

Ceny hoteli w Amsterdamie

Natalia Dyś

8 czerwca 2024

Spis treści

1	Wstęp	1
1.1	Cel raportu	1
1.2	Zbieranie danych	2
2	Wstępna analiza danych	2
2.1	ilość zebranych rekordów	2
2.2	zbierane wartości	2
2.3	usunięcie niepotrzebnych kolumn	3
2.4	brakujące dane	3
2.5	analiza danych	3
2.5.1	cena	3
2.5.2	ocena użytkowników	5
2.5.3	liczba ocen	7
2.5.4	liczba gwiazdek	8
2.5.5	liczba osób	9
2.5.6	data	10
2.5.7	odległość od centrum	12
2.5.8	śniadanie	14
2.5.9	macierz korelacji	15
3	Testowane modele	15
4	Wyniki i wnioski	16
4.1	propozycje rozwoju	17

1 Wstęp

1.1 Cel raportu

Raport ma na celu oszacowanie ceny per osoba jednonocnego pobytu w hotelu w Amsterdamie z użyciem metod regresji. Cena będzie szacowana na podstawie parametrów takich jak:

- liczba osób
- data pobytu
- ocena hotelu (opinie użytkowników oraz ich ilość)
- liczba gwiazdek
- dystans od centrum

1.2 Zbieranie danych

Analizowane przeze mnie dane zostały zebrane poprzez generowanie losowych requestów do strony booking.com, a następnie wydobycie poszczególnych ofert z html otrzymanego przy pomocy selenium i BeautifulSoup. Dane z każdego wyszukania są zapisywane do odrębnego pliku csv (umożliwia to szybka analiza mniejszego skrawka danych).

Losowo generowane parametry wyszukiwań to:

- liczba dorosłych osób - od 1 do 10
- data - od 30 do 400 dni od obecnej (w chwili wywołania scrapera) daty.

Ceny są zbierane i zapisywane w euro.

Zakładamy pobyt jednodniowy, ponieważ długość pobytu nie ma bezpośrednio wpływu na cenę, a dzięki temu maksymalizujemy ilość dostępnych obiektów.

Amsterdam jest arbitralnie wybrana lokalizacja. Jest on zapisywany w kolumnie 'city' aby działalność algorytmu można było potencjalnie rozszerzyć na inne miasta.

Nie uwzględniam dat z najbliższego miesiąca, ponieważ większość miejsc będzie już prawdopodobnie zajęta, co może generować niemiernodajne dane.

Zdecydowałam się też w moich wyszukaniach nie generować liczby dzieci (wynosi zawsze 0), jednak w kodzie istnieje taka możliwość.

2 Wstępna analiza danych

2.1 ilość zebranych rekordów

Ostatecznie dane zostały zebrane dla 400 requestów. Ilość listingów dla każdego wyszukania waha się od 26 do 485 (booking wyświetla tylko obiekty dostępne w danym terminie dla podanej liczby osób) Łączna ilość zebranych ofert to 42488, co daje średnio 106,22 per wyszukanie.

2.2 zbierane wartości

- name - nazwa hotelu
- city - miasto (parametr wyszukiwania)
- stars - liczba gwiazdek
- rating - ocena hotelu
- opinions - liczba ocen hotelu (czasami referowane jako reviews)
- distance_from_centre - odległość hotelu od centrum
- free_cancellation - czy istnieje możliwość darmowego odwołania rezerwacji (bool)
- breakfast - czy hotel oferuje śniadanie (bool)
- price - pełna cena rezerwacji
- price_per_person - cena rezerwacji podzielona przez ilość osób (dorośli + dzieci)
- date - data rezerwacji (parametr wyszukiwania)
- adults - liczba dorosłych (parametr wyszukiwania)
- children - liczba dzieci (parametr wyszukiwania)

2.3 usunięcie niepotrzebnych kolumn

Ponieważ dane zebrałam w tym przypadku tylko dla jednego miasta oraz nie uwzględniam udziału dzieci, odpowiednie kolumny zostały usunięte.

Nazwa hotelu nie będzie miała wpływu na jego cenę. Z późniejszej analizy wynikało, że darmowe odwołanie jest we wszystkich rekordach oznaczone jako prawda (prawdopodobnie przez lekko przekłamany sposób wyświetlania bookingu), więc możemy zignorować obie te kolumny.

W tym raporcie skupiam się na cenie per osoba, dlatego kolumna z łączną ceną rezerwacji również powinna zostać usunięta. (inaczej wyliczilibyśmy cenę na podstawie samej siebie)

2.4 brakujące dane

Column	Count Missing	Percentage Missing
rating	332	0.78
opinions	258	0.61

Brakujących danych jest na szczęście niewiele i znajdują się tylko w dwóch kolumnach, więc nie powinny znacząco wpłynąć na wyniki.

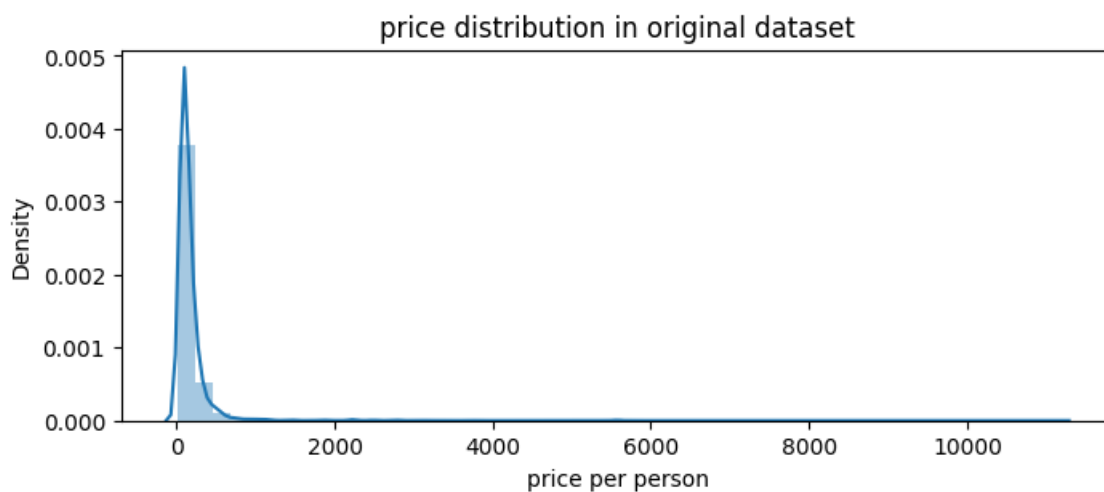
Według mojej analizy, liczby opinii będzie brakować kiedy jest ona równa 0 (booking nie wyświetla wtedy liczby, dlatego parser jej nie wyłapuje). Zakładam, że ten fakt jest odpowiedzialny za większość brakujących wartości, dlatego zostaną one wypełnione zerami.

Z ratingiem będzie podobnie, jednak tutaj tutaj logiczniejsze wydaje się uzupełnienie wartości losową liczbą z zakresu średnia \pm odchylenie standardowe. (ma spore szanse być poprawną wartością jeśli obiekt rzeczywiście miałby wystawione opinie)

2.5 analiza danych

2.5.1 cena

min	max	średnia	odchylenie std.
14.00	11147.00	181.15	400.67

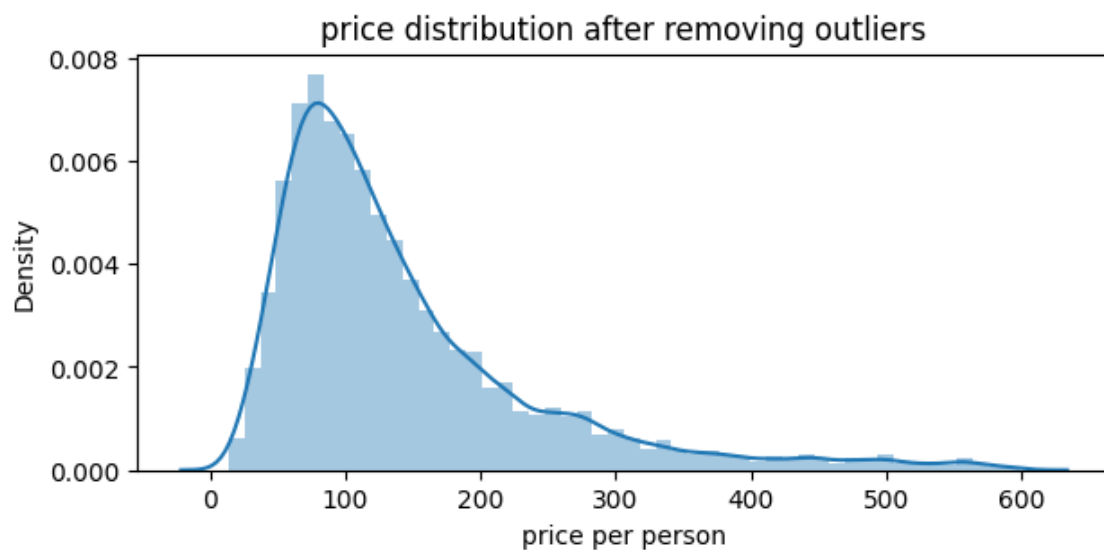


Rysunek 1: dystrybucja cen w oryginalnym zbiorze danych

Jak widać na wykresie, niewielka liczba skrajnych wartości jest odpowiedzialna za spore odchylenie, które nie oddaje rzeczywistego rozkładu większości danych. Z tego powodu zdecydowałam się usunąć część wyników. Za punkt odciecia uznałam cenę 600 € per osoba.

W wyniku tej operacji tracimy 2,34 % oryginalnych danych, jednak odchylenie standardowe spada czterokrotnie (a zakres standardowych cen przestaje zawierać ujemne wartości).

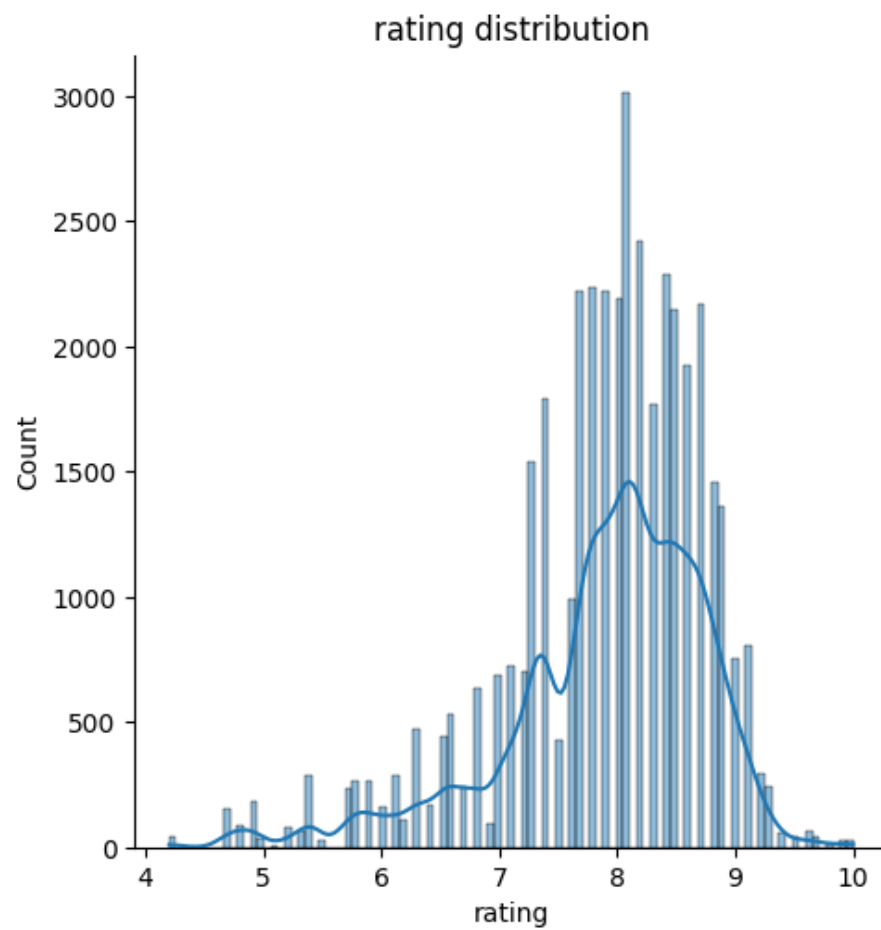
min	max	średnia	odchylenie std.
14.00	598.60	142.90	99.03



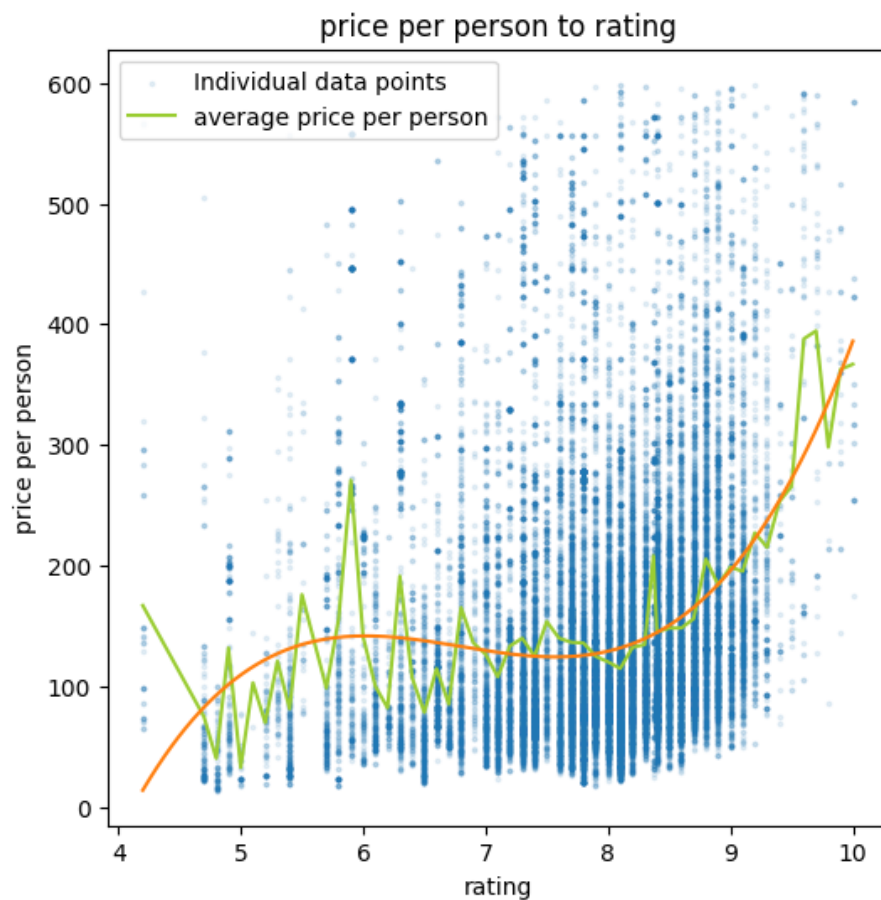
Rysunek 2: dystrybucja cen po usunięciu skrajnych wartości

2.5.2 ocena użytkowników

min	max	średnia	odchylenie std.
4.2	10	7.89	0.86



Rysunek 3: dystrybucja ocen

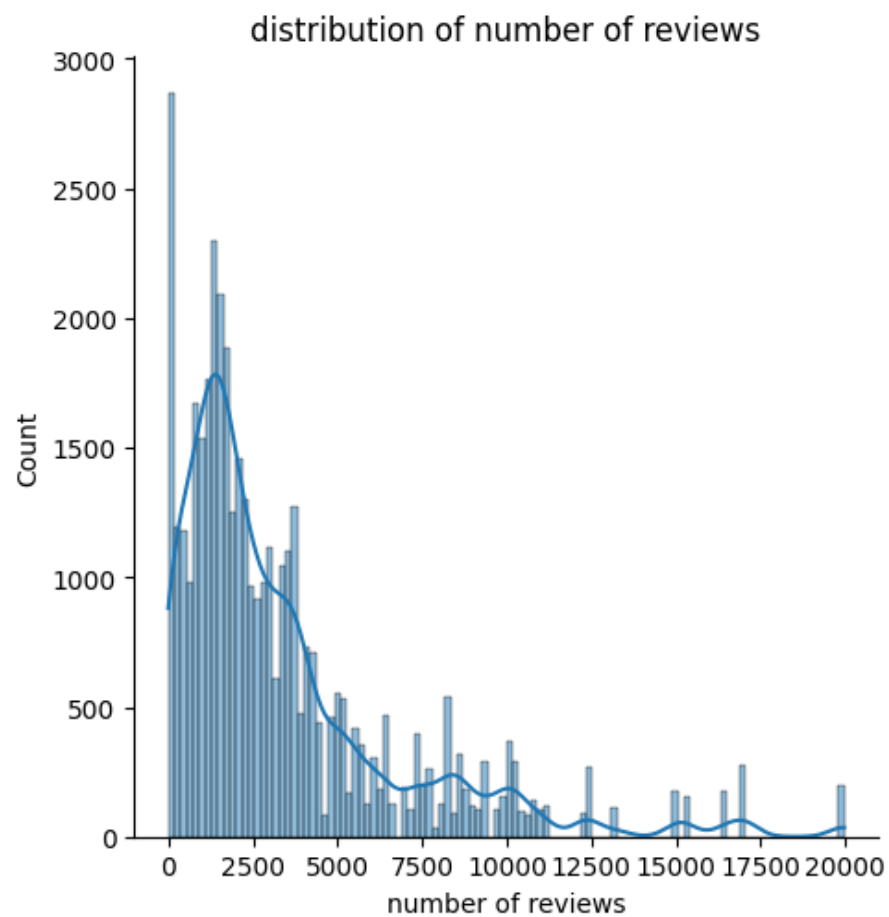


Rysunek 4: zależność ceny rezerwacji od oceny hotelu

Na wykresie można zobaczyć że istnieje pozytywna korelacja między oceną hotelu a ceną jaką musimy zapłacić za pobyt, jednak zależność nie jest liniowa i da się ją zauważyć głównie przy ocenach powyżej 8.

2.5.3 liczba ocen

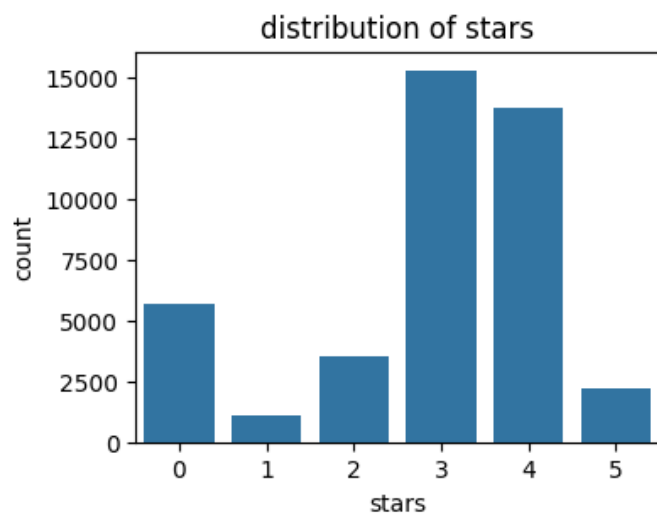
min	max	średnia	odchylenie std.
0.00	19967.00	3424.81	3474.33



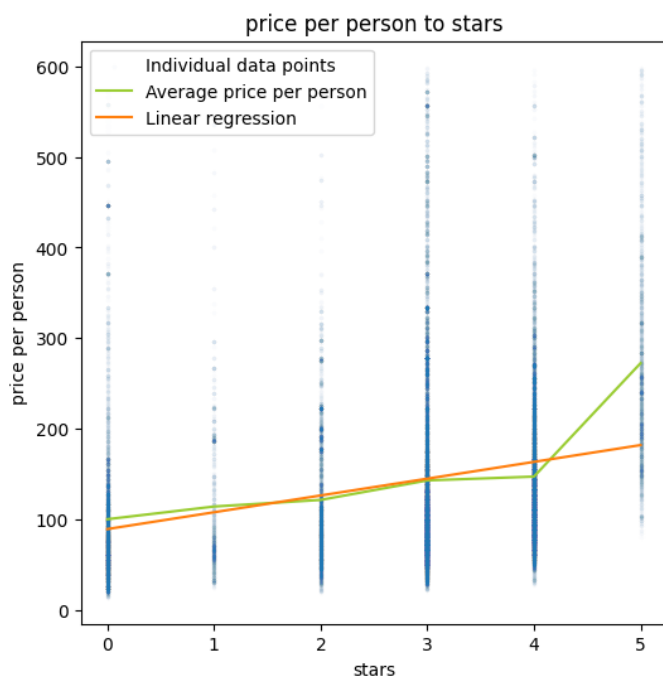
Rysunek 5: dystrybucja liczby ocen

2.5.4 liczba gwiazdek

min	max	średnia	odchylenie std.
0	5	2.89	1.39

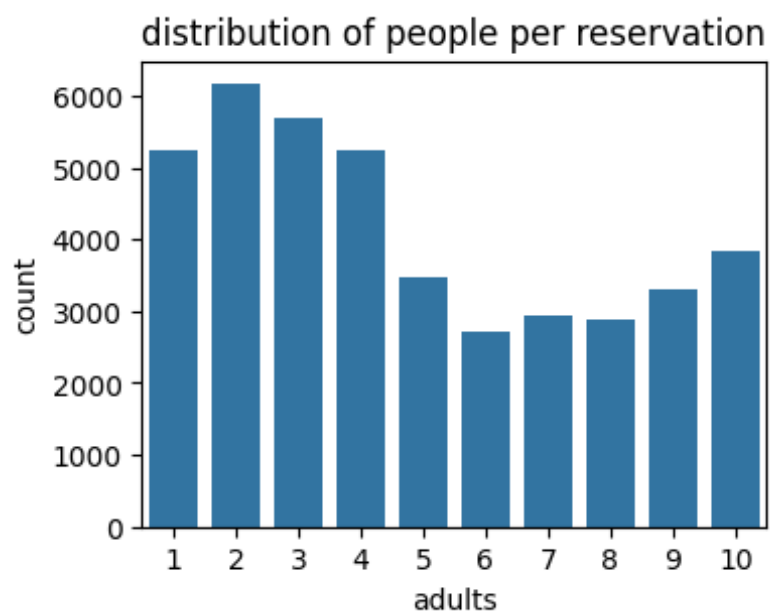


Rysunek 6: dystrybucja liczby gwiazdek

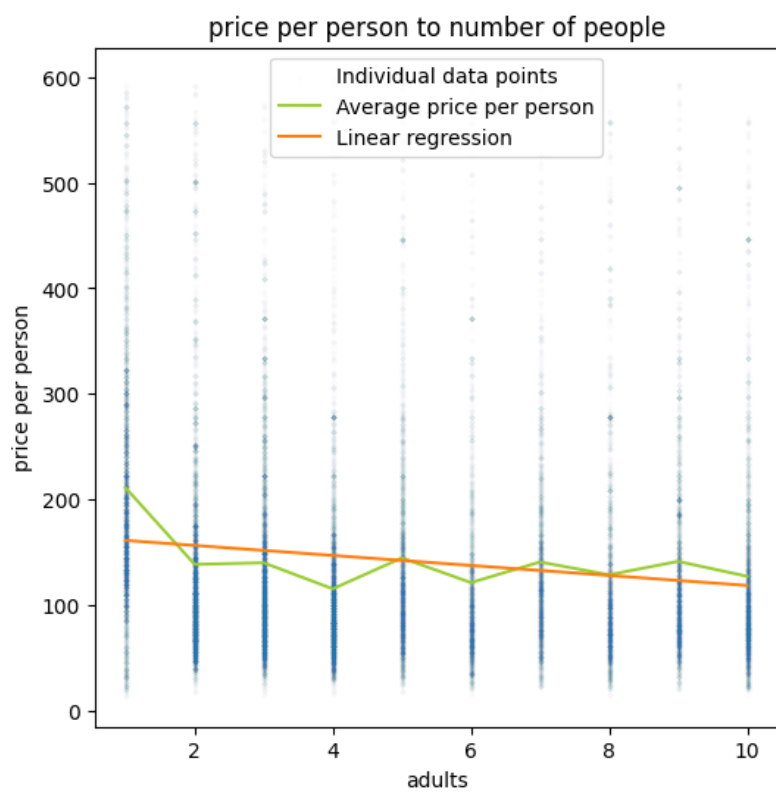


Rysunek 7: zależność ceny od liczby gwiazdek

2.5.5 liczba osób



Rysunek 8: dystrybucja liczby osób



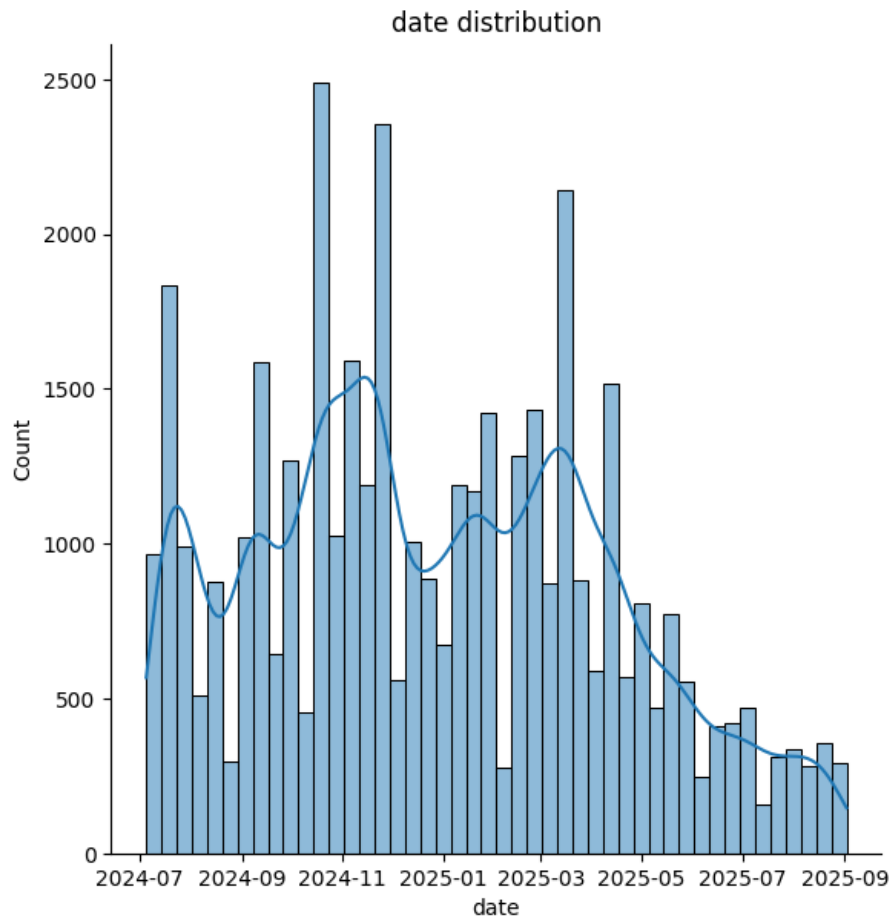
2.5.6 data

Daty mieszczą się w zakresie od 2024-07-05 do 2025-09-02

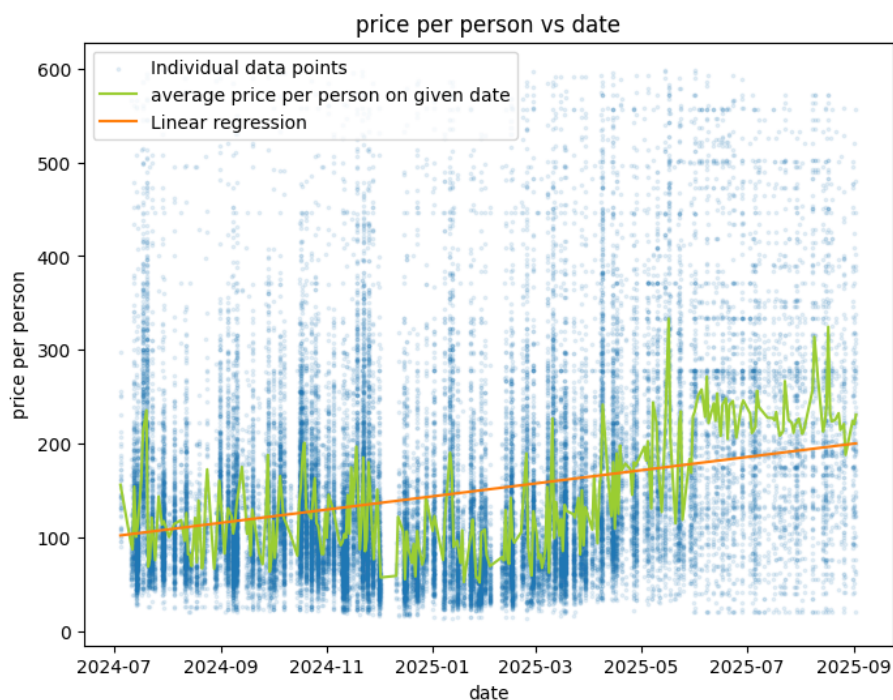
Rozkład wydaje się dość losowy, prawdopodobnie jest on po części dyktowany zainteresowaniem (jeśli data jest oblegana pokaże się mniej ofert). Potencjalnie można w ten sposób wytłumaczyć spadek liczby ofert w okolicy świąt i sylwestra. Tutaj wpływ mogło mieć także odfiltrowanie rekordów ze skrajnymi cenami.

Liczba ofert zaczyna od pewnego momentu maleć wraz z czasem. Może to być powiązane zarówno z nakładającym się okresem wakacyjnym, jak i odległym terminem, przez co hotele mogą jeszcze nie przyjmować rezerwacji.

W celu poprawnego szacowania ceny prawdopodobnie trzeba by uwzględnić zarówno datę pobytu, jak i datę składania rezerwacji oraz analizować dane zebrane w większym przedziale czasu.



Rysunek 9: dystrybucja dat



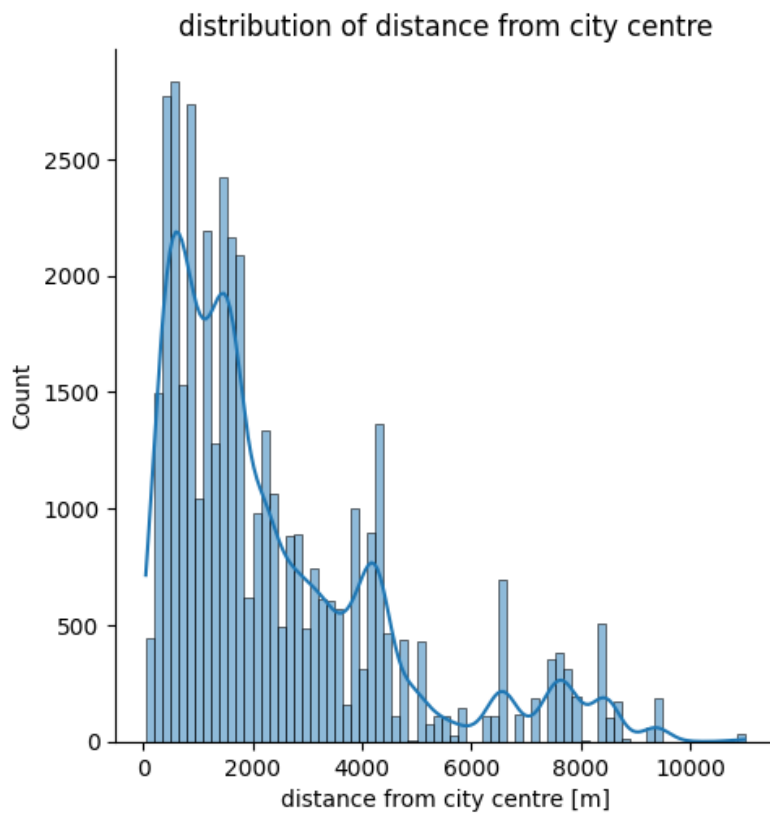
Rysunek 10: zależność ceny od daty rezerwacji

Na wykresie zależności ceny od daty możemy zobaczyć spory skok cen w okresie wakacyjnym. Potencjalnie na wyższe ceny może także wpływać mniejsza liczba dostępnych ofert. Aby jednoznacznie określić powód należałoby uwzględnić dane zbierane na przestrzeni czasu i dodać kolumny odpowiadające za datę zebrania danych oraz liczbę dostępnych ofert dla określonych parametrów wyszukiwania.

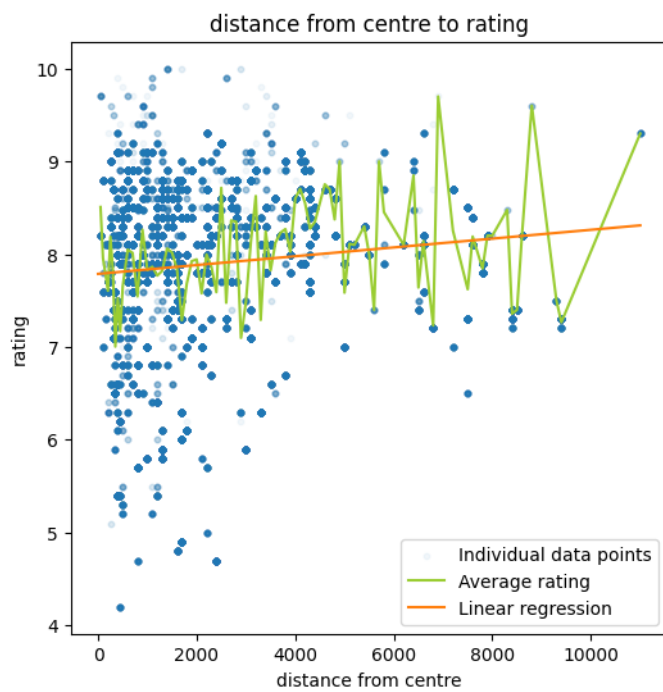
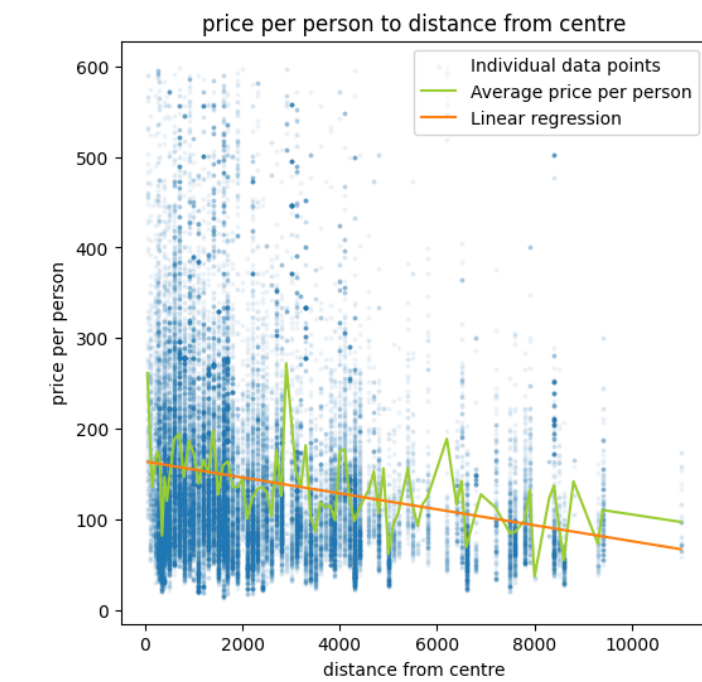
2.5.7 odległość od centrum

min	max	średnia	odchylenie std.
50.00	11000.00	2378.90	2079.47

Zdecydowanie najwięcej ofert znajduje się w okolicy kilometra od centrum. Liczba ofert maleje wraz z dystansem. Występuje negatywna korelacja między ceną a dystansem. Wbrew moim podejrzaniom, odległość od centrum nie wydaje się mieć zbyt dużego wpływu na ocenę hotelu (zachowujemy obie kolumny).

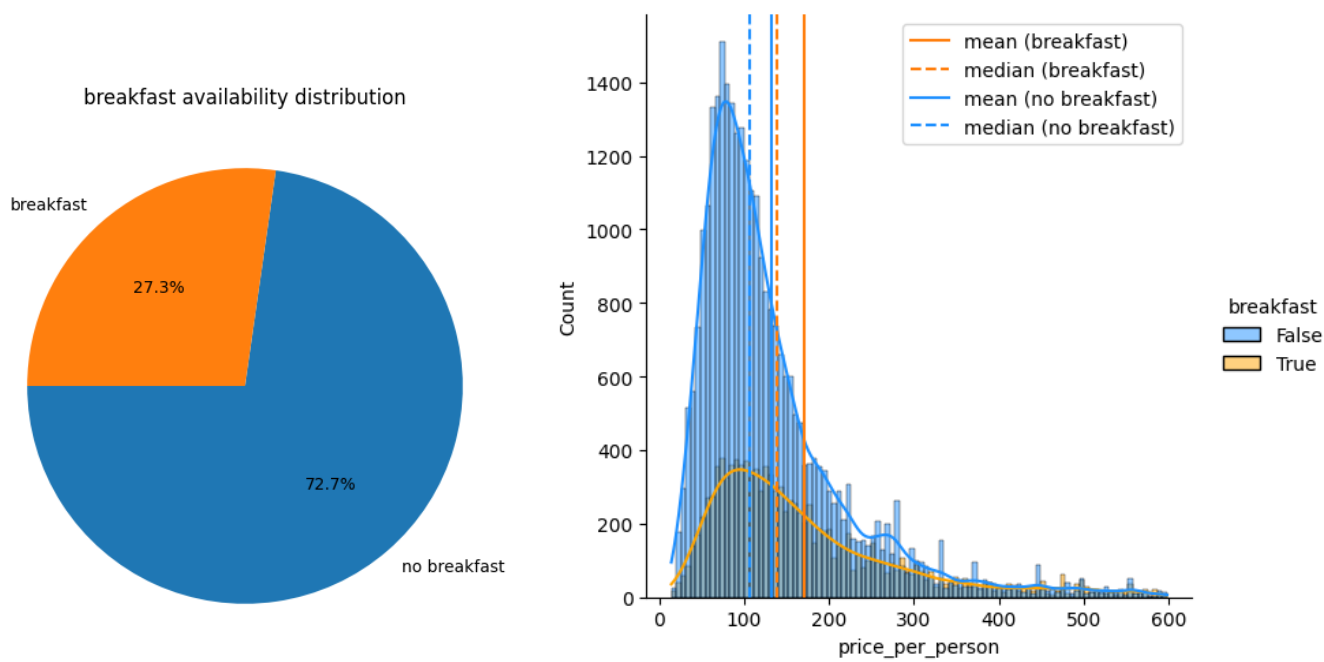


Rysunek 11: dystrybucja liczby osób



2.5.8 śniadanie

Wraz z wzrostem ceny rośnie procent hoteli oferujących śniadanie, co przekłada się na wyższą średnią cenę w tej grupie.



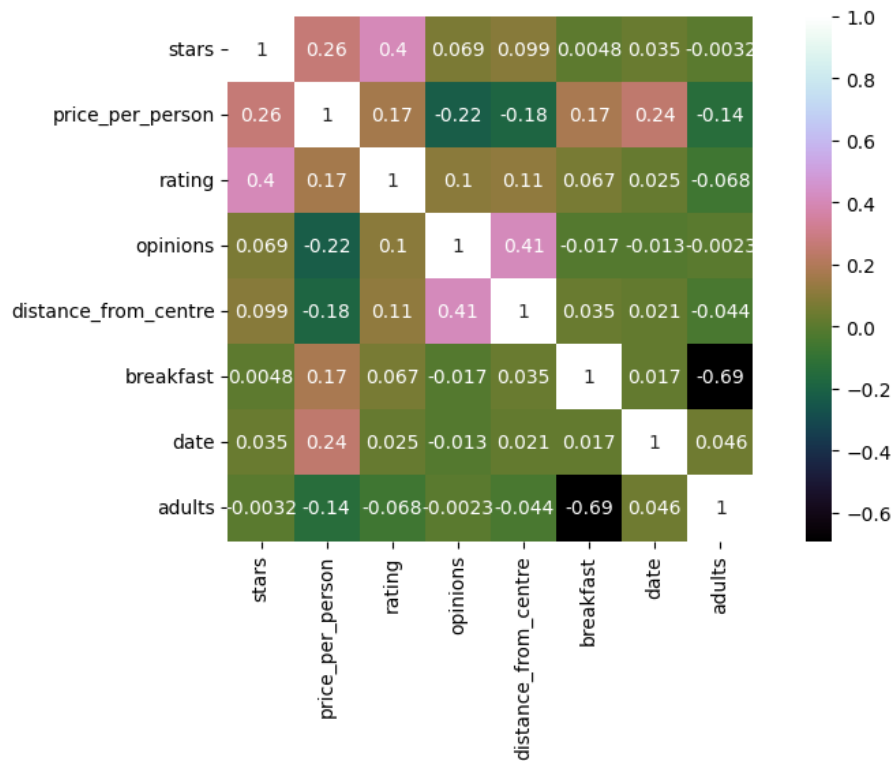
Rysunek 12: dystrybucja cen w zależności od dostępności śniadania

2.5.9 macierz korelacji

Na macierzy korelacji możemy zobaczyć, że wszystkie zmienne mają wpływ na cenę. Większość kolumn ma bardzo małą korelację ze sobą nawzajem. Wyjątkiem jest rating, na który reszta parametrów ma lekko większy wpływ (co było dość spodziewane).

Widać też zależność między liczbą opinii a odległością od centrum (prawdopodobnie większy popyt na hotele w środku miasta). Nie jest ona jednak na tyle duża aby uznać którąś z kolumn za niepotrzebną.

Lekka anomalia jest silna ujemna korelacja między dostępnością śniadania a liczbą osób. Trudno stwierdzić czym jest spowodowana.



Rysunek 13: macierz korelacji

3 Testowane modele

Dane zostały przeskalowane za pomocą StandardScaler (data wcześniej przekonwertowana na czas porządkowy), a następnie losowo podzielone na zbiór treningowy i testowy w proporcji 70:30. Przetestowane zostały:

- 4 algorytmy regresji liniowej: least squares (lls), non-negative ls, ridge i lasso.
- least squares na wielomianach większego stopnia
- random forest z różną liczbą drzew decyzyjnych
- k-neighbours z różną liczbą uwzględnianych sąsiadów

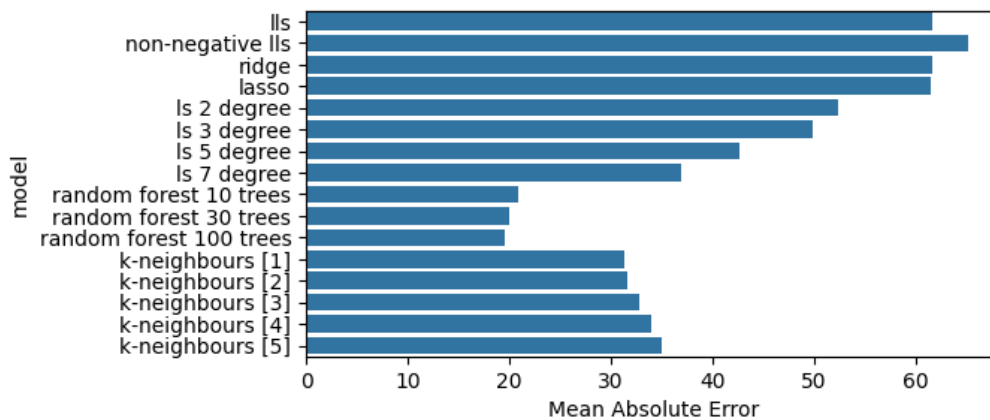
4 Wyniki i wnioski

Zdecydowanie najlepiej wypada algorytm random forest, nawet przy małej liczbie drzew decyzyjnych. Zwiększenie liczby drzew powoduje niewielką poprawę kosztem szybko rosnącego czasu obliczeń.

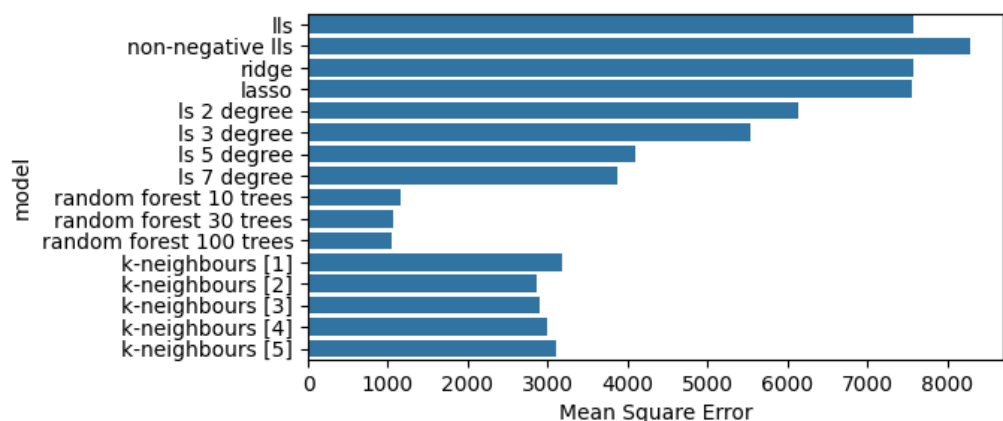
Drugi w kolejności jest algorytm k najbliższych sąsiadów z różnicą we współczynniku determinacji r^2 równą około 0,2. Algorytm działa najlepiej przy uwzględnieniu 3 sąsiadów jeśli za metrykę uznamy współczynnik determinacji i MSE. Najmniejszy średni błąd bezwzględny (MAE) osiąga jednak dla jednego sąsiada.

Algorytmy optymalizujące wielomian pierwszego stopnia wypadają najgorzej i wszystkie osiągnęły niemal identyczne wyniki (różnica przy zmianie parametrów również była prawie niezauważalna). Jedyne, który wypadł lekko gorzej od reszty to Non-negative LLS. Generalne niepowidzenie modeli liniowych ma sens, ponieważ nawet na wykresach możemy zobaczyć, że nie wszystkie zmienne mają liniowy (lub zbliżony do liniowego) wpływ na cenę.

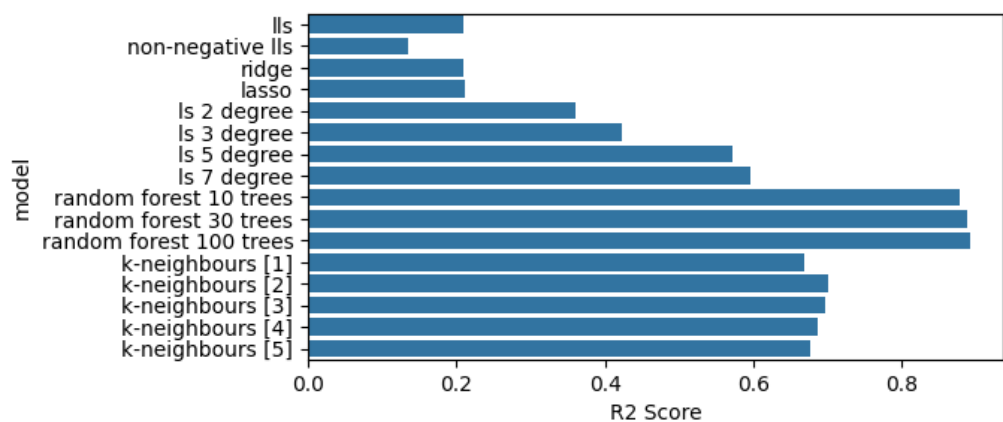
Zwiększanie stopnia optymalizowanego wielomianu poprawia wyniki (zależność między stopniem wielomianu a współczynnikiem dopasowania bliska liniowej). Niestety, nawet drugiemu w kolejności algorytmowi k najbliższych sąsiadów, ls zaczyna dorównywać dopiero przy wielomianie 7 stopnia, podczas gdy czas wywołania robi się ponad 200-krotnie większy



Rysunek 14: średni błąd bezwzględny



Rysunek 15: błąd średniokwadratowy



Rysunek 16: współczynnik determinacji

4.1 propozycje rozwoju

Ponieważ booking pokazuje tylko obecnie dostępne listingi, a ceny mogą się dynamicznie zmieniać w zależności od czasu i dostępności miejsc, sadzę że zbieranie danych na przestrzeni czasu oraz uwzględnienie zarówno daty rezerwacji jak i daty wyszukania mogłoby dawać bardziej użyteczne wyniki. (z obecnym zbiorem danych wiarygodnie da się przewidzieć jedynie propozycje cen, które otrzymalibyśmy składając rezerwacje około 30 maja 2024)

Aby rozszerzyć zakres eksperymentu w generowanych zapytaniach można też uwzględnić różne miasta. (Przeprowadzić ten sam test na innym mieście i porównać wyniki poszczególnych algorytmów lub uwzględnić miasto jako zmienną regresji)