

Bank Marketing data (with social/economic context)

Maciej Maecki

25 października 2019

Streszczenie

W pliku Bank Marketing data.csv znajduj si dane charakteryzujce klientw pewnego banku oraz kampanie marketingowe skierowane do tych klientw. Doczone s ponadto wskaniki spoeczne i ekonomiczne. Na podstawie tych danych naley zbudowa model prognozujcy szans, e klient w wyniku prowadzonej kampanii zaoy lokat terminow.

Spis treści

1	Wprowadzenie	2
1.1	Opis problemu	2
1.2	Opis danych	2
1.3	Wstpna eksploracja danych	2
2	Analiza eksploracyjna	4
2.1	Age	4
2.2	Job	5
2.3	Marital status	5
2.4	Education	7
2.5	Has credit in default?	8
2.6	Has personal loan?	9
2.7	Has credit in default?	10

1 Wprowadzenie

1.1 Opis problemu

W ramach kampani marketingowej organizowanej przez pewien bank w latach między majem 2008 rok, a listopadem 2010 roku, były zbierane informacje na temat klientów tego banku. Na podstawie tych danych planowane jest przewidzenie, czy i jaki rodzaj klientów kupi lokat terminową w tym banku.

1.2 Opis danych

Nasze dane zawierają 21 kolumn danych. Kolumny możemy podzielić na 3 grupy:

I: Zmienne związane z danymi klienta bankowego:

1. Wiek (age): wiek klienta.
2. Praca (job): rodzaj pracy klienta.
3. Stan cywilny (marital): stan cywilny klienta.
4. Edukacja (education): edukacja klienta.
5. Domylnie (default): Klient wcześniej domylnie miał kredyt.
6. Mieszkanie (housing): Klient ma kredyt mieszkaniowy.
7. Pożyczka (loan): Klient ma osobistą pożyczkę.

II: Zmienne związane z ostatnim kontaktem bieżącej kampanii marketingowej:

8. Kontakt (contact): Typ komunikacji kontaktowej (telefonicznej lub komrkowej).
9. Miesiąc (month): Ostatni kontakt miesiąca roku.
10. Dzień tygodnia (day of week): dzień ostatniego kontaktu tygodnia.
11. Czas trwania (duration): czas trwania ostatniego kontaktu w sekundach. Jeśli czas trwania wynosi 0, nigdy nie skontaktowaliśmy się z klientem, aby założyć konto lokaty terminowej.
12. Kampania (campaign): liczba kontaktów wykonanych podczas tej kampanii i dla tego klienta
13. Liczba dni (pdays): liczba dni, które upłynęły od ostatniego kontaktu klienta z poprzedniej kampanii (warto liczbowa; 999 oznacza, że klient wcześniej się nie skontaktował)
14. Poprzedni (previous): liczba kontaktów wykonanych przed tą kampanią i dla tego klienta (numerycznie)
15. Outcome: wynik poprzedniej kampanii marketingowej (kategorycznie: porażka, nieistniejąca, sukces)

III: Atrybuty kontekstu społecznego i gospodarczego:

16. Emp.var.rate: wskaźnik zmienności zatrudnienia - wskaźnik kwartalny
17. Cons.price.idx: wskaźnik cen konsumpcyjnych - wskaźnik miesięczny
18. Cons.conf.idx: wskaźnik zaufania konsumentów - wskaźnik miesięczny
19. Euribor3m: stawka 3-miesięczna euribor - wskaźnik dzienny
20. Liczba zatrudnionych (nr employed): liczba pracowników - wskaźnik kwartalny

Zmienna wyjściowa (podany cel):

21. y - czy klient subskrybował lokatę? (dwukowy: tak, nie)

1.3 Wstępna eksploracja danych

Badane dane zawierają 4119 wierszy oraz 21 kolumn o następujących nazwach:

```
## [1] "age"          "job"          "marital"      "education"
## [5] "default"      "housing"      "loan"        "contact"
## [9] "month"       "day_of_week"  "duration"    "campaign"
## [13] "pdays"      "previous"     "poutcome"    "emp.var.rate"
## [17] "cons.price.idx" "cons.conf.idx" "euribor3m"   "nr.employed"
## [21] "y"
```

Struktura danych:

```
str(df_bank)

## 'data.frame': 4119 obs. of  21 variables:
## $ age          : int  30 39 25 38 47 32 32 41 31 35 ...
## $ job          : chr  "blue-collar" "services" "services" "services" ...
## $ marital      : chr  "married" "single" "married" "married" ...
## $ education    : chr  "basic.9y" "high.school" "high.school" "basic.9y" ...
## $ default      : chr  "no" "no" "no" "no" ...
## $ housing      : chr  "yes" "no" "yes" "unknown" ...
## $ loan         : chr  "no" "no" "no" "unknown" ...
## $ contact      : chr  "cellular" "telephone" "telephone" "telephone" ...
## $ month        : chr  "may" "may" "jun" "jun" ...
## $ day_of_week  : chr  "fri" "fri" "wed" "fri" ...
## $ duration     : int  487 346 227 17 58 128 290 44 68 170 ...
## $ campaign     : int  2 4 1 3 1 3 4 2 1 1 ...
## $ pdays       : int  999 999 999 999 999 999 999 999 999 999 ...
## $ previous     : int  0 0 0 0 0 2 0 0 1 0 ...
## $ poutcome     : chr  "nonexistent" "nonexistent" "nonexistent" "nonexistent" ...
## $ emp.var.rate : num  -1.8 1.1 1.4 1.4 -0.1 -1.1 -1.1 -0.1 -0.1 1.1 ...
## $ cons.price.idx: num  92.9 94 94.5 94.5 93.2 ...
## $ cons.conf.idx: num  -46.2 -36.4 -41.8 -41.8 -42 -37.5 -37.5 -42 -42 -36.4 ...
## $ euribor3m    : num  1.31 4.86 4.96 4.96 4.19 ...
## $ nr.employed  : num  5099 5191 5228 5228 5196 ...
## $ y           : chr  "no" "no" "no" "no" ...
```

Czy w danych znajdują się wartości typu NaN lub Na ?

```
## [1] FALSE
```

Jednakże wiemy, że w danych występują wartości brakujące i są one opisane "unknown". W danych znajduje się 30 rekordów o wartości "unknown" rozmieszczonych w 1029 różnych wierszach. To stanowi 24.98% wszystkich wierszy w naszej bazie danych, więc możemy pozwolić na usunięcie tych wszystkich informacji. W tabeli 1 znajdują się informacje na temat liczby nieznanymi wartości w każdej z kolumn z osobna.

```
## Error: nie znaleziono obiektu 'Number_of_unknown'
## Error in eval(expr, envir, enclos): nie znaleziono obiektu 'table_unknown'
```

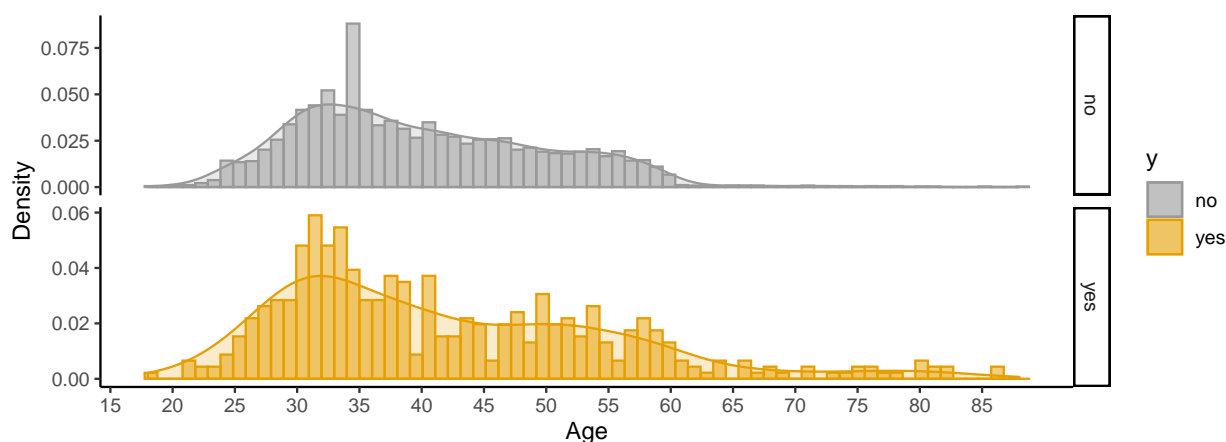
2 Analiza eksploracyjna

W tej sekcji zostanie omnany każdy parametr z osobna. Następnie dane zostaną odpowiednio przygotowane do wykorzystania ich w modelach predykcyjnych.

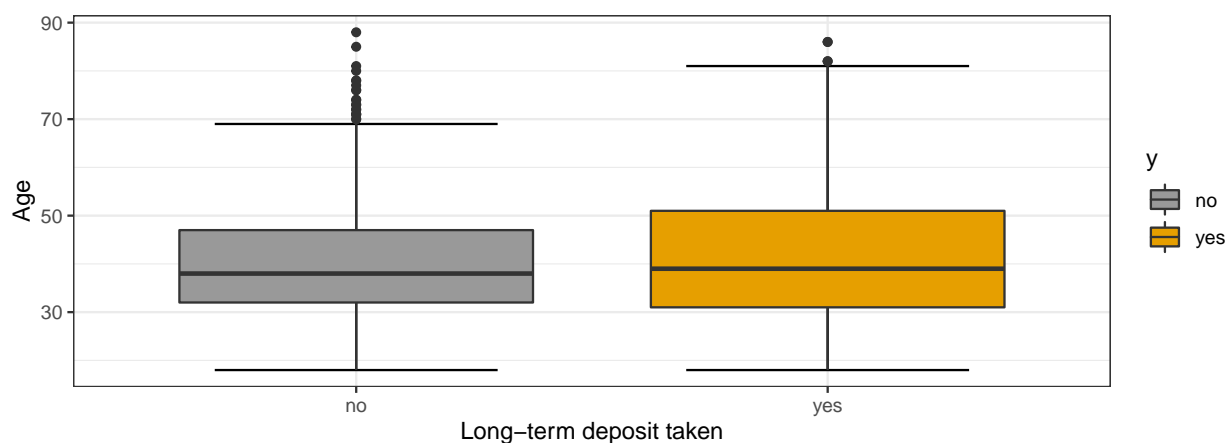
2.1 Age

Przedział wiekowy oszczędzających kredyt szacuje się na 18-88 lat. Jednakże można zauważyć osoby około 60 roku życia, które nie robią lokat, natomiast te, które robią, są w wieku 30-65 lat. Średni wiek utrzymuje się na poziomie 40 lat. Wiadomo, że osoby odkładają na lokaty fundusze wtedy, kiedy dobrze zaczynają zarabian. Podzieliłbym ludzi ze względu na wiek. Młodym [MIN, 30] - młodzi, [30, 65] - pracownicy, [65, MAX] - emeryci. Taki podział powinien ułatwić algorytm

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	18.00	32.00	38.00	40.11	47.00	88.00



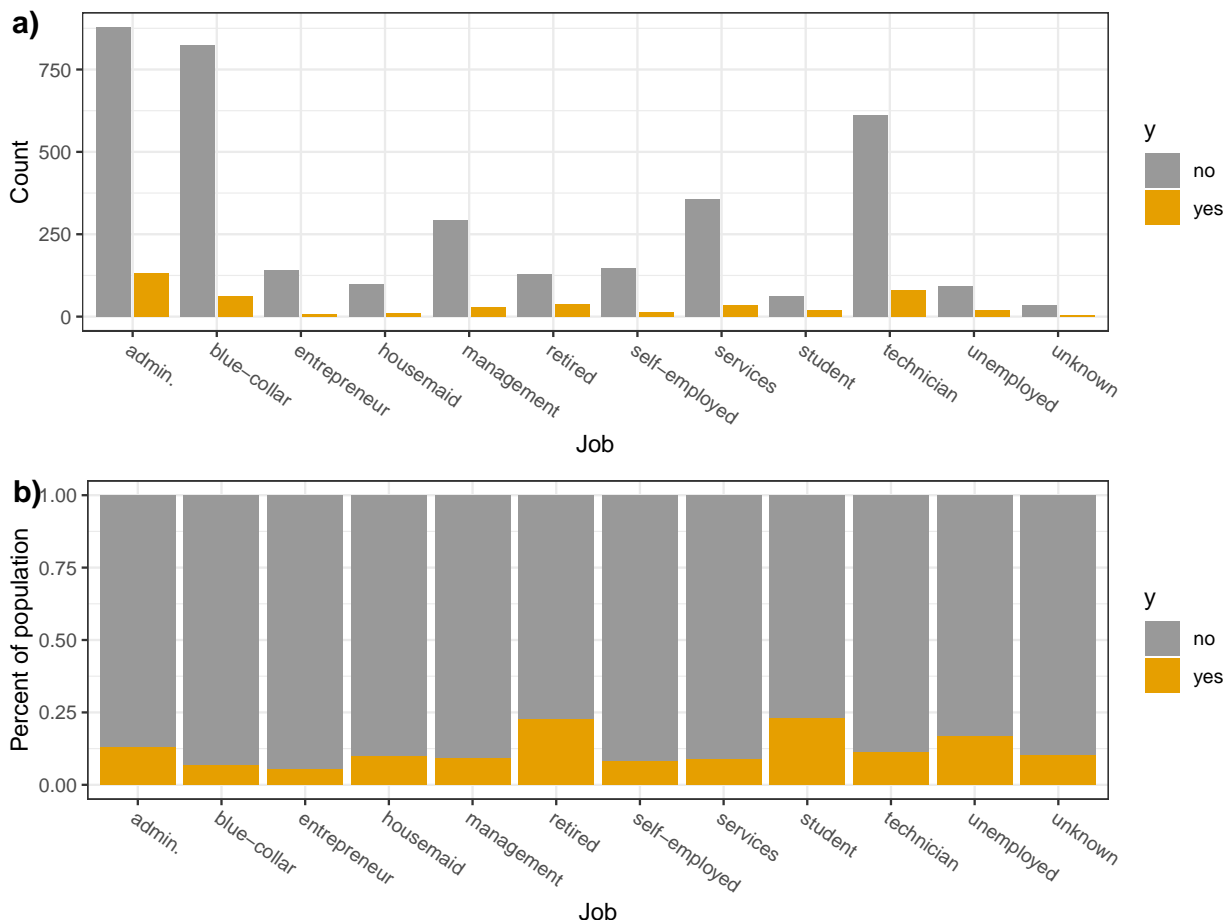
Rysunek 1: Histogram wieku klientów w zależności od wzięcia lokaty długoterminowej.



Rysunek 2: Boxplot wieku klientów w zależności od wzięcia lokaty długoterminowej.

2.2 Job

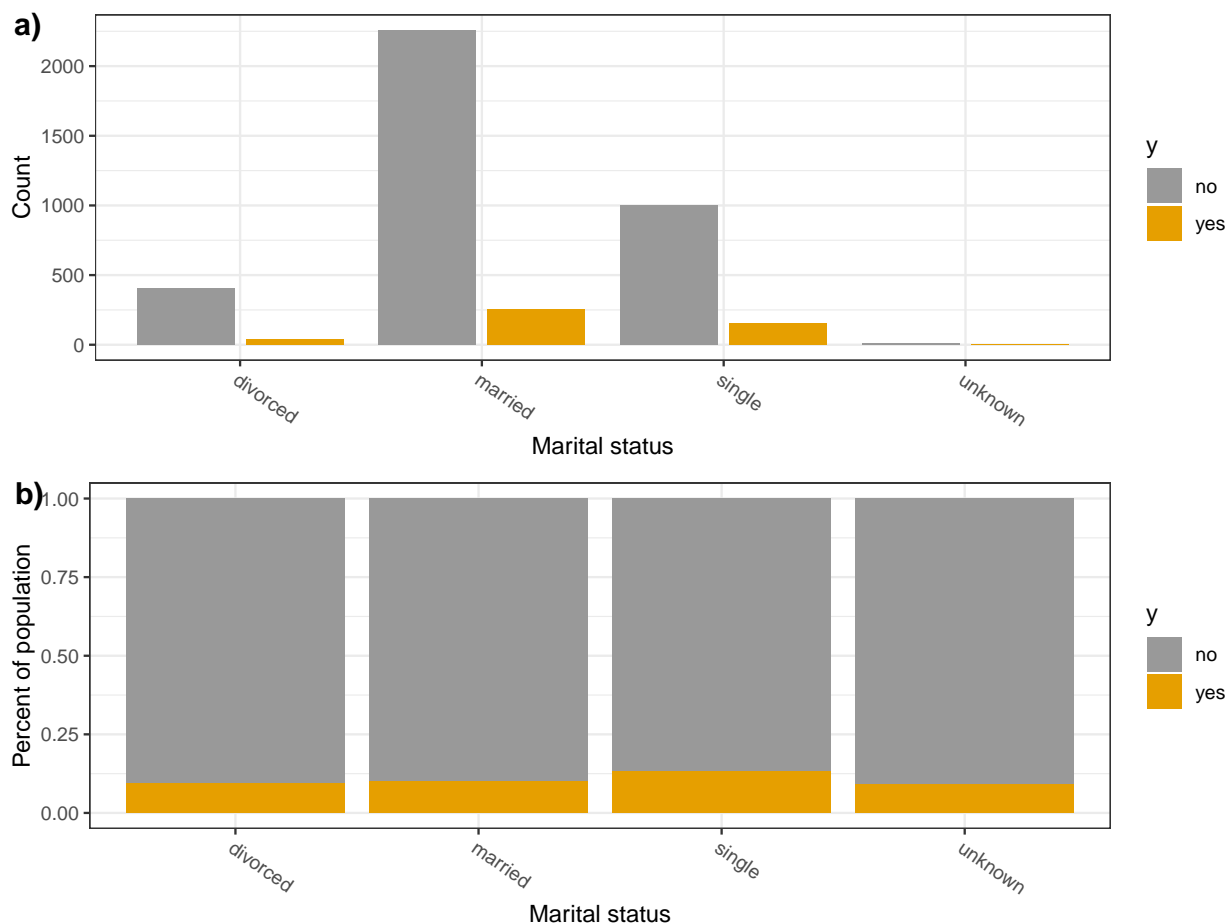
W tej kolumnie mamy 39 wartoci nieznanych, co stanowi ledwo 1% caego zbioru, wiozbywamy sierszy, ktawiera] tformacj



Rysunek 3: Barplot typu a) przedstawiający jak wiele osłabienia za osłabieniem b) przedstawiający stosunek procentowy osłabienia b), kt. reza osłabieniem y lokalizacji z osłabieniem no ci od pracy.

2.3 Marital status

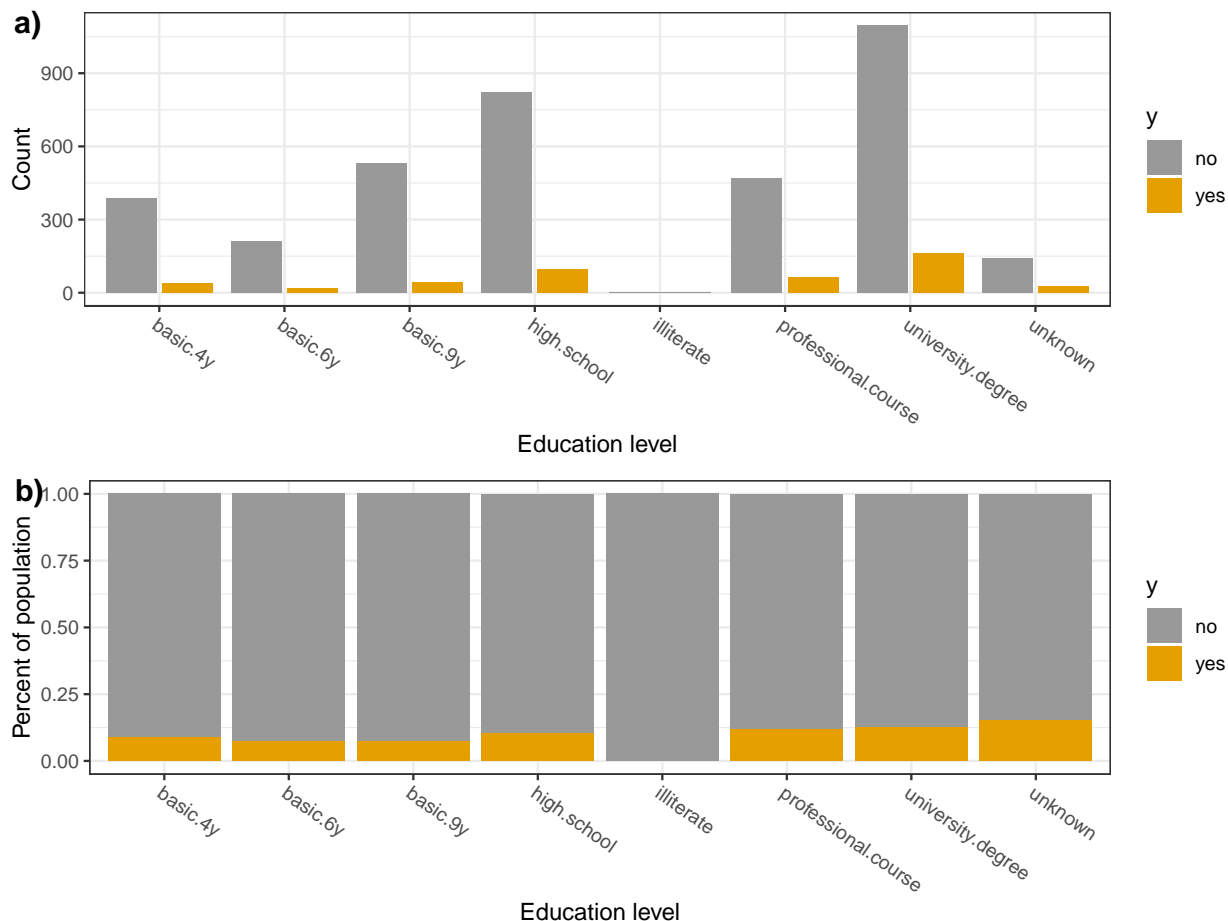
Sytuacja taka sama jak przy kolumnie 'job'. Mamy tutaj nieznane wartoci, ale stanowi one tylko 0.3% wszystkich danych, wie usuwamy te wiersze.



Rysunek 4: Barplot typu a) przedstawiający, jak wiele osób w danej kategorii stanu cywilnego ma odpowiedź 'no' na pytanie o posiadanie dzieci; b) przedstawiający stosunek procentowy osób z odpowiedziami 'no' i 'tak' do stanu cywilnego.

2.4 Education

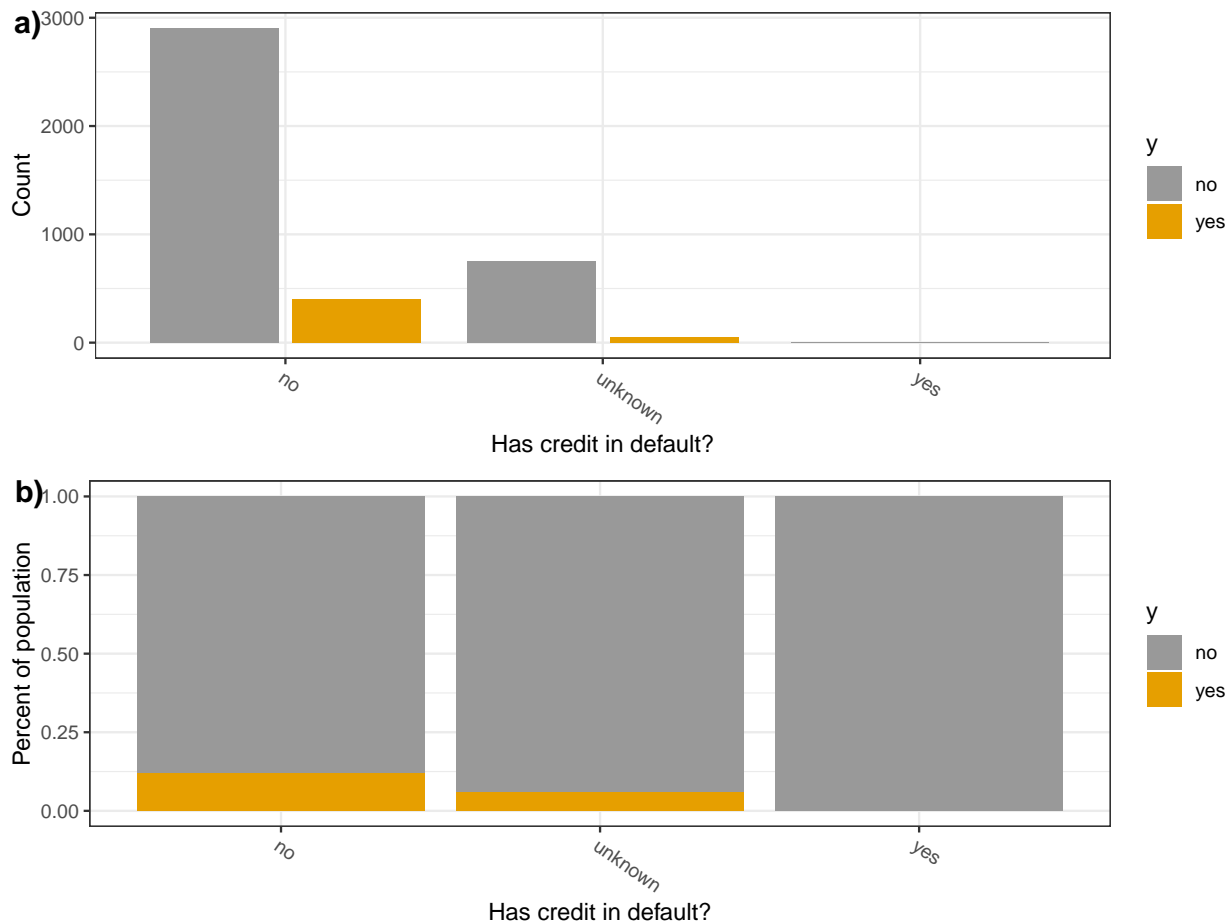
Pula klientów zmienną 'illiterate' zawiera tylko jedna obserwacja, nie ma więc sensu dodawać jej do tej puli klientów. Natomiast w tym przypadku mamy problem z nieznanymi wartościami. Po pierwsze stanowią one 4.1% wszystkich badanych. Najbardziej podobne proporcje danych między 'yes' i 'no' ma kategoria klientów ukończyła uniwersytet, w wszystkich klientach znaną wartość dodamy do tej puli klientów.



Rysunek 5: Barplot typu a) przedstawia, jak wiele osób zależy od poziomu wykształcenia; b) przedstawia stosunek procentowy osób zależących od poziomu wykształcenia.

2.5 Has credit in default?

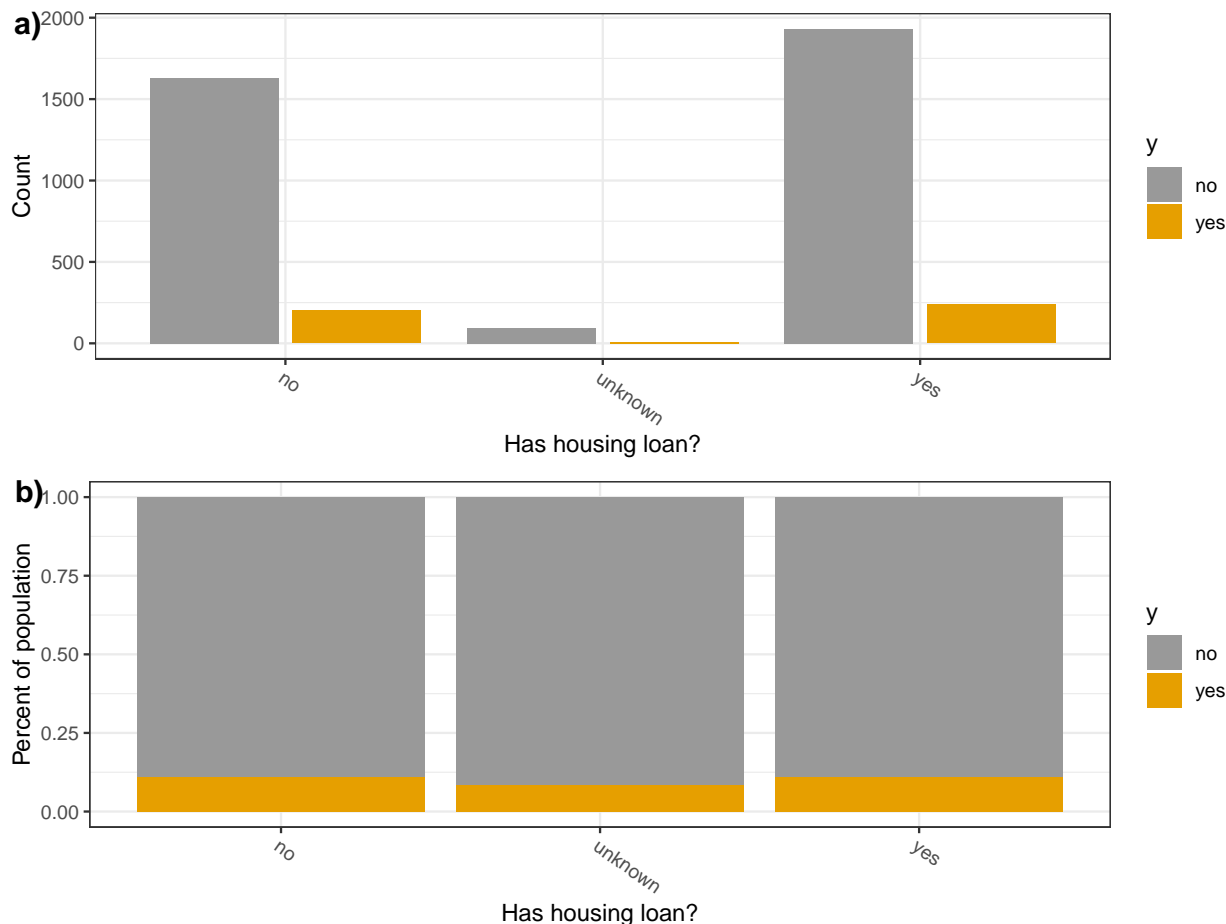
Ta kolumna nie przedstawia wystarczającej ilości danych o osobach, które wzięły ten kredyt. Z tego powodu ta kolumna nie będzie miała żadnego wpływu na nasze modele, dlatego ją usuwamy.



Rysunek 6: Barplot typu a) przedstawia, jak wiele osób założyło kredyt; b) przedstawia stosunek procentowy osób, które złożyły kredyt.

2.6 Has personal loan?

W tej kolumnie znajduj si informacje na temat posiadania kredytu hipotecznego (kredytu na dom). Ilonych nieznanych odpowiada, 2,5% wszystkich obserwacji, Nie moemy pozwolobie na usunie tak duiej liczny wierszy, a podaczenie do jakie innej opcji nie wchodzi w grz przeprowadmy testy na niezalenoienych kategoriycznych. Wykonamy test chisq w zelu zbadania niezalenoci miy 2 zmiennymi.



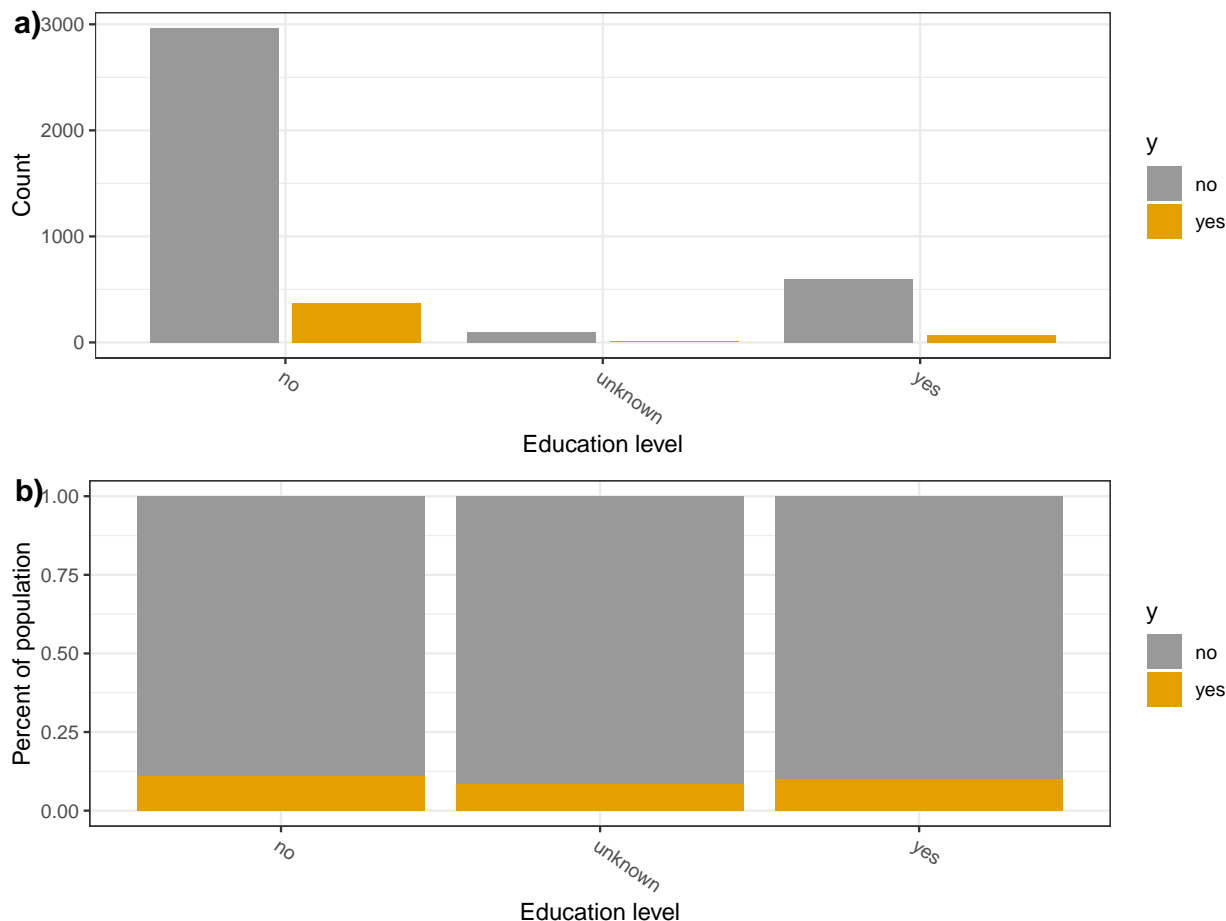
Rysunek 7: Barplot typu a) przedstawiajcy jak wiele os<U+663C><U+3E33>b za<U+623C><U+3E33>o<U+623C><U+3E66>y<U+623C><U+3E33>o lokat<U+653C><U+3E61> w zale<U+623C><U+3E66>no<U+393C><U+3E63>ci od ju<U+623C><U+3E66> posiadanej po<U+623C><U+3E66>yczki; b) przedstawiajcy stosunek procentowy os<U+663C><U+3E33>b, kt<U+663C><U+3E33>re za<U+623C><U+3E33>o<U+623C><U+3E66>y<U+623C><U+3E33>y lokat<U+653C><U+3E61> z zale<U+623C><U+3E66>no<U+393C><U+3E63>ci od posiadanej po<U+623C><U+3E66>yczki.

```
##
## Pearson's Chi-squared test
##
## data: df_bank$housing and df_bank$y
## X-squared = 0.62865, df = 2, p-value = 0.7303
```

Niestety poziom istotnoci(p-value) na poziomie 73% wiadczy o duej zalenoci miy danymi, wiej zmiennej re nie b bra pod uwag

2.7 Has credit in default?

W tej kolumnie znajduj siformacje na temat posiadania kredytu. Sytuacja jest ta sama co w przypadku kredytu hipotecznego. Nie moemy pozwolbie na usunie tak duzej liczny wierszy (2.5%), a podaczenie do jakie innej opcji nie wchodzi w grykonamy test chisq w zelu zbadania niezalenoci miy 2 zmiennymi.



Rysunek 8: Barplot typu a) przedstawiaj<U+623C><U+3E39>cy jak wiele os<U+663C><U+3E33>b za<U+623C><U+3E33>o<U+623C><U+3E66>y<U+623C><U+3E33>o lokat<U+653C><U+3E61> w zale<U+623C><U+3E66>no<U+393C><U+3E63>ci od domy<U+393C><U+3E63>lnie posiadanego kredytu; b) przedstawiaj<U+623C><U+3E39>cy stosunek procentowy os<U+663C><U+3E33>b, kt<U+663C><U+3E33>re za<U+623C><U+3E33>o<U+623C><U+3E66>y<U+623C><U+3E33>y lo- kat<U+653C><U+3E61> z zale<U+623C><U+3E66>no<U+393C><U+3E63>ci od domy<U+393C><U+3E63>lnie posiadanego kredytu.

```
##
## Pearson's Chi-squared test
##
## data: df_bank$loan and df_bank$y
## X-squared = 1.123, df = 2, p-value = 0.5703
```

Poziom istotnoci(p-value) na poziomie 56.8% wiadczy o duzej zalenoci miy danymi, wiej zmiennej re nie b bra pod uwag