

Przewidywanie ceny samochodu

Denys Troshchylo

12 czerwca 2023

Spis treści

1	Definicja problemu badawczego	3
2	Zdobycie danych	3
3	Przygotowanie, wyczyszczenie danych	3
4	Wstępna analiza danych	4
5	Przygotowanie modelu	4
6	Wnioski, komentarze, ulepszenia	5
7	Link do projektu	5

1 Definicja problemu badawczego

W tym projekcie badawczym celem jest stworzenie modelu, który będzie przewidywać ceny samochodów na podstawie dostępnych danych. Celem jest stworzenie modelu, który będzie dokładnie przewidywał ceny na podstawie różnych atrybutów samochodów. Cały projekt jest napisany w języku python z wykorzystaniem różnych bibliotek

2 Zdobycie danych

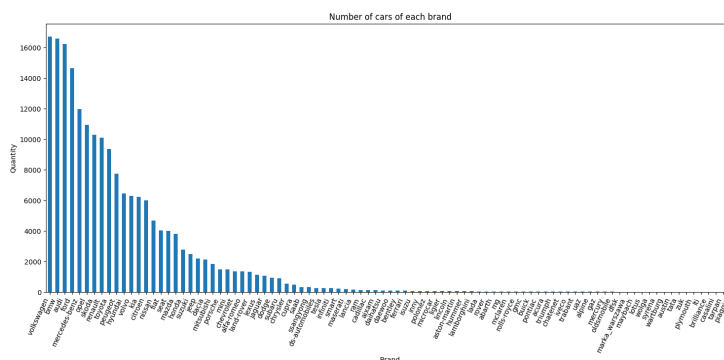
Do realizacji pobierania danych dla tego projektu użyłem biblioteki *scrapy*. Dane pobierałem ze strony <https://www.otomoto.pl/>. Podczas implementacji danej części napotkałem kilka problemów. W pierwszej wersji tego skryptu otwierałem każdą stronę z ofertami samochodów, a następnie otwierałem każdą stronę każdego samochodu po to żeby pobierać więcej atrybutów dla każdego samochodu. W tym przypadku miałem problem z tym, że skrypt działał tak szybko że strona z której pobierałem dane blokowała mnie po 5 minutach skrapowania. Żeby wyeliminować ten problem zrobiłem losowe opóźnienie przed każdym pobieraniem. Potem wystąpił taki problem, że po pobraniu mniej więcej 15 tysięcy stron predkość pobierania nagle spadała i nic się nie zapisywało do bazy samochodów, później się okazało, że tak naprawdę ta strona nie miała 6 tysięcy stron z ofertami, a tylko 500 i potem jak skrypt próbował wejść np. na 501 stronę to zostawał przekierowany na 500. Wtedy zmieniłem strategię i pobrałem dla każdej marki samochodu dostępne modele i przechodziłem już po stronach modeli i nie wchodzić na stronę z samochodem skrapowałem dane tego samochodu.

3 Przygotowanie, wyczyszczenie danych

Przed przystąpieniem do analizy danych, przeprowadziłem kilka kroków w celu przygotowania i wyczyszczenia danych. Najpierw to było sprawdzić czy mam wszystkie atrybuty dla każdego samochodu jak się okazało to występowały takie oferty w których nie było przebiegu, to oznacza że samochód jest nowy, więc po prostu wstawiałem 0 w pole przebiegu. Również były takie oferty które nie miały pojemności skokowej, ale to występowało wyłącznie u samochodów które miały rodzaj paliwa Elektryczny lub Hybryda, w tym przypadku, jak i w poprzednim, zamieniałem na 0. W sumie udało mi się pobrać 195000 opisów aut. Kolejny problem był taki, że niektóre modele miały dużo stron i jak się przechodziło przez te strony czasami występowały powtórzenia ofert, żeby poradzić sobie z tym problemem usunąłem duplikaty. Po wyeliminowaniu tych problemów zostało +100 samochodów które zawierały jakieś błędne dane, dlatego usunąłem je. Niektóre ceny były w EUR, więc po prostu rzeliczyłem walute z euro na złotówki. Na wyjściu z tego dostałem 192000 poprawnych jednostek z danymi.

4 Wstępna analiza danych

Przeprowadziłem wstępną analizę danych w celu zrozumienia zawartości zbioru i identyfikacji ewentualnych zależności. Wykorzystałem techniki wizualizacji danych za pomocą biblioteki *matplotlib*. Przykładowo, wygenerowałem wykres przedstawiający stosunek ilości pobranych aut do marek. Inne wykresy do analizy danego zbioru umieściłem w Jupyter Notebook'u



Rysunek 1: Stosunek ilości pobranych aut do marek

5 Przygotowanie modelu

Przygotowanie modelu składało się z kilku etapów. Wybrałem sobie model drzewa decyzyjnego z biblioteki *sklearn*. Już gotowy sformatowany zbiór danych umieściłem w DataFrame z biblioteki *pandas*, żeby móc wytrenować model. Następnie wytrenowałem model oraz zapisałem go do pliku żeby nie robić tej samej czynności wiele razy.

6 Wnioski, komentarze, ulepszenia

Po przeprowadzeniu analizy danych i wytrenowaniu modelu, doszedłem do następujących wniosków. Mój model osiągnął dokładność na poziomie 10-30 procent i uważam, że dla takiego, dość małego, zbioru atrybutów to jest dobry rezultat. Najtrudniejszą częścią tego projektu dla mnie było skrapowanie danych bo było wiele niuansów które trzeba uwzględnić. Żeby polepszyć precyzyjność tego programu, moim zdaniem, można by było:

1. Pobierać więcej atrybutów wchodząc na stronę każdego auta.
2. Pobrać dane również z innych stron które sprzedają samochody.
3. (Zaawansowane) Pobierać zdjęcia samochodów oraz wykrywać widoczne na nich uszkodzenia.

7 Link do projektu

Projekt jest dostępny na platformie GitHub pod adresem:
<https://github.com/troshchylo/predicting-car-prices>