

Metody Eksploracji Danych - Projekt - Analiza Eksploracyjna

Jakub Matłacz, Kamil Trysiński

3 stycznia 2022

Streszczenie

Termin – 3 stycznia 2022r. Ilość punktów do zdobycia – 20 pkt. Analiza eksploracyjna posiadanego zbioru/wycinka zbioru danych i postawienie tezy/zadania badawczego mającego na celu opracowanie modelu regresji dla opisanych danych. Celem analizy eksploracyjnej danych jest sprawdzenie zależności między posiadanymi danymi (5 pkt.), zbadanie ich zakresów i stopnia zmienności (5 pkt.), analiza stopnia wypełnienia danych (5 pkt.), wizualizacja (5 pkt.). Wynikiem tych analiz, ma być raport zakończony postawieniem hipotezy badawczej mającej na celu znalezienie relacji między zmiennymi objaśniającymi, a zmienna objaśniana.

1 Tematyka Projektu

Celem projektu było znalezienie właściwego zbioru danych do jego analizy i postawienia hipotezy. Dalszym celem tych działań jest budowa efektywnego modelu regresji przewidującego wartość wybranej zmiennej opisywanej na podstawie wartości pozostałych cech.

2 Omówienie wybranego zbioru danych

2.1 Źródło

Wybrano zbiór <https://data.world/data-society/student-alcohol-consumption>. Strona data.world udostępnia wiele ciekawych zbiorów danych i stanowi dobre ich źródło.

2.2 Omówienie cech

Poniżej przedstawiono omówienie znaczenia cech w zbiorze oryginalnym. Na dalszym etapie część z nich nie będzie potrzebna. Jak widać zbiór oryginalny zawiera 34 cechy z czego wiele jest kategoriowych. Poza tym nawet jeśli liczbowe to są to wartości dyskretne z jakiegoś zbioru możliwych wartości.

1. student's school (binary: 'GP' - Gabriel Pereira or 'MS' - Mousinho da Silveira)
2. sex - student's sex (binary: 'F' - female or 'M' - male)
3. age - student's age (numeric: from 15 to 22)
4. address - student's home address type (binary: 'U' - urban or 'R' - rural)
5. famsize - family size (binary: 'LE3' - less or equal to 3 or 'GT3' - greater than 3)
6. Pstatus - parent's cohabitation status (binary: 'T' - living together or 'A' - apart)
7. Medu - mother's education (numeric: 0 - none, 1 - primary education (4th grade), 2 - 5th to 9th grade, 3 - secondary education or 4 - higher education)
8. Fedu - father's education (numeric: 0 - none, 1 - primary education (4th grade), 2 - 5th to 9th grade, 3 - secondary education or 4 - higher education)
9. Mjob - mother's job (nominal: 'teacher', 'health' care related, civil 'services' (e.g. administrative or police), 'at_home' or 'other')

10. Fjob - father's job (nominal: 'teacher', 'health' care related, civil 'services' (e.g. administrative or police), 'at_home' or 'other')
11. reason - reason to choose this school (nominal: close to 'home', school 'reputation', 'course' preference or 'other')
12. guardian - student's guardian (nominal: 'mother', 'father' or 'other')
13. traveltime - home to school travel time (numeric: 1 - 15 min., 2 - 15 to 30 min., 3 - 30 min. to 1 hour, or 4 - 1 hour)
14. studytime - weekly study time (numeric: 1 - 2 hours, 2 - 2 to 5 hours, 3 - 5 to 10 hours, or 4 - 10 hours)
15. failures - number of past class failures (numeric: n if 1=n3, else 4)
16. schoolsup - extra educational support (binary: yes or no)
17. famsup - family educational support (binary: yes or no)
18. paid - extra paid classes within the course subject (Math or Portuguese) (binary: yes or no)
19. activities - extra-curricular activities (binary: yes or no)
20. nursery - attended nursery school (binary: yes or no)
21. higher - wants to take higher education (binary: yes or no)
22. internet - Internet access at home (binary: yes or no)
23. romantic - with a romantic relationship (binary: yes or no)
24. famrel - quality of family relationships (numeric: from 1 - very bad to 5 - excellent)
25. freetime - free time after school (numeric: from 1 - very low to 5 - very high)
26. goout - going out with friends (numeric: from 1 - very low to 5 - very high)
27. Dalc - workday alcohol consumption (numeric: from 1 - very low to 5 - very high)
28. Walc - weekend alcohol consumption (numeric: from 1 - very low to 5 - very high)
29. health - current health status (numeric: from 1 - very bad to 5 - very good)
30. absences - number of school absences (numeric: from 0 to 93)
31. These grades are related with the course subject, Math or Portuguese:
32. G1 - first period grade (numeric: from 0 to 20)
33. G2 - second period grade (numeric: from 0 to 20)
34. G3 - final grade (numeric: from 0 to 20, output target)

2.3 Wyzwania zbioru

Wybraliśmy ten zbiór, ponieważ posiada kilka ciekawych cech stanowiących ciekawe wyzwanie.

1. Duża liczba cech.
2. Zawiera cechy kateryczne o wielu możliwych wartościach.
3. Zawiera cechy binarne.
4. W ogólności wszystkie cechy (również po przekształceniach) są dyskretne, a nie ciągłe (liczbowe).
5. Zbiór nie jest zbalansowany względem interesującej nas docelowej cechy opisywanej.

3 Wstępne przygotowanie danych

Dane są tylko połączone w zbiór, który jest nam potrzebny, ale należy traktować je jako dane oryginalne, czyli niepoddane normalizacji, czy innym inwazyjnym technikom. Jedyne co zrobiono to usunięto/połączono pewne kolumny, zamieniono atrybuty katagoryczne na ilościowe w celu umożliwienia dalszej analizy zbioru i uczenia modelu, zmieniono nazwy kolumn i posortowano w celu większej przejrzystości w dalszych wizualizacjach, zamieniono typy danych na int64, ponieważ wszystkie wartości były całkowite (dane składają się z wielu atrybutów dobieranych przez człowieka w wyniku rozmowy z uczniem, więc są to pewne wskaźniki). Ponadto połączono picie w tygodniu i weekend we wspólny wskaźnik problemu alkoholowego z wyższą wagą dla picia w tygodniu.

3.1 Edycja zbioru oryginalnego

Wczytano oba pliki tabelki dla 2 przedmiotów obserwowanego kursu. Następnie połączono je ignorując index. Usunięto `columns=["G1","G2","G3","paid"]`, ponieważ są to cechy różne w obu tabelkach dla tego samego ucznia, a chcemy w następnym kroku usunąć duplikatów uczniów (w dokumentacji zbioru jest opisane, że uczniowie w obu tabelach się powtarzają, ale różnią kursami). Następnie łączę picie alkoholu w tygodniu z piciem w weekend jako sumę ważoną, gdzie picie w tygodniu ma wagę 2.5, a w weekend wagę 1. Po operacjach połączenia picia oraz usunięcia duplikatów usuwam resztę niepotrzebnych cech `columns=["school","Dalc","Walc"]`. Gdzie `school` oznacza szkołę, do której chodzi uczeń, co dla nas nie ma znaczenia. Ponieważ chcemy, aby nasz model zwracał ryzyko alkoholowe dla każdego ucznia (nie tylko chodzącego do jednej z tych kilku szkół ze zbioru).

3.2 Binarizacja cech katagorycznych

Część cech było łatwo zbinaryzować `['sex','address','famsize','Pstatus']`, ponieważ należało jedynie zamienić nazwy typu tak/nie na 1/0, bo były one naturalnie binarne, a trzeba było jedynie zamienić je na numeryczne. Dla cech `features.to_encode = ['Mjob','Fjob','reason','guardian']` należało użyć techniki `one-hot-encode`. Biblioteka `sklearn` posiada większość narzędzi do budowy klasycznych modeli. Nic więc dziwnego, że miała również to i bardzo dobrze opisane w dokumentacji. Technika ta zamienia każdy pojedynczy atrybut katagoryczny na zbiór atrybutów binarnych. Po jednym dla każdej możliwej wartości cechy oryginalnej. Technika ta mocno poszerza `dataframe`, jednak jest o tyle lepsza od zwykłej zamiany kategorii w kolejne liczby, że minimalizuje zwiększanie wpływu jednej wartości cechy nad inną.

3.3 Zmiany dla wygody

Następnie zamieniłem dla pewności wszystkie dane na int64 (po sprawdzeniu, że taki jest nasz zbiór). Zamieniłem również nazwy cech na z małej litery, aby następnie posortować kolumny `dataframe` alfabetycznie, bo ułatwia dalszą analizę, np. w przypadku analizy macierzy korelacji łatwiej szybko znaleźć cechy, której się szuka.

4 Braki w danych

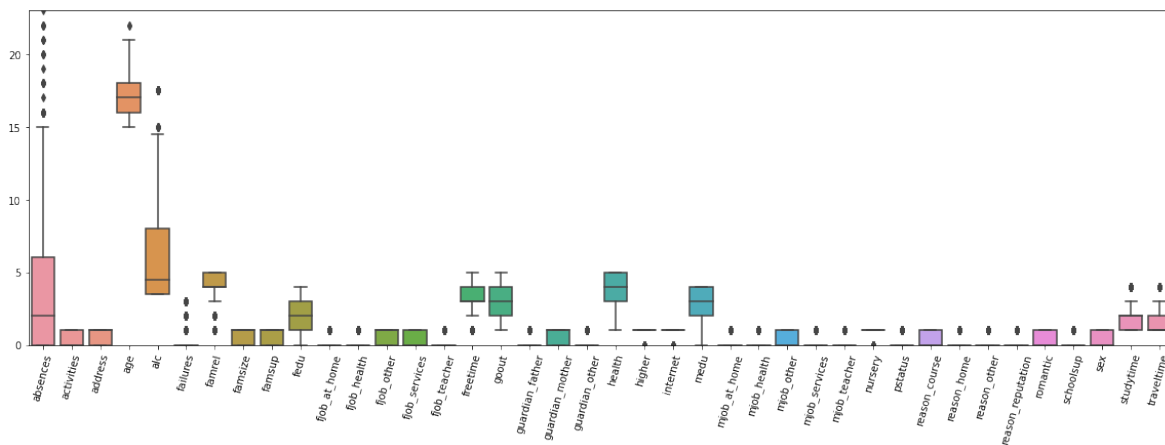
W zbiorze nigdy nie było i również po prostych edycjach nie ma braków danych.

5 Wizualizacja danych

Dane, które wizualizuje pochodzą już ze zbioru nieco przekształconego do naszych potrzeb jak opisano powyżej. Nie poddano ich jednak żadnym inwazyjnym metodom.

5.1 Wykresy pudełkowe

Najwięcej punktów oddalonych ma cecha nieobecności i trzeba było je przyciąć, aby nie popsuć widoczności wykresu. Inną opcją była zamiana skali na log ale uważam, że w tym przypadku nie miałyby



Rysunek 1: Wykresy pudełkowe cech

to sensu, bo ciężko jest odczytać dokładne zakresy zmienności cech, które są dość małe. Naturalnie wiele atrybutów binarnych (powstałych z kategoriycznych) ma zakres zmienności 0,1 bo przyjmuje tylko te 2 wartości. Niezdania mają sporo punktów oddalonych, co ma sens, bo to jedna z tych rzeczy, które na poziomie liceum dotyczą tylko nielicznych uczniów. Tak samo relacje rodzinne. Co ciekawe czas wolny ma punkt oddalony. Prawdopodobnie uczeń z dużą ilością zajęć dodatkowych. Dodatkowo warto wspomnieć, że skoro dane były pobierane przez nauczyciela od ucznia na podstawie rozmowy to na pewno są mocno zaszumione przez subiektywne odczucia ucznia odnośnie np. poziomu zdrowia czy relacji z rodziną. Ludzie mają tendencję do narzekania w tych sprawach co może niesłusznie pogarszać wyniki itp. Jak widać najwięcej punktów oddalonych mają cechy takie jak:

1. nieobecności - większość uczniów prawie nie ma nieobecności,
2. niezadania do następnej klasy - większości uczniów nigdy nie zdarzyło się niezdać do następnej klasy,
3. relacje rodzinne - większość uczniów ma dobre relacje rodzinne.

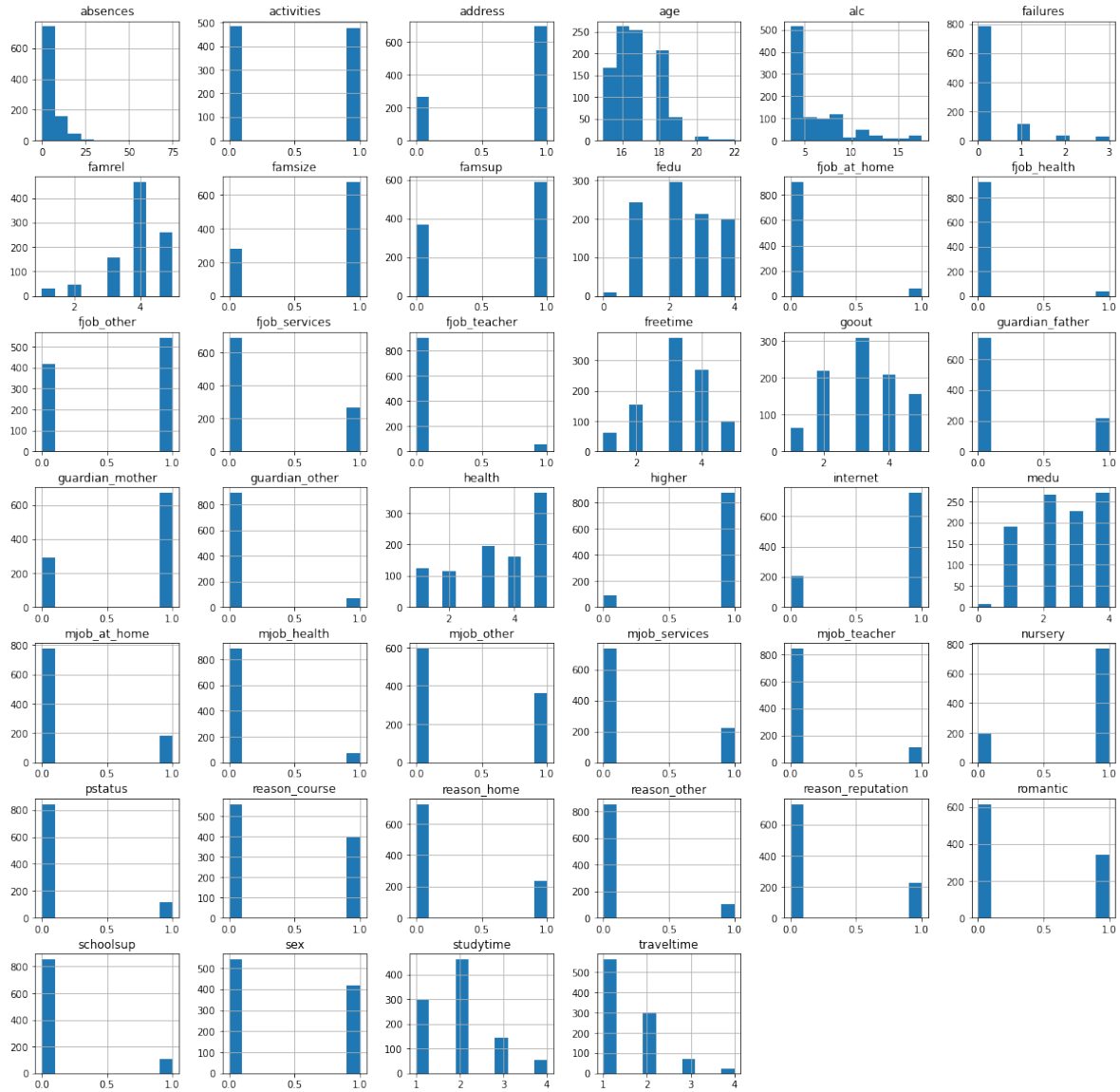
Możliwe, że odchylenia od normy w tak ważnych dla rozwoju dziecka metrykach będą mieć spory wpływ na ryzyko alkoholowe.

5.2 Histogramy cech

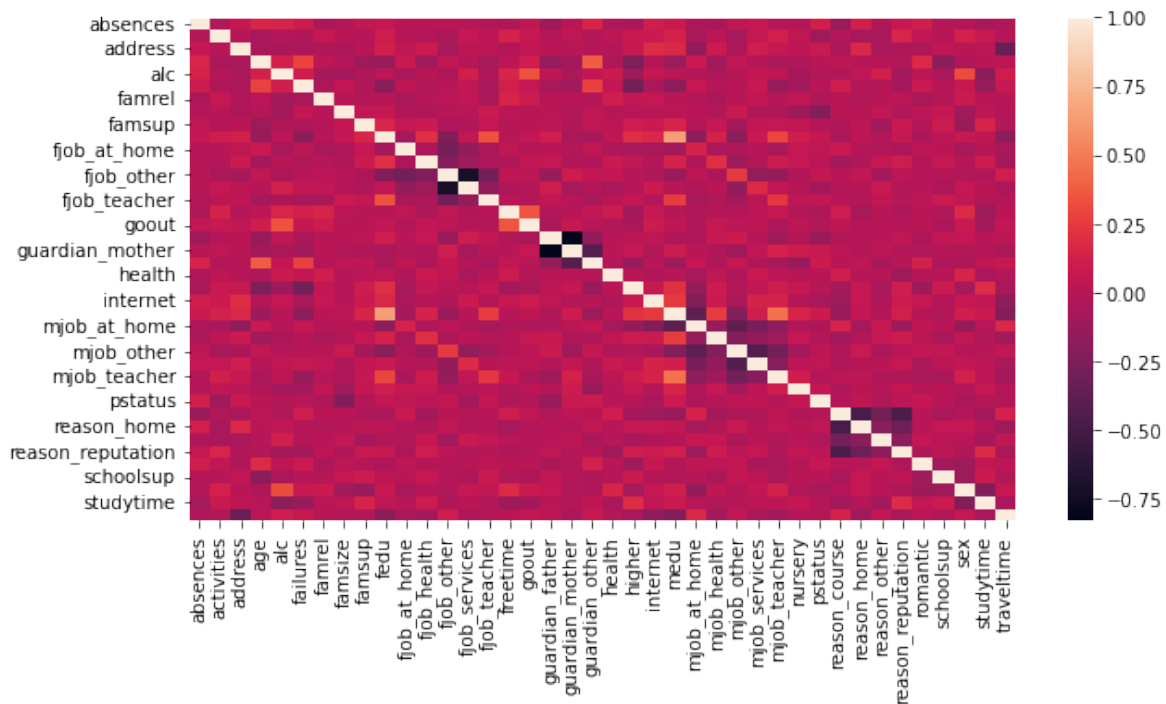
Poniżej widać rozkłady wszystkich cech. Jak widać wiele z nich jest binarnych lub ogólniej dyskretnych. Cechy nie są ciągłe, co może stanowić wyzwanie trenowania modelu, co nas zachęciło do wyboru tego zbioru.

6 Zbalansowanie zbioru

Można powiedzieć, że zbiór jest niezbalansowany względem naszej zmiennej opisywanej (odpowiednik klasy w klasyfikacji). Wynik maksymalny to 17.5. Połowa osób ma wynik mniejszy od 4.5. Čwierć osób ma wynik 4.5,8. Tylko pozostałe ćwierć osób ma wynik wyższy niż 8. Jest to jednak zrozumiałe i na całe szczęście tylko nieliczni mają problem w badanej próbie. Jak widać mamy około tysiąca obiektów w zbiorze. Zbiór naturalnie nie jest zbalansowany pod kątem zmiennej opisywanej `alc`. Przyjmuje ona wartości z zakresu 3.5,17.5. Wartości rosną wraz ze stopniem alkoholizmu. Tylko 25 procent uczniów pije więcej niż około połowa tego zakresu, czyli średni stopień problemu. Jeden uczeń (punkt oddalony) wyznacza górny limit skali stanowiąc dla nas naturalny punkt odniesienia (skala empiryczna). Poniżej widać opis zmiennej opisywanej.



Rysunek 2: Histogramy cech



Rysunek 3: Macierz korelacji

Cecha opisywana	
mean	6.094369
std	3.308938
min	3.500000
25 procent	3.500000
50 procent	4.500000
75 procent	8.000000
max	17.500000

7 Stopień korelacji cech

7.1 Opis najbardziej znaczących korelacji

1. opieka ojca albo matki w rozwiedzionych parach
2. edukacja matki i ojca
3. praca nauczyciela i edukacja dla matek i ojców ale w mniejszym stopniu
4. praca matki w domu, a wykształcenie
5. tak zwany inny opiekun dla uczniów, których wiek jest nietypowo wysoki w szkole średniej (dom dziecka, etc.)
6. wyjścia z domu, a alkohol i czas wolny
7. czas podróży tym wyższy jeśli mieszkają poza miastem
8. płeć meska i alkohol
9. wybór szkoły ze względu na reputację kontra ze względu na wygodę czyli bliskość od domu

Size of Correlation	Interpretation
.90 to 1.00 (−.90 to −1.00)	Very high positive (negative) correlation
.70 to .90 (−.70 to −.90)	High positive (negative) correlation
.50 to .70 (−.50 to −.70)	Moderate positive (negative) correlation
.30 to .50 (−.30 to −.50)	Low positive (negative) correlation
.00 to .30 (.00 to −.30)	negligible correlation

Rysunek 4: Interpretacja stopnia korelacji

7.2 Dokładne wartości

1. guardian_mother guardian_father -0.830369
2. fjob_services fjob_other -0.711511
3. medu fedu 0.637790
4. reason_home reason_course -0.477166
5. reason_reputation reason_course -0.467695
6. mjob_teacher medu 0.450380
7. mjob_services mjob_other -0.430532
8. guardian_other guardian_mother -0.425006
9. mjob_at_home medu -0.380988
10. mjob_other mjob_at_home -0.379828
11. guardian_other age 0.376413
12. goout freetime 0.342537
13. fjob_teacher fedu 0.341598
14. goout alc 0.336565
15. traveltime address -0.335628
16. sex alc 0.320500
17. reason_reputation reason_home -0.314566

8 Teza badawcza

Teza naturalnie pojawiała się już w powyższej części jednak powtórzę ją raz jeszcze. Celem będzie przewidywanie cechy alćstanowiacej sume ważona picia w tygodniu i weekendy na podstawie cech ucznia. Cechy to można traktować jako zmienne środowiskowe i psychologiczne. Mamy nadzieję, że model będzie na tyle dokładny, że posłuży jako narzędzie dla psychologów do automatycznego obliczenia stopnia ryzyka alkoholowego ucznia na podstawie badanych czynników pobranych w formie rozmowy/wywiadu z psychologiem szkolnym.