

GenderVoice

Drazkowski Hubert, Wilk Marcin

Kwiecień 2021

1 Wstęp

Poruszonym przez nas w projekcie problemem było przewidywanie płci na podstawie właściwości akustycznych ich głosu. Na potrzeby stworzenia bazy danych zostały nagrane 3168 próbki dźwiękowe pochodzące od kobiet lub mężczyzn. Nagrania zostały przetworzone przez pakiety `tuneR` i `seewave` w pakiecie R na przedziale częstotliwości 0hz - 280 hz. Analiza statystyczna tak stworzonych danych została wykonana w programie Python. Dla czytelności i zwiezłości raportu niektóre opisy metod lub rysunki wykorzystywane podczas analizy zostały pominięte. Można je znaleźć w notatniku jupyter.

Badanie rozpoczyna się od analizy eksploracyjnej, usuwamy zmienne współliniowe. Na początku zmienne deterministycznie współliniowe następnie zmienne wysoko ze sobą skorelowane. Następnie dokonujemy przekształceń kilku zmiennych w celu nadania im rozkładom brzegowym kształtu bardziej zbliżonego do rozkładu normalnego. Kolejnym krokiem jest poszukiwanie obserwacji odstających na dwa sposoby. Pierwszym sposobem jest analiza ich rozkładów brzegowych, drugim algorytm grupujący DBSCAN. Następnie tworzymy kilka ramek danych odzwierciedlających kolejne etapy przetwarzania danych, aby znaleźć odpowiedź jaki wpływ mają te etapy przetwarzania danych na wynik różnych klasyfikatorów. Kończymy nasze badanie aplikując beterie klasyfikatorów i porównując ich błąd predykcji mierzony odsetkiem poprawnie zaklasyfikowanych przypadków posługując się 10 krotną krosvalidacją.

2 Opis bazy danych

Baza danych to 21 zmiennych. Jedna zmienna jest zmienna nominalna wskazująca na płeć pozostałe 20 zmiennych to zmienne pochodzące z rozkładów ciągłych. Zmienne informują o właściwościach akustycznych, to jest można wśród nich znaleźć średnia częstotliwość fali, odchylenie standardowe częstotliwości, statystyki spektralne itp. W zbiorze danych nie występują braki danych. Dokładny opis zmiennych można znaleźć miedzy innymi pod adresem <https://data.world/ml->

```

RangeIndex: 3168 entries, 0 to 3167
Data columns (total 21 columns):
#   Column      Non-Null Count  Dtype
---  -
0    meanfreq    3168 non-null   float64
1    sd           3168 non-null   float64
2    median      3168 non-null   float64
3    Q25         3168 non-null   float64
4    Q75         3168 non-null   float64
5    IQR         3168 non-null   float64
6    skew        3168 non-null   float64
7    kurt        3168 non-null   float64
8    sp.ent      3168 non-null   float64
9    sfm         3168 non-null   float64
10   mode        3168 non-null   float64
11   centroid    3168 non-null   float64
12   meanfun     3168 non-null   float64
13   minfun      3168 non-null   float64
14   maxfun      3168 non-null   float64
15   meandom     3168 non-null   float64
16   mindom      3168 non-null   float64
17   maxdom      3168 non-null   float64
18   dfrange     3168 non-null   float64
19   modindx     3168 non-null   float64
20   label       3168 non-null   object

```

	meanfreq	sd	median	Q25	Q75	IQR	skew	kurt	sp.ent	sfm
count	3168.000000	3168.000000	3168.000000	3168.000000	3168.000000	3168.000000	3168.000000	3168.000000	3168.000000	3168.000000
mean	0.180907	0.057126	0.185621	0.140456	0.224765	0.084309	3.140168	36.568461	0.895127	0.408216
std	0.029918	0.016652	0.036360	0.048680	0.023639	0.042783	4.240529	134.928661	0.044980	0.177521
min	0.039363	0.018363	0.010975	0.000229	0.042946	0.014558	0.141735	2.068455	0.738651	0.036876
25%	0.163662	0.041954	0.169593	0.111087	0.208747	0.042560	1.649569	5.669547	0.861811	0.258041
50%	0.184838	0.059155	0.190032	0.140286	0.225684	0.094280	2.197101	8.318463	0.901767	0.396335
75%	0.199146	0.067020	0.210618	0.175939	0.243660	0.114175	2.931694	13.648905	0.928713	0.533676
max	0.251124	0.115273	0.261224	0.247347	0.273469	0.252225	34.725453	1309.612887	0.981997	0.842936

research/gender-recognition-by-voice. Próba jest zbilansowana, jako że mężczyźni i kobiety są reprezentowani po równo.

3 Eksploracyjna analiza i przetwarzanie danych

Większość tych zmiennych jest statystykami wygenerowanymi na podstawie nagrań głosów poszczególnych osób, stąd musimy uważać na (niekoniecznie unikać) tworzenie statystyk z innych statystyk, takich jak średnia ze średnich. Może to prowadzić do nieścisłości, gdyż nie mamy dostępu do początkowych próbek.

Zmienne mają w większości dość małe zakresy, mieszczą się co najwyżej między 0 a 1, z kilkoma wyjątkami. Duże wartości maksymalne mają również zmienne skew, maxdom i dfrange i kurt.

	mode	centroid	meanfun	minfun	maxfun	meandom	mindom	maxdom	dfrange	modindx
count	3168.000000	3168.000000	3168.000000	3168.000000	3168.000000	3168.000000	3168.000000	3168.000000	3168.000000	3168.000000
mean	0.165282	0.180907	0.142807	0.036802	0.258842	0.829211	0.052647	5.047277	4.994630	0.173752
std	0.077203	0.029918	0.032304	0.019220	0.030077	0.525205	0.063299	3.521157	3.520039	0.119454
min	0.000000	0.039363	0.055565	0.009775	0.103093	0.007812	0.004883	0.007812	0.000000	0.000000
25%	0.118016	0.163662	0.116998	0.018223	0.253968	0.419828	0.007812	2.070312	2.044922	0.099766
50%	0.186599	0.184838	0.140519	0.046110	0.271186	0.765795	0.023438	4.992188	4.945312	0.139357
75%	0.221104	0.199146	0.169581	0.047904	0.277457	1.177166	0.070312	7.007812	6.992188	0.209183
max	0.280000	0.251124	0.237636	0.204082	0.279114	2.957682	0.458984	21.867188	21.843750	0.932374

3.1 Współliniowość

W celu zbadania współliniowości posłużyliśmy się macierzami korelacji Spearmana i Pearsona oraz współczynnikiem VIF. Z oryginalnego zbioru danych zostały usunięte zmienne, które okazały się być kombinacjami liniowymi pozostałych zmiennych w zbiorze danych. Zarówno macierze korelacji jak i wyniki VIF obliczeń można znaleźć w notatniku jupyter. Poniżej znajduje się kilka statystyk, które wykluczyliśmy z dalszej analizy.

1. IQR jako, że jest to kombinacja liniowa Q25 i Q75
2. dfrange która okazała się współliniowa z maxdom
3. centroid będąca współliniowa względem meanfreq

Zostały także usunięte wysoce skorelowane ze sobą zmienne, jako że potencjalnie przekazywały te same informacje.

1. skew była mocno skorelowana z kurt
2. sfm jest najsilniej skorelowana z sp.ent oraz ma pozostałe korelacje dużo większe od sp.ent, zatem możemy się jej pozbyć z naszego zbioru.
3. Zmienna meanfreq jest z kolei bardzo silnie skorelowana z median i Q25. Teoretycznie statystyki te nie muszą być w ogóle ze sobą powiązane, lecz w naszym przypadku bardzo możliwe jest, że meanfreq jest kombinacją liniową tych dwóch zmiennych i możemy usunąć ją z modelu.
4. Niespodzianka jest silna korelacja ujemna występuje pomiędzy sd a Q25, jednak przy usunięciu meanfreq oraz centroid musimy zostawić w naszym zbiorze Q25, a możemy usunąć sd.

3.2 Przekształcenia zmiennych

Dokonaliśmy kilku przekształceń zmiennych, by były bardziej podobne do rozkładów normalnych. Takimi przekształceniami objeśliśmy zmienna kurt, maxdom, meandom

$$kurt := \log(\log(kurt)),$$

$$maxdom := \log(maxdom + 1).$$

$$meandom := \sqrt{meandom}$$

Dodatkowo ze względu na fakt, że w przypadku zmiennej *mode* jej rozkład dla obu płci ma zupełnie różne maksima, kształty oraz punkty koło których gromadziło się najwięcej obserwacji. Stworzyliśmy zmienną nominalną indykatorową według reguły

$$mode_trans := \mathbb{1}_{(mode \in [0.08, 0.15])}(mode)$$

3.3 Obserwacje potencjalnie odstające

Okazuje się, że dobrym przekształceniem dla zmiennej *meandom* jest pierwiastek, wtedy rozkład ten staje się niemal symetryczny i przypomina nawet w pewnym sensie krzywa Gaussa. Martwi nas jedynie duży słupek w okolicach zera. Sprawdziliśmy czy nie jest to jedna wartość co wskazywałoby na to, że dane zostały w pewien sposób ucięte do tej wartości (np. była to minimalna czułość przyrządów pomiarowych). Aż 61 wartości przyjmuje 0.088388. Drugi najliczniejszy wynikiem jest 0.265165 z 4 reprezentantami tej wartości. Podobne zjawisko pojawia się w przypadku zmiennej *modindx*. Tam 0 jest reprezentowane 65 razy. W obu przypadkach w proporcjach 1:3 są to głosy kobiet i mężczyzn.

W naszym zbiorze mamy sporo skośnych rozkładów, więc tam skrajne obserwacje wcale nie muszą być odstające, wynikają one raczej ze specyfiki rozkładów. Jedynymi pozostałymi zmiennymi gdzie można zauważyć jakieś obserwacje odstające są zmienne *Q75* oraz *kurt*. *Q75* ma bardzo cienki ogon z lewej strony, za to brak go z prawej, dodatkowo jest to rozkład mocno skupiony na krótkim odcinku, stąd nasze podejrzenie. *kurt* ma z kolei bardzo dziwna "górkę" w histogramie z prawej strony, zupełnie nie pasuje ona do rozkładu. Sprawdźmy czy rozkłady te są w jakiś sposób normalne testem Shapiro-Wilka i czy możemy zastosować regułę 3 sigm. Zarówno *Q75* jak i *kurt* mają bardzo małe *p* wartości, rzędu 10^{-31} . Jak widać z *p*-wartości obu testów nie mamy żadnego prawa uznać je za normalne.

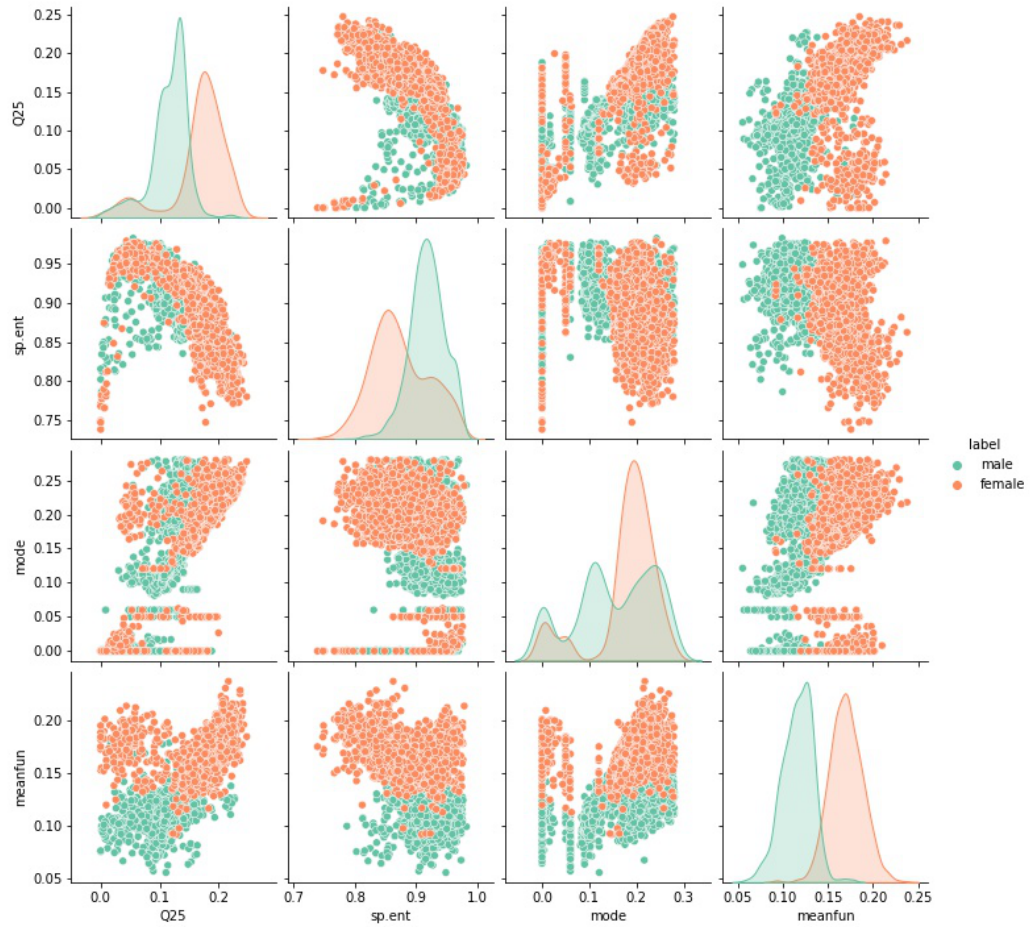
Na następnym etapie posłużymy się nierównością Markowa. Z tego powodu wycentrowaliśmy najpierw te zmienne. Chcemy sprawdzić czy odstające są obserwacje mniejsze od -0.07 dla *Q75* oraz większe od 1 dla *kurt*. Z nierówności Markowa wiemy, że

$$\mathbb{P}(|X| > eps) \leq \frac{\mathbb{E}(X^2)}{eps^2} = \frac{Var(X)}{eps^2} = p_X$$

Wartości dla *Q75* i *kurt* to odpowiednio 0.114 oraz 0.129. Zatem tak się szacują prawdopodobieństwa, że wartości bezwzględne wychodzą poza te wartości. Nie są one może rzędu 0.05, należy jednak pamiętać, że jest to zazwyczaj dość grube oszacowanie i sporo powyżej faktyczna wartość prawdopodobieństwa.

Posługując się kryteriami z przykładu nierówności Markowa, z analizy uciętych wartości *meandom* i *modindx* uzyskaliśmy 121 obserwacji potencjalnie odstających. Detale ujęte są w notatniku jupyter.

W poszukiwaniu ciekawych algorytmów wyróżniania obserwacji odstających natrafiliśmy na algorytm DBSCAN. Nie przyjmuje on założeń na temat rozkładu z



którego pochodzą zmienne. Wyróżnia on jedynie obserwacje znacznie oddalone od reszty przy założeniu pewnej metryki mierzącej odległość punktów w wielu wymiarach. Algorytm dokonuje klastrowania punktów do siebie zbliżonych i wyróżnia te obserwacje, które nie zostają przyporządkowane do żadnej z klas, które są od pozostałych punktów znacznie oddalone. Dzięki algorytmowi również zostało wyróżnionych 121 obserwacji odstających. Jest to natomiast zbieg okoliczności ponieważ tylko 3 obserwacje pokrywają się pomiędzy zbiorami obserwacji odstających.

3.4 Ważne zmienne

Niektóre zmienne wydają się bardzo dobrze separować klasy kobieta i mężczyzna. Są to zmienne meanfun, Q25, sp.ent. Ciekawa zależność widać tutaj w przy-

padku zmiennej mode. Rozkłady warunkowe kobiet i mężczyzn w tym przypadku różnią się tym, że w zakresie zmiennej mode od 0.08 a 0.15 kobiece głosy pojawiają się bardzo sporadycznie, zaś jeżeli chodzi o rozkład mężczyzn to osiągnięta jest w tym przedziale jego moda.

4 Modelowanie

4.1 Opis badania

Stworzyliśmy kilka wersji ramek danych. Zobaczmy czy element preprocessingu danych istotnie wpływa na jakość klasyfikacji różnych algorytmów, jeżeli tak, to na które algorytmy.

Nazwa ramki	Operacje
df1	surowy zbiór danych z usuniętymi deterministycznie współliniowymi zmiennymi
df2	df1 z dodatkowo usuniętymi wysoce skorelowanymi zmiennymi
df3	df2 z przekształceniami w celu normalizowania
df4	df3 z usuniętymi obserwacjami odstającymi na podstawie analizy rozkładów brzegowych
df5	df3 z usuniętymi obserwacjami na podstawie DBSCAN

Drugą częścią badania jest porównanie jakości klasyfikatorów. Porównujemy ze sobą kilka prostych modeli regresji logitowej jako modele benchmarkowe. Następnie model logitowy z regularyzacją, lasy losowe, SVM, XGBoost, oraz komitet modeli QDA, regresji logistycznej i naiwnego klasyfikatora bayesa. Porównanie przebiega za pomocą prostego zliczania odsetka poprawnych wskazań w 10-krotnej krosvalidacji.

4.2 Wyniki badania

Po pierwsze zauważyliśmy, że transformacja rozkładów w celu ich znormalizowania nie wpłynęła na jakość klasyfikacji używanych przez nas algorytmów. Odkryliśmy także, że wystandaryzowanie zmiennych polepsza jakość klasyfikacji. W tabeli poniżej przedstawiamy nazwę modelu i średnie accuracy dla df1, df2, df4 i df5. Pierwsza tabela zawiera wyniki przed standaryzacją kolumn, druga tabela zawiera wyniki po standaryzacji kolumn. Wyniki to mediana accuracy. Więcej statystyk i dokładniejsza analiza znajduje się w notebooku. Przed wystandaryzowaniem kolumn

Nazwa modelu	df1	df2	df4	df5
Regresja logistyczna sam intercept	0.498423	0.498423	0.510672	0.511475
Regresja logistyczna sam meanfun	0.954259	0.954259	0.958941	0.957302
Regresja logistyczna z kara l2	0.894147	0.908373	0.880118	0.875210
Lasy losowe	0.971609	0.966877	0.975410	0.972131
Komitet	0.958991	0.954259	0.970492	0.970492
SVM	0.802215	0.958936	0.960591	0.963881
XGBoost	0.973186	0.973186	0.975383	0.977022

Po wystandaryzowaniu kolumn

Nazwa modelu	df1	df2	df4	df5
Regresja logistyczna z kara l2	0.970032	0.966877	0.980328	0.978689
Lasy losowe	0.971609	0.966877	0.975410	0.972131
Komitet	0.966877	0.965300	0.980296	0.977017
SVM	0.976341	0.962145	0.981946	0.975388
XGBoost	0.973186	0.973186	0.977022	0.975383

5 Wnioski

Standaryzowanie zmiennych wpłynęło w następujący sposób na jakość klasyfikacji poszczególnych modeli:

1. Regresja logistyczna znacząco zyskała na wystandaryzowaniu zmiennych. Mediana przed zabiegiem wynosiła na poszczególnych zbiorach od [88%, 91%] teraz jest na poziomie [97%,99%]
2. Lasy losowe w żaden sposób nie zyskały na standaryzacji
3. Mediana jakości klasyfikacji dla komitetu modeli wzrosła od [95% ; 97%] do [96% ; 97%]
4. Mediana jakości klasyfikacji dla SVM po wszystkich ramach danych wzrosła od [80%;96%] do [97%;98%]
5. Dla XGBoosta mediana jakości klasyfikacji nie zmieniła się znacząco.

Wydaje się, że prosty model z samym meanfun jest potężny osiągając niebagatelne 95% jakości stabilnie na wszystkich ramach danych.

Standaryzacja kolumn nie zmienia wyników dla lasów losowych, xgboosta. Nieznacznie wpływa na klasyfikacje komitetów, natomiast ma kolosalne znaczenie w przypadku SVM i regresji logistycznej z kara RIDGE.

Jakie znaczenie ma pozostały preprocessing ?

Średnio usuwanie obserwacji odstających na podstawie ich rozkładów brzegowych radzi sobie lepiej niż wielowymiarowy DBSCAN. Usuwanie obserwacji odstających pozwala także na poprawę klasyfikacji względem zbioru tego przekształcenia pozbawionego, ale nie zawsze znacząco. Usuwanie zmiennych silnie

skorelowanych, ale nie idealnie współliniowych pogarsza jakość klasyfikacji.

Najlepiej klasyfikuje mniej więcej 1. Regresja logistyczna z regularyzacją 2.SVM 3. XGBoost 4. Komitet 5. Lasy losowe . Aczkolwiek nie jest to żadna ścisła dominacja. Są przypadki ramek danych dla których taka kolejność mogłaby być myląca