

1 Zadanie

W naszej analizie, przyjrzymy się zbiorowi zakwaterowań hotelowych, a naszym zadaniem będzie przewidzieć, czy pojedyncza wizyta zostanie odwołana. Naszymi modelami uczenia maszynowego będą LightGBM oraz regresja logistyczna. Dokładność modelu LightGBM na zbiorze treningowym oraz testowym wynoszą odpowiednio 89% oraz 88%.

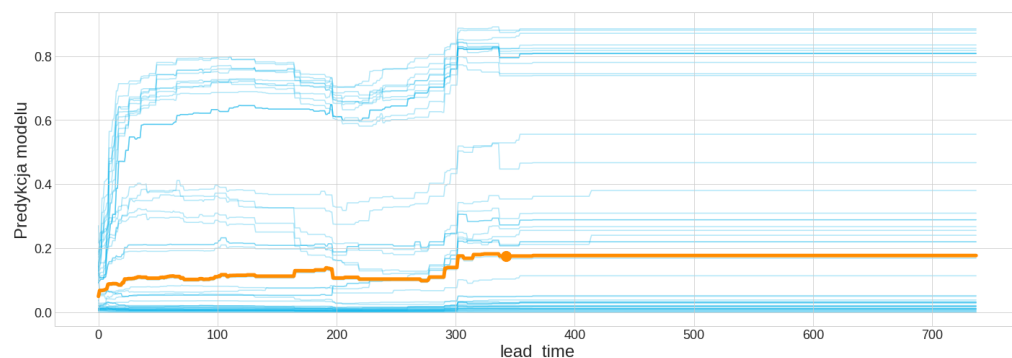
2 Zadanie

Zdaniem naszego modelu z prawdopodobieństwem 18% nasz klient odwoła rezerwację i jest to poprawna decyzja, gdyż nasz gość tej rezerwacji nie odwołał.

3 Zadanie

Spójrzmy teraz na wykresy metody Ceteris Paribus dla kilku najciekawszych zmiennych objaśniających. Nasza obserwacja będzie podkreślona kolorem pomarańczowym. Dodatkowo prawdziwa wartość zmiennej będzie oznaczona punktem, jeżeli takiego punktu nie ma to oznacza to, że nasza obserwacja ma brakującą wartość dla szukanej zmiennej. Pozostałe 49 losowych obserwacji dadzą nam więcej informacji o działaniu modelu z konkretną zmienną i będą oznaczone kolorem błękitnym.

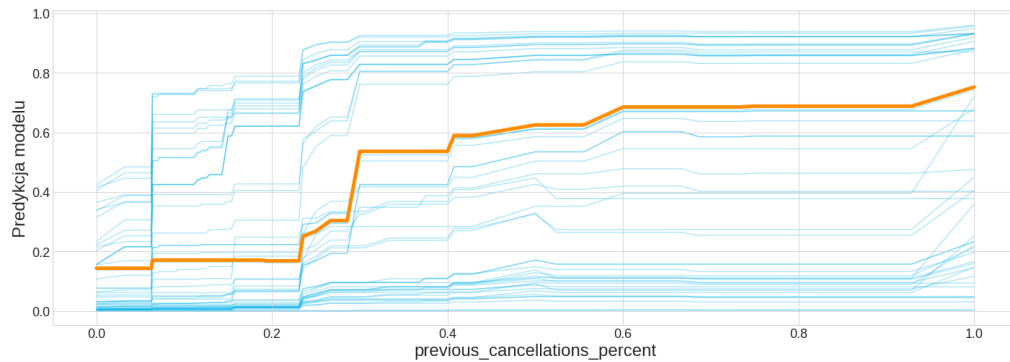
Pierwszym wykresem jest CP dla zmiennej oznaczającej różnicę między datą rezerwacji, a przyjazdu. Idąc od prawej strony widzimy, że osoby zamawiające zakwaterowanie z rocznym wyprzedzeniem częściej odwołują zamówienia. Ponadto, największy spadek odwołań jest umiejscowiony blisko miesiąc przed przybyciem. Ten efekt może być związany z całkowitym zwrotem kosztów jeżeli odwołamy rezerwację odpowiednio wcześniej.



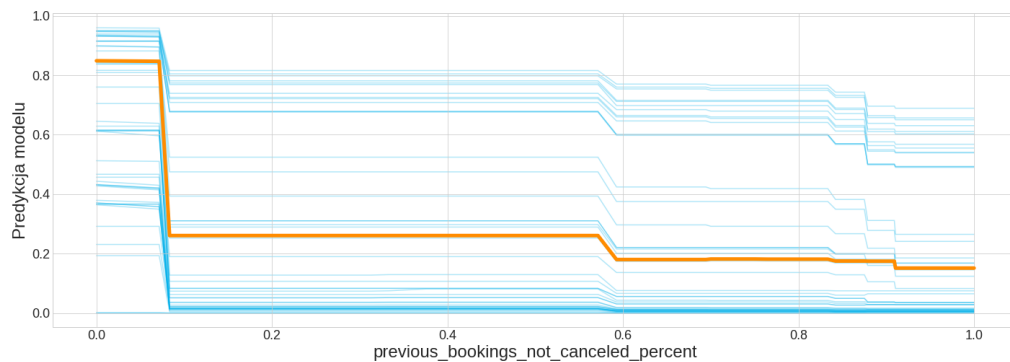
Ceteris paribus dla odległości między dniem zarezerwowania, a dniem przyjazdu.

Kolejne dwie bliźniacze zmienne to procent poprzednich wizyt, które zostały

odwołane, bądź nie. W obydwu przypadkach widzimy mocny trend między procentem odwołanych rezerwacji, a szansa odwołania wizyty.



Ceteris paribus dla procentu poprzednich, odwołanych rezerwacji.



Ceteris paribus dla procentu poprzednich, nieodwołanych rezerwacji.

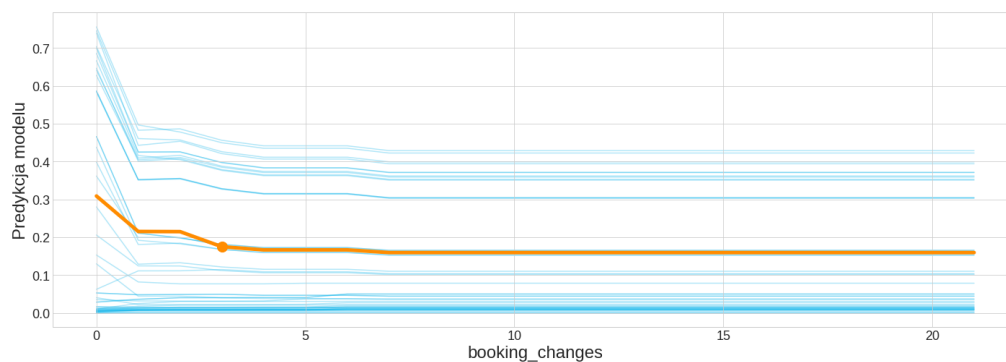
Ostatnie dwie zmienne, które rozważymy to liczba zmian w zamówieniu od początku jego zatwierdzenia oraz liczba miejsc parkingowych, wymaganych przez klienta. W obydwu przypadkach szansa na anulowanie pobytu drastycznie spada dla wartości większych niż zero.

Dla zmiennej *booking_changes* może mieć to związek z tym, że w przypadku niezadowolenia z rezerwacji mamy dwie opcje: anulowanie, albo zmiana. W przypadku, gdy prosimy o zmiany, może nam bardziej zależeć na rezerwacji.

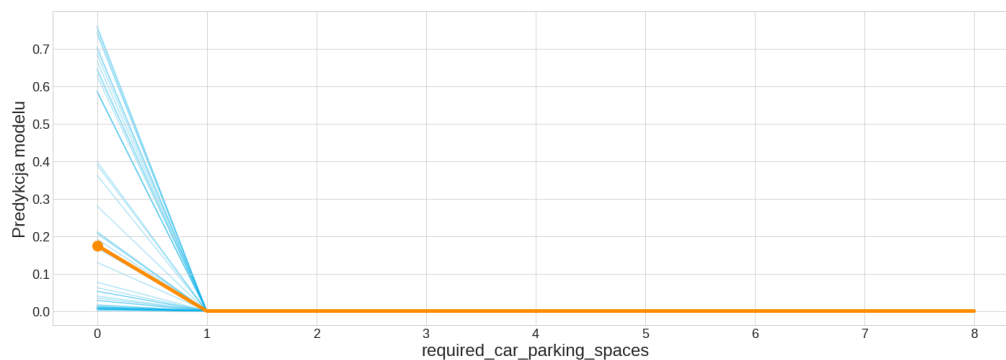
W drugim przypadku liczba miejsc parkingowych może mieć związek z charakterem naszego pobytu. Jeżeli jest to pobyt w większej grupie osób, to liczba miejsc parkingowych powinna wzrosnąć (na przykład dla autokaru). Ponadto, rezerwacje dla większej liczby osób może być bardziej przemyślana i podjęta przy konsultacji większej liczby osób.

W obydwu powyższych przypadkach, dodanie zmiennych oznaczających czy

wyżej wymienione predykaty są równe zero, czy nie może być przydatne przy kolejnej iteracji budowania modelu.



Ceteris paribus dla liczby zmian w zakwaterowaniu.



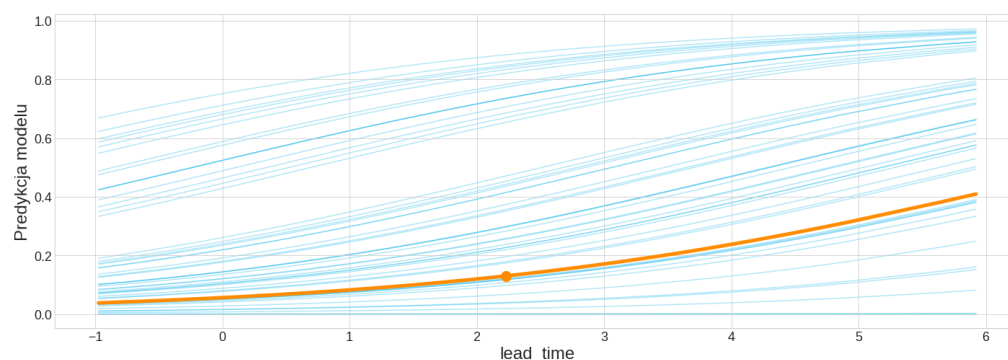
Ceteris paribus dla wymaganej liczby miejsc parkingowych.

4 Zadanie

Zmieniając zmienną *is_repeated_guest*, zazwyczaj otrzymamy większe prawdopodobieństwo przy zmianie z wartości zero na jeden, gdzie predykcja modelu może się różnić nawet o 0.4. Jest jednak przypadek, gdzie zmiana tej wartości z jedynki na zero zwiększa prawdopodobieństwo anulowania o 0.2. Powodem tego może być, niebalansowana liczba powtórnych gości wynosząca zaledwie 3%. Dodając fakt, iż nasz model miał mniejszy wynik na zbiorze testowym, niż na treningowym, może oznaczać początek przeuczenia się modelu.

5 Zadanie

To co możemy zauważyć w modelu logistycznym to gładkie wykresy CP, co oczywiście jest związane z tym, że decyzja naszego modelu jest funkcja ciągła, różniczkowalna. Największe różnice można zauważyć dla zmiennej *lead_time*. W przypadku użycia metod drzewiastych widoczny jest punkt przegięcia funkcji, gdzie z wykresu wklęsłego zmienia się on na wypukły, podczas gdy w nowym modelu tego nie ma. Nie mniej jednak, sens wyjaśnienia tej zmiennej pozostaje prawdziwy (Osoby, które wcześniej zamawiają, częściej anulują rezerwacje). Jednak w przypadku prostego modelu liniowego, lokalne zmiany nie zostały wyłapane. Dodanie zmiennych wielomianowych od tej zmiennej może poprawić powyższy model logistyczny.



Ceteris paribus modelu regresji logistycznej dla odległości między dniem zarezerwowania, a dniem przyjazdu.