

WUM - Raport z projektu 1 - Klasyfikacja

Przemysław Chojecki, Michał Wdowski

28 04 2020

1 Opis

W ramach projektu zdecydowaliśmy się na analizie zbioru danych medycznych zawierającego dane o kobietach ze szpitala z Wenezueli w różnym wieku. U niektórych z nich zidentyfikowano raka szyjki macicy. Dane zawierają informacje z przeprowadzonych ankiet oraz badań medycznych.

Celem projektu było stworzenie modelu przewidującego wynik testu na którykolwiek z raków na podstawie takich informacji jak: Wiek, ilość partnerów seksualnych w przeszłości, wiek inicjacji seksualnej, intensywność palenia papierosów, czas korzystania z antykoncepcji hormonalnej oraz wyników przeprowadzonych na pacjentkach badań na choroby AIDS, syfilis(kiła), Hepatitis B, kłykcina i tym podobne.

Projekt skupia się na trzech obszarach: analizie zbioru i danych o pacjentach, inżynierii cech i przekształcenie zbioru jak najlepiej dla modelowania oraz modelowanie i wybór najlepszego klasyfikatora na podstawie kilku różnych miar. Dobrze przewidzenie choroby u pacjenta jest ważne, dlatego należy znaleźć najlepszy klasyfikator. Jako najważniejsza uznaliśmy miarę "Weighted TPR-TNR Measure" opisana w załączonym artykule (pod nazwą W_R.pdf), gdyż jest ona tam opisana jako najlepsza do oceny danych niezbalansowanych.

2 Analiza zbioru

Oryginalny zbiór zawiera wyniki ankiet 858 pacjentek ze szpitala w Caracas w Wenezueli. Informacje dotyczące każdej z kobiet podzieliliśmy na 3 typy:

1. numeryczne, czyli takie, których wartości są różnymi liczbami dodatnimi. W oryginalnym zbiorze jest takich 12.
2. kategoryczne, których wartościami jest prawda lub fałsz. W oryginalnym zbiorze jest takich 20.
3. celu, wynik danego badania na chorobę szyjki macicy. W oryginalnym zbiorze jest takich 4.

Do analizy dla każdej z tych grup podchodziliśmy trochę inaczej.

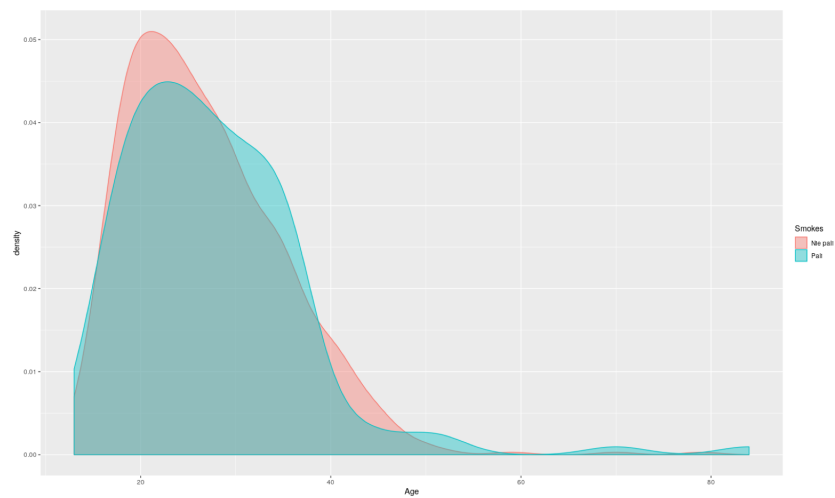


Figure 1: Rozkład wieku w zależności od palenia

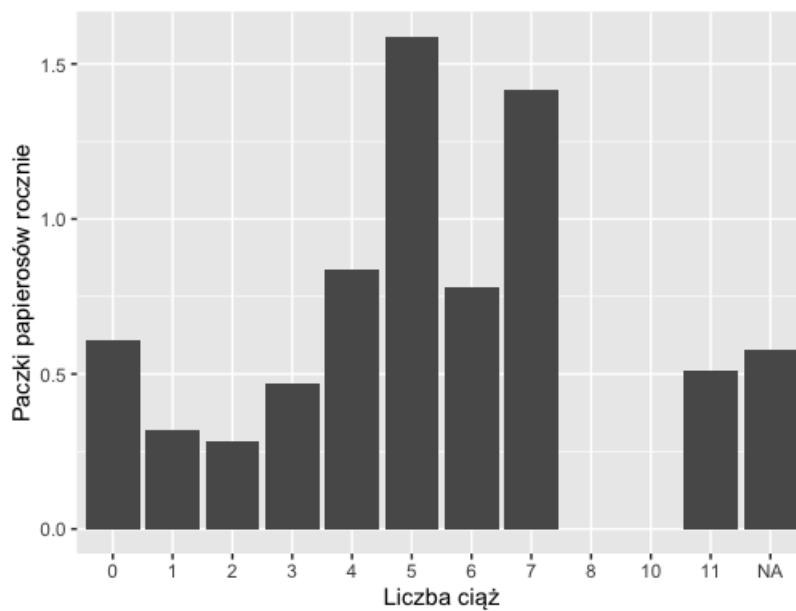


Figure 2: Rozkład palenia w zależności od liczby ciąż

Na pierwszy wykresie widzimy rozkład wieku w zależności od tego, czy dana kobieta pali, czy nie. Spodziewaliśmy się zobaczyć dużo starsze kobiety wśród palaczek, jednakże tak niema.

Na drugim wykresie widzimy jak dużo średnio pala kobiety w zależności jak

dużo razu były w ciąży. Wygląda na to że ludzie bez dzieci palą trochę, jak mają mało dzieci to palą mniej, a potem stresu jest za dużo i palą więcej. Uwaga - dla osób z ciążami 8 i więcej są tylko 1-2 osoby w każdej grupie.

Jesteśmy bardzo zdziwieni niską ilością palonych przez kobiety paczek papierosów. Wnioskujemy, że być może w Wenezueli standardy palenia są niższe niż w Polsce i stąd ta dysproporcja.

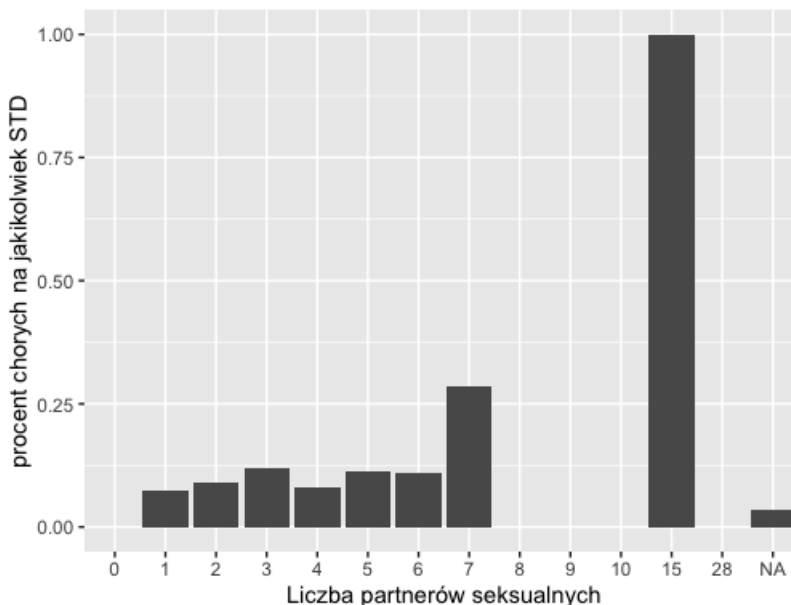


Figure 3: Rozkład wieku w zależności od palenia

Na trzecim wykresie widzimy jaką część ludzi ma choroby STD w zależności od liczby partnerów seksualnych. Uwaga - dla osób z liczbą partnerów seksualnych 8 i więcej robi się po jednej osobie na daną liczbę partnerów

Wiele z danych są brakujące (kwestja ta będzie dokładnie omówiona w następnym paragrafie). W celach późniejszej analizy i modelowania zamieniliśmy je na wartości za pomocą algorytmu z pakietu "mice".

Na koniec ciekawostka. Jeden z wierszy przedstawia się następująco: 16-sto latka posiadająca 28 byłych partnerów seksualnych, z czego pierwszy z kontaktów seksualnych odbyła w wieku 10-ciu lat. Co jeszcze ciekawsze, nigdy nie brała antykoncepcji, ale tylko raz była w ciąży. Od 5-ciu lat pali. Testy nie wykryły choroby na żadną z chorób.

3 Inżynieria cech

Najważniejszą kwestją w tej części było poradzenie sobie z danymi brakującymi. Jak się okazało, najważniejsze dla modelowania cechy, czyli te zawierające

wyniki testów, są wybrakowane. Dokładnie 100 kobiet w posiadanych przez nas danych nie posiada wyników testów na żadne z chorób, co oznacza, że jedyne informacje o tych 100 kobietach jakie posiadamy, to wyniki ich ankiet. Zdecydowaliśmy się usunąć te dane, co było trudną decyzją, gdyż jest to aż 12% posiadanych przez nas danych.

Zauważyliśmy istnienie kilku danych niepoprawnych, pozbyliśmy się ich. Polegało to na przykład na tym, że pacjetka w ankiecie twierdziła, że miała inicjację seksualną w wieku większym niż jest jej obecny.

Postanowiliśmy również pozbyć się kolumn, które powielają posiadane już informacje, lub były jej bardzo zbliżone.

Dwie z kolumn posiadały aż 91% danych brakujących, zastanawialiśmy się więc nad jej usunięciem. Jednakże doszliśmy do wniosku, że mają one jednak sens. Kolumny te nosiły nazwy: "czas od ostatniej diagnozy" oraz "czas od pierwszej diagnozy". Brak danych w nich zinterpretowaliśmy jako niedbanie przez pacjetkę żadnej diagnozy, więc zastąpiliśmy te braki wartościami "0". Poza tym pozbyliśmy się kolumny "czas od ostatniej diagnozy", gdyż była praktycznie taka sama jak "czas od pierwszej diagnozy". Pozbyliśmy się również kolumn z informacją o chorobie AIDS, gdyż żadna z pacjetek nie była chora, oraz kilku innych niewnoszących wiele do analizy kolumn.

W tym miejscu zdecydowaliśmy się wiele ze zmiennych ciągłych (numerycznych) zamienić na zmienne kategoryczne, tworząc tak zwane "kubelki". Decyzje o punktach tak zwanych "cięć" podejmowaliśmy później na podstawie wyników dostosowanych modeli.

Dopiero po zredukowaniu wymiaru danych, na sam koniec etapu inżynierii cech, pozostałe brakujące dane wypełniliśmy za pomocą funkcji z pakietu "mice".

4 Modelowanie

Większa część tego procesu przeprowadzaliśmy w załączonej, utworzonej przez nas aplikacji shiny. Umożliwiła ona nam łatwe i szybkie analizowanie wyników naszych modeli i dostosowywanie parametrów i kubelków z etapu inżynierii cech.

TODO(MOZE, wkleic zdjecia a appki?) Działanie aplikacji jest następujące:

1. Aplikacja wczytuje dane i wykonuje obróbkę opisaną w poprzednich dwóch paragrafach na podstawie wybranych przez użytkownika parametrów
2. Aplikacja dzieli zbiór danych na treningowy (70%) walidacyjny (15%) i testowy (15%). Wykonuje standaryzację danych treningowych, po czym za pomocą średnich i wariancji ze zbioru treningowego, stara się ustandaryzować dane ze zbiorów treningowego i walidacyjnego.
3. Jeśli użytkownik tego sobie życzy (widnieje to pod nazwą "Create more positive data") aplikacja sztucznie zmnoży i trochę zaburzy nowe dane kobiet ze zdiagnozowanym rakiem. My w naszych modelach zdecydowaliśmy się na tę opcję, gdyż oryginalna kolumna celu w danych

jest bardzo niebalansowana, a to utrudnia modelom proces uczenia. Dzięki takiemu małowemu oszustwu model miał łatwiej w zrozumieniu, że ważniejszym dla nas było, żeby model wychwytywał kobiety chore, a mniej, żeby udawało mu się ze zdrowymi. TODO(przetłumaczyć to zdanie na Polski)

4. Użytkownik decyduje się który z modeli chciałby nauczyć. Do wyboru ma algorytmy:
 - (a) kkm
 - (b) las losowy
 - (c) regresję logistyczną TODO(Sprawdź, czy tak to się na pewno nazywa)
 - (d) drzewo decyzyjne
5. Następnie aplikacja wykonuje tyle modeli, o ile użytkownik poprosił w opcji "Find best of" i porównuje je ze sobą na podstawie wyników miary Weighted TPR-TNR liczonej na zbiorze walidacyjnym. Każdy z modeli jest liczony w systemie kroswalidacji takiej o jaką prosił użytkownik. TODO(Przetłumacz końcówkę na Polski)
6. Na koniec wyświetlane są użytkownikowi wyniki testów modelu wykonane na zbiorze testowym, czyli takim, którego model nigdy nie widział i nie dostosowywał się do niego. Wylistowane są wyniki miar:
 - (a) AUC
 - (b) FBeta score
 - (c) Weighted TPR-TNR
 - (d) AUPRC
 - (e) tabela liczby poprawnych/niepoprawnych klasyfikacji

Z ważnych rzeczy: każdy z dostępnych modeli zwracał jako wynik ciąg prawdopodobieństw, że dany rekord będzie chory. Zdecydowaliśmy się na "uciecie" chorych w punkcie, który jest m -tym od góry prawdopodobieństwem, gdzie $m = n * \text{procentowy udział targetu w zbiorze treningowym}$, gdzie n to liczba zmiennych w zbiorze testowym, lub walidacyjnym (wartości te są takie same). TODO(spytać Michała, czy to jest dobrze wytłumaczone) Dzięki temu symulujemy w testowanym zbiorze taką ilość chorych, jaka była w danych uczących. Taki rodzaj interpretacji sprawdza się przy sprawozdaniach tego typu, jednakże jest bezużyteczny "na produkcji", gdzie dane napływają w sposób ciągły, a nie w paczkach, co jest wymagane w naszym sposobie. Aktualny wkład procentowy chorych można podejrzeć z lewej strony aplikacji, a aktualnie wybrane prawdopodobieństwo graniczne w części środkowej nad wynikami modelu.

Na koniec ostrzeżenie dla użytkownika. Aplikacja czasem się buguje, dlatego zamieszczony został łatwodostępny przycisk "RESET", którego naciśnięcie zazwyczaj naprawia problem.

Za pomoca aplikacji ustaliliśmy najbardziej optymalny układ "kubeków" i innych parametrów, oraz uznaliśmy, że najlepiej sprawdzającymi sie modelami sa knn i las losowy.

Ostatecznie dostrajaliśmy parametry modeli za pomoca randomsearch i wyniki sa nastepujace: TODO(Wkleić zdjecia krzywej ROC i cyferki miar)

5 Wnioski

TODO(Wymyślić wnioski, najlepiej nie kompletnie negatywne xd)