

Picipolo - dokumentacja

April 9, 2022

1 Czyszczenie danych i redukcja wymiarowości

W celu wyczyszczenia danych:

- usuneliśmy zmienne identyfikujące (numer budynku, numer ulicy itp.)
- usuneliśmy zmienne skorelowane, zostawiliśmy najbardziej predykcyjne, użyliśmy korelacji spearmana i pearsona
- usuneliśmy zmienne, które mogliśmy uzyskać z innych zmiennych
- usuneliśmy zmienne, w których występuje jedna wartość oraz te które w większości były brakami
- w niektórych kolumnach NaN oznaczał 0, więc wystarczyło zrobić One Hot Encoding, aby poradzić sobie z brakami.

2 Wybór i wytrenowanie modelu uczenia maszynowego

Wykonaliśmy preprocessing:

- usuneliśmy zmienne, które były mało znaczące w naszym modelu
- imputowaliśmy rekordy, w których powierzchnia działki i powierzchnia budynku były równe 0 (KNNImputer), nie mogliśmy ich usunąć, bo stanowiły 20 % ramki
- wykonaliśmy OneHot Encoding zmiennych
- zmieniliśmy kolumnę YEAR z przedziałów na rok
- zmieniliśmy wymiary budynków i działek na ich powierzchnie
- skalowanie zmiennych za pomocą MinMax na przedział $[0,1]$
- usuneliśmy outliery
- wykonaliśmy transformację logarytmiczną, kolumn o wysokiej skośności

- w celu oszacowania ceny wytrenowaliśmy RandomForestRegressor (ze wszystkich modeli, które testowaliśmy, ten radził sobie najlepiej na domyślnych ustawieniach), próbowaliśmy także dobrać hiperparametry, jednak te domyślne okazały się być najlepsze

3 Zrozumienie i interpretacja wyników.

Otrzymaliśmy średni błąd na poziomie 100 000, jednak model mylił się także nawet o 22 000 000. Ramka danych nie była najlepsza do predykcji, ponieważ jedna zmienna miała udział w modelu równy 0.95. Dodatkowo ta zmienna też oznaczała cenę i była skorelowana z etykietą na poziomie 1. Pozostałe zmienne nie miały dużego udziału w predykcji, przez co model mylił się znacząco w niektórych przypadkach. Dodatkowo wiele ważnych zmiennych takich jak wymiary budynku, czy działki w wielu przypadkach były równe 0, a z pewnością ta zmienna miałaby duży wpływ na cenę.

4 Prezentacja i zastosowanie

Uważamy, że przygotowany przez nas model mógłby zostać wykorzystany w obydwu podejściach biznesowych - do oszacowania cen nieruchomości w Technopolis takiej jak Nowy Jork zarówno dla potencjalnych nabywców, jak i sprzedających. W naszej aplikacji użytkownik mógłby podać parametry nieruchomości, a następnie otrzymać jej wycenę. Mógłby on także zaznaczyć miejsce na mapie, a nasza aplikacja ściągnęłaby informacje na temat pobliskich nieruchomości na sprzedaż ze strony aukcyjnej, a także sprawdziłaby na podstawie pobranych parametrów, czy nieruchomości te są dobrze wycenione. Pozwoliłoby to kupującemu na nieprzeplacanie i uczciwym sprzedającemu na wyznaczanie odpowiedniej ceny.