

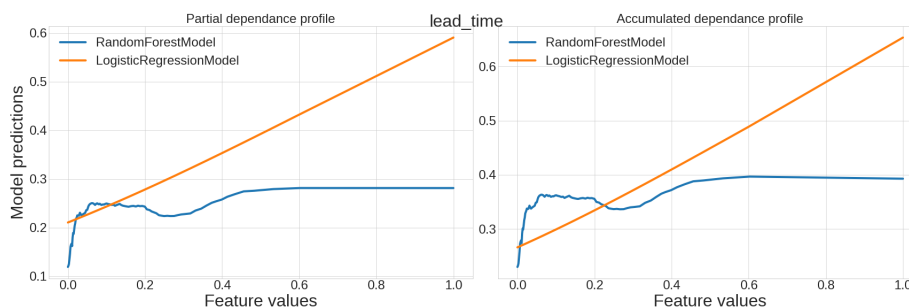
XAI - Praca domowa #6

Miłosz Kacper Michta

May 14, 2020

W naszej analizie, przyjrzymy się zbiorowi zakwaterowań hotelowych, a naszym zadaniem będzie przewidzieć, czy pojedyncza wizyta zostanie odwołana. Naszymi modelami uczenia maszynowego będą las losowy oraz regresja logistyczna.

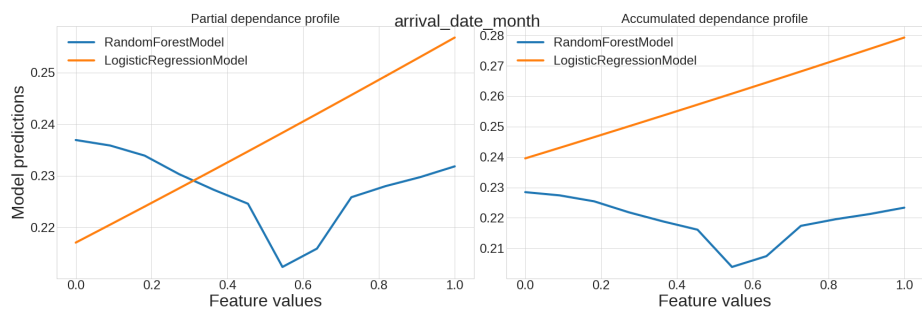
Spójrzmy teraz na profile czastkowe i kumulatywne, naszych modeli pod kątem kilku cech. Zaczniemy od zmiennej *lead_time* oznaczającej różnicę w dniach między data zakwaterowania, a rezerwacji.



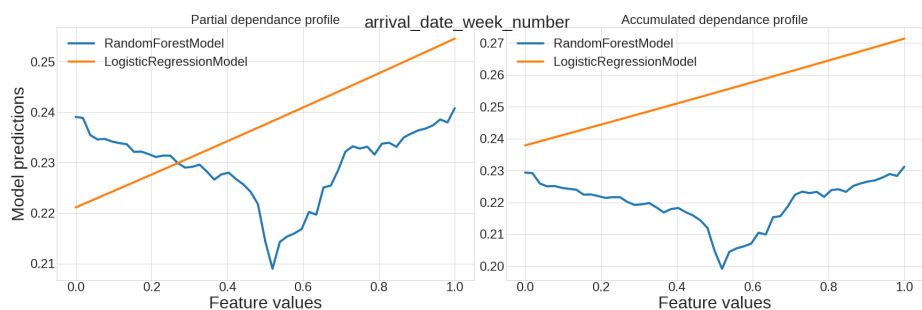
Profile dla odległości między dniem zarezerwowania, a dniem przyjazdu.

To co możemy zaobserwować patrząc na ten wykres, to przewaga modelu drzewiastego, który jest w stanie lepiej się dopasować do danych z nieliniowym trendem. Kolejną różnicą jest to, że dla opcji kumulatywnej, predykcje obu modeli przyjmują większe wartości, co mówi nam o tym, że modele zawierają interakcje, które osłabiają wpływ powyższej cechy samej w sobie.

Kolejne dwie bliźniacze zmienne to czasy przybycia do hotelu, odpowiednio w miesiącach i tygodniach.



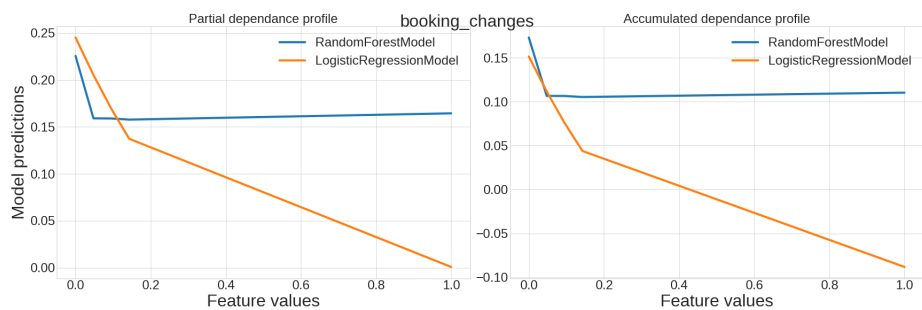
Profile dla miesiąca zakwaterowania.



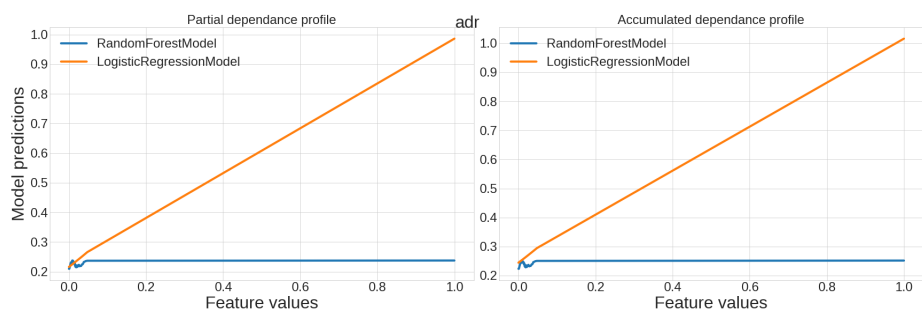
Profile dla tygodnia zakwaterowania.

Podobnie jak w poprzednim przykładzie, model drzewiasty ma przewagę nad modelem logistycznym w postaci wykonania podziału na danych niemonotonicznych/nieliniowych. Uwzględniając wakacyjną porę roku, las losowy uznał, że w tym okresie prawdopodobieństwo anulowania rezerwacji jest mniejsze z uwagi na zajętość miejsc w większości wakacyjnych kurortów.

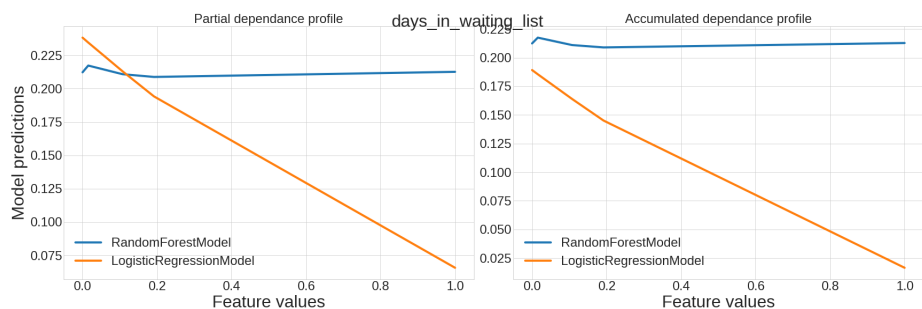
W powyższych przykładach, mogliśmy zaobserwować dużą skuteczność modelu drzewiastego, jednak poniżej zaprezentujemy przykłady zmiennych, które przechylają szalę w stronę modelu regresji logistycznej.



Profile dla liczby zmian w zamówieniu.



Profile dla średniej ceny za pokój w wybranym hotelu.



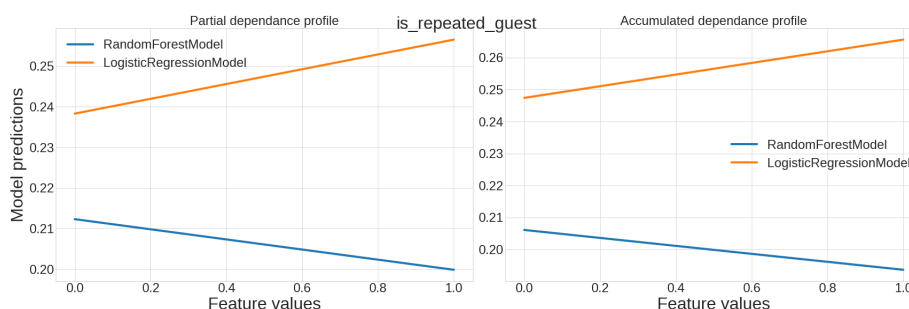
Profile dla długości oczekiwania na akceptację, liczona w dniach.

Pierwsza rzecz rzucająca się w twarz, to brak reakcji lasu losowego na zmiany wartości w powyższych zmiennych, w przeciwieństwie do regresji logistycznej.

Wszystkie trendy w pokazanych, powyżej grafach mają sensowną interpretację. Jeżeli składamy zmiany w rezerwacji, to zapewne zależy nam na danym zakwaterowaniu. Jesteśmy bardziej skłonni zrezygnować z drogiego wywczasu,

niż z atrakcyjnej oferty. Im dłużej czekamy na zwolnienie miejsca w hotelu, tym większe jest nasze zaangażowanie.

Jednak dlaczego model drzewiasty nie jest zdolny nauczenia się tych kilku bardzo prostych reguł? Wyjaśnień może być wiele, jedno z nich to fakt, że modelom drzewiastym zdecydowanie trudniej jest aproksymować funkcje liniowe, niż modelom liniowym. Kolejnym bardzo znaczącym powodem (moim zdaniem) mogą być pierwsze trzy przykłady, które obrazują jak las losowy wybrał cechy związane z datą jako bardziej istotne, niż zmienne ciągłe opisujące rezerwacje. Spójrzmy jeszcze na różnice przy zmiennej opisującej czy ktoś jest powracającym gościem.



Profile dla cechy *is_repeated_guest*.

To co możemy zaobserwować to, wzajemnie przeczące sobie profile dla zmiennej *is_repeated_guest*. Model drzewiasty poprawnie znajduje trend w odróżnieniu od regresji logistycznej. Ponadto, negatywny trend modelu liniowego ma zdecydowanie mniejszy wpływ na model patrząc na profil cząstkowy, a kumulatywny jeszcze bardziej pogłębia te różnice. Powodem tego może być niezbalansowanie przykładów dla danej zmiennej (gości powracających jest około 3%), przez co maksymalizacja entropii przy podziale może mieć większe znaczenie, niż logistyczna funkcja straty.

To co można wykorzystać z powyższej analizy, to wzmocnienie modelu logistycznego poprzez dodanie zmiennych wielomianowych, bądź innego kodowania zmiennych opisujących datę (choćby, przez kodowanie cykliczne na okręgu jednostkowym).