

Klasyfikacja płci głosu

Maciej Chylak, Maciej Gryszkiewicz

Kwiecień 2021

Contents

1	Wstęp	2
2	EDA - badania eksploracyjne	2
3	Preprocessing	3
4	Wybór modelu	4
5	Strojenie modelu	4
6	Test modelu na własnych danych	5

1 Wstep

Nasz projekt opiera sie na danych ze strony <https://apispreadsheets.com/datasets/119>. Owy zbiór danych składa sie z 3168 rekordów opisujących różne charakterystyki głosu ludzkiego (np. czestotliwość, ton dominujący itd.). Celem projektu jest utworzenie modelu klasyfikacyjnego, który jak najdokładniej będzie przewidywał czy dany głos należy do kobiety, czy do mężczyzny.

2 EDA - badania eksploracyjne

Na początku badań eksploracyjnych zgłębiliśmy temat reprezentowania dźwięku w postaci liczb, w szczególności głosu ludzkiego. Pozwoliło to nam to lepiej zrozumieć dane, na których operowaliśmy. Oprócz tego, stworzyliśmy skrypt w R, którego zadaniem jest transformowanie nagrań głosu ludzkiego na ramki danych analogiczne do tych dostępnych na stronie apispreadsheets.com. Dzięki temu zabiegowi mogliśmy na koniec przetestować nasz model na samodzielnie nagranych próbkach.

Dalsza część EDA stanowiła wizualizacja danych. Zapoznaliśmy się z rozkładami każdej ze zmiennej przy pomocy histogramów. Następnie zbadaliśmy korelacje zmiennych przy pomocy macierzy korelacji reprezentowanej przez heatmapę. W kolejnym kroku sporządziliśmy pairploty dla czterech najciekawszych zmiennych. Wybraliśmy te zmienne ze względu na to, że różnica w dystrybucjach męskich i żeńskich była w ich wypadku największa. Na koniec tego rozdziału stworzyliśmy 2 ridgeploty ilustrujące zmienne, których nie udało nam się odtworzyć przy pomocy wyżej wspomnianego skryptu w R. Na nasze szczęście, zmienne te nie wносиły zbyt wiele do naszego modelu, za ich pomocą nie dało rozróżnić się głosu męskiego od żeńskiego.

3 Preprocessing

Początkowo przyjrzelśmy się rozkładowi zmiennych. W trakcie przeglądania heatmapy zauważyliśmy, że kilka zmiennych, tj. centroid, kurt, dfrange jest mocno skorelowanych z innymi, więc ostatecznie postanowiliśmy je wyrzucić. Następnie dokonaliśmy transformacji zmiennych skew, modindx oraz maxfun, gdyż wykresy tych zmiennych wykazywały duże skośności.

Zajeliśmy się także analizą ciszy, gdyż zauważyliśmy, że większość obserwacji zawiera bardzo długie fragmenty, w których do mikrofonu nic nie mówiono. Szczególnie widać to w przypadku zmiennej Q25, gdzie bardzo duża część naszych obserwacji posiada tę zmienną o wartości mniejszej niż 80Hz, a jak wiadomo głos ludzki rozpoczyna się od około 80Hz. Postanowiliśmy zatem usunąć wszystkie obserwacje, które mają wartość tej zmiennej poniżej 60Hz (Aby nie wyrzucić przypadkiem zbyt wielu obserwacji).

Sprawdziliśmy także, czy któraś z naszych zmiennych ma dystrybucję rozkładu normalnego. Zrobiliśmy to w celu łatwej możliwości usunięcia outlierów. W tym celu użyliśmy biblioteki scipy, lecz jak się okazało żadna z tych zmiennych nie wykazała tej cechy (było to spowodowane tym, że test ten w sposób bardzo surowy podchodzi do definicji bycia rozkładem normalnym).

Ostatnią część preprocessingu stanowił encoding wartości label. W naszym wypadku 0 - oznaczać będzie głos żeński, natomiast 1 - męski.

4 Wybór modelu

Przetestowaliśmy następujące modele:

- SVM
- KNNNeighbors
- Regresja logistyczna
- Naiwny Klasyfikator Bayesowski
- Drzewo decyzyjne
- Model stackowany (DecisionTree, NB, LogisticRegression — > LogisticRegression)
- XGBoost

Głównym kryterium wyboru modelu był wynik "accuracy". Zdecydowaliśmy się na użycie tej metryki, ponieważ nasz zbiór danych był idealnie zrównoważony. Do tego, z naszego punktu widzenia nie ma różnicy pomiędzy predykcją typu false-positive, a false-negative. Oprócz tego, jest to miara najbardziej intuicyjna. Dodatkowo wspomagaliśmy się metryką ROC AUC. Używaliśmy jej jako ogólny i uniwersalny wyznacznik jakości modelu, ale "accuracy" pozostała decydującym czynnikiem.

Najlepiej wypadł model XGBoost i to jego zdecydowaliśmy się użyć w dalszej części projektu. Bardzo dobrym kandydatem wydawał się również model Decision Tree. Jego "accuracy" stało na bardzo przyzwoitym poziomie, a dodatkowo był w pełni przejrzysty i jawny.

5 Strojenie modelu

Aby polepszyć jakość wybranego modelu, zdecydowaliśmy się na dostrojenie jego hiperparametrów. W tym celu użyliśmy metody RandomizedSearchCV. Skupiliśmy się na najistotniejszych parametrach takich jak:

- learning_rate
- gamma
- max_depth
- nround
- n_estimators
- min_child_weight

Strojenie to nie zmieniło w istotnym stopniu jakości naszego modelu, niemniej jednak delikatnie poprawiło osiągnięte przez model "accuracy".

6 Test modelu na własnych danych

Na samym końcu postanowiliśmy sprawdzić, czy nasz model będzie również działał w przypadku nagranych przez nas dźwięków. Do przetworzenia dźwięku posłużyliśmy się wcześniej utworzona przez nas funkcja. Zamienia ona dźwięk z formatu wav na ramkę danych, zawierająca (prawie) te same zmienne co te zawarte na stronie <https://apispreadsheets.com/datasets/119>. Na nasze nieszczęście, model źle przewiduje ostateczne wyniki. Postanowiliśmy więc dokonać dekompozycji predykcji dla naszych obserwacji i zauważyliśmy, że zmienna która najbardziej definiowała ostateczny wynik była zmienna `meanfun`, która w przypadku utworzonej przez nas ramki danych miała wartości nieadekwatne do tych, które były w przypadku pierwotnej ramki. To samo tyczy się wszystkich zmiennych dotyczących częstotliwości fundamentalnej oraz dominującej. Być może jest to kwestia sprzętu nagrywającego lub złego sposobu obliczania wartości tych zmiennych.