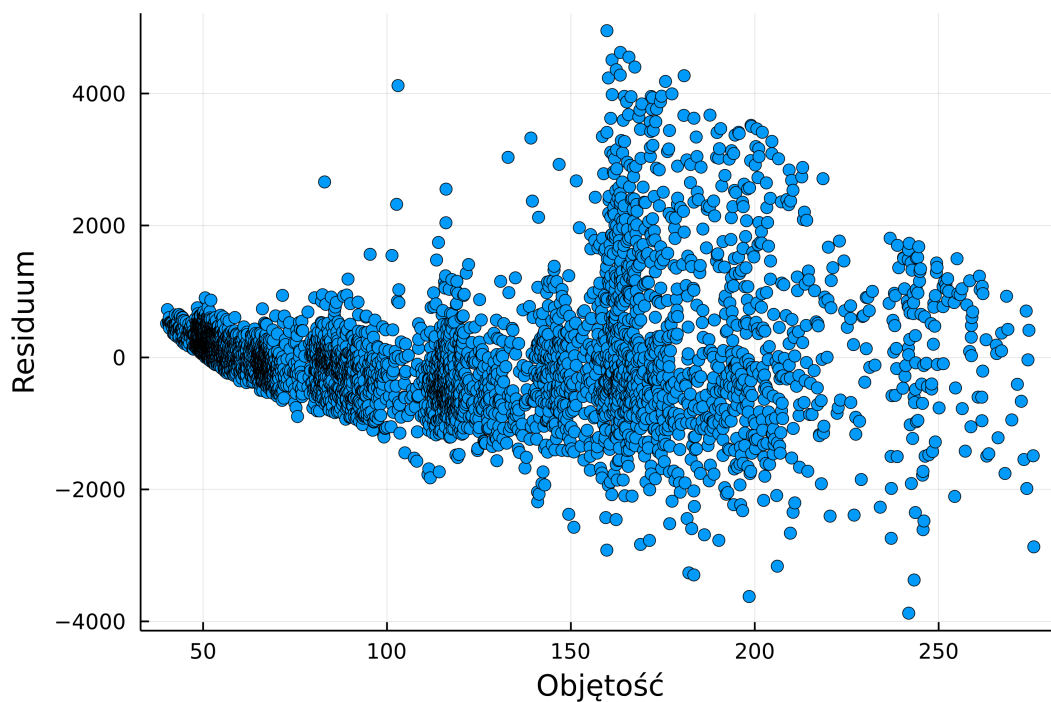
Do sprawdzenia tych warunków wykonamy analizę residuów e_i danych wzorem

$$e_i = y_i - \hat{y}_i$$

Zacniemy od sprawdzenia średniej rozkładu. W naszej próbie wynosi

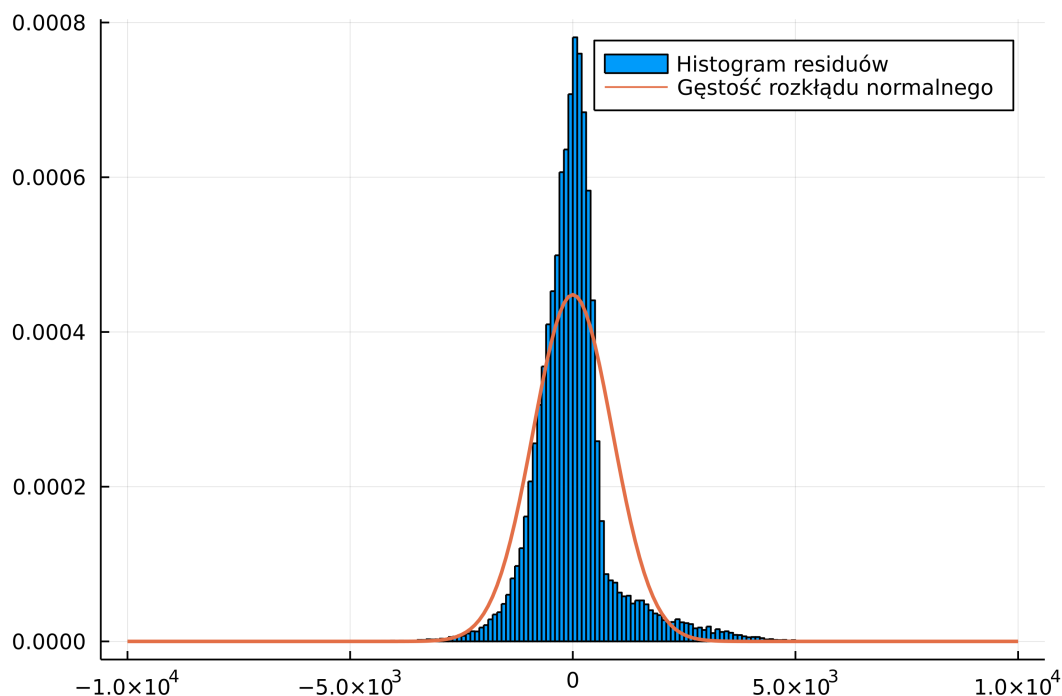
$$\mathbb{E}\xi_i \approx \bar{e} = 1.18e - 8,$$

zatem średnia residuów jest równa zero, co nie jest zaskakujące, uwzględniając, że skorzystaliśmy z metody najmniejszych kwadratów. Przeanalizujemy jeszcze wykres residuów na poniższym wykresie.



Rysunek 9: Zależność residuów od objętości.

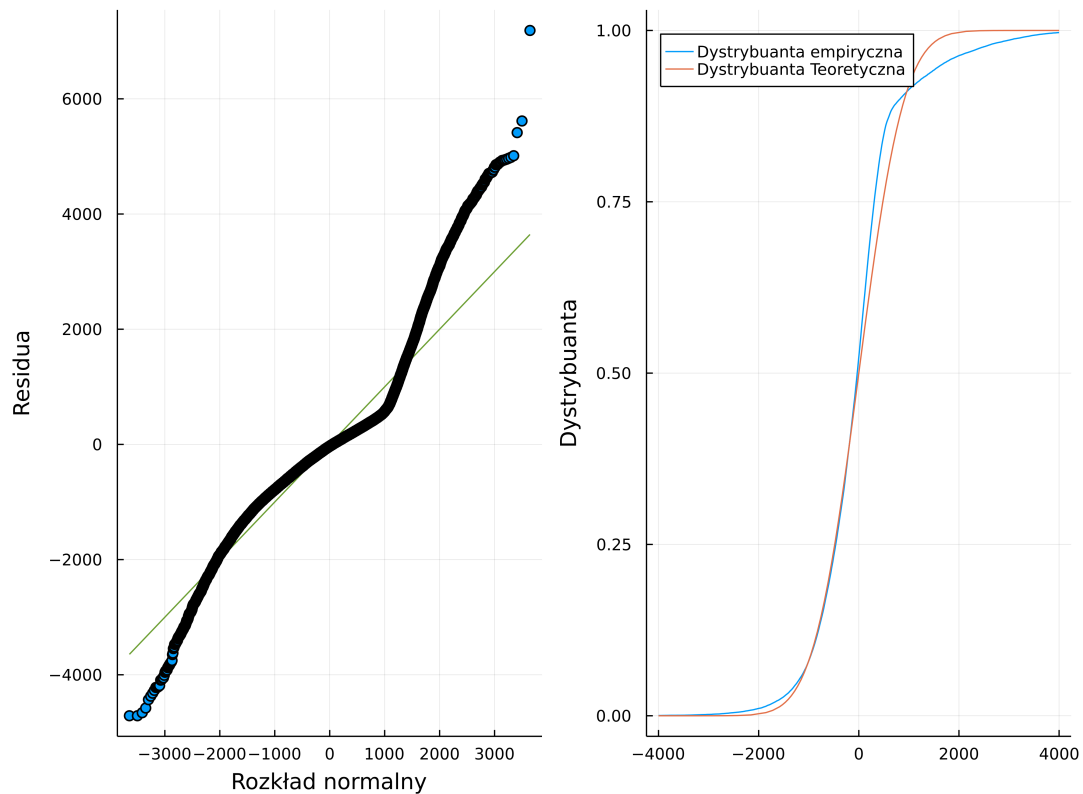
Na wykresie wyraźnie widać, że dla małych objętości wszystkie residua są znacznie większe od zera, więc ich wartość oczekiwana nie może być bliska zera. Dodatkowo możemy zauważyć, że wraz ze zwiększającą się objętością, zwiększa się wariancja. Zatem pierwsze, jak i drugie założenie, nie jest spełnione. Do sprawdzenia założenia trzeciego o normalności rozkładu sprawdzimy wykorzystując histogram residuów.



Rysunek 10: Histogram residuów z nałożoną gęstością rozkładu normalnego.

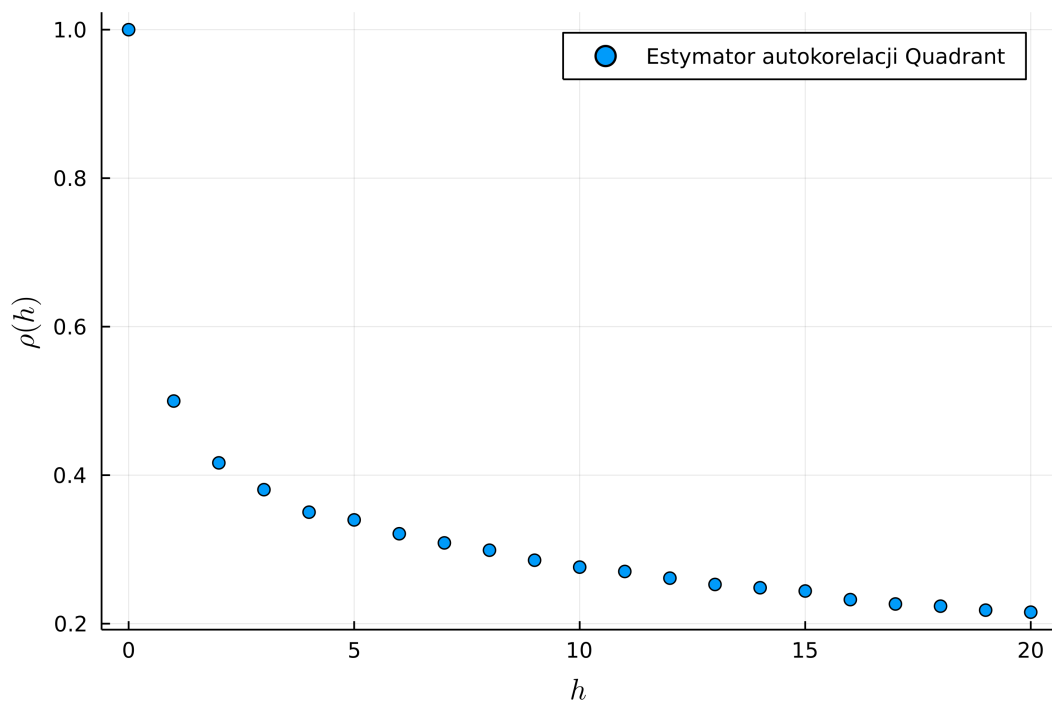
Na wykresie widzimy, że histogram nie pokrywa się z gęstością teoretyczną. Przed odrzuceniem założenia trzeciego, porównajmy jeszcze dystrybucję naszych residuów oraz rozkładu normalnego oraz oba

rozkłady na qqplocie.



Rysunek 11: Porównanie rozkładu residuów z rozkładem normalnym.

Również na tych wykresach widzimy znaczną różnicę między rozpatrywanymi rozkładami. Możemy więc z pewnością stwierdzić, że rozkład residuów nie jest normalny. Do sprawdzenia pozostaje nam założenie o niezależności błędów. W tym celu sprawdzimy, czy nasze zmienne są skorelowane.



Rysunek 12: Estymator funkcji autokowariancji dla residuów.

Wartości w estymatorze autokowariancji nie są zerowe na prawo od początku układu współrzędnych, więc nasze residua są skorelowane. Dlatego możemy wywnioskować, że nasze dane są zależne, zatem czwarte założenie również musimy odrzucić.

6 Wnioski autorów

Początkowo wydawało się, że to model liniowy jednej zmiennej może dobrze estymować cenę. Niestety wynik końcowy zaprzeczył temu. Pomimo nawet współczynnika korelacji Pearsona bliskiego wartości 1. Możemy wywnioskować stąd, że cena diamentów nie zależy liniowo od ich objętości, ale ma ona znaczący wpływ na ich wartość, w szczególności dla małych rozmiarów. Wraz z rozmiarem rosła wariancja drogocennego kruszcu, a cena była rozłożona bardziej nieprzewidywalnie. Naszą hipotezą jest, że planując zakup diamentu o małej cenie, głównym kryterium będzie jego rozmiar. Jeśli nasz budżet na zakup jest większy, rozmiar schodzi na kolejny plan, ustępując miejsca innym własnościom, takim jak na przykład przejrzystość, kolor, czy doskonałość cięcia.