

Bank Marketing data (with social/economic context)

Maciej Maecki

25 października 2019

Streszczenie

W pliku Bank Marketing data.csv znajdują się dane charakteryzujące klientów pewnego banku oraz kampanie marketingowe skierowane do tych klientów. Dołączone są ponadto wskaźniki społeczne i ekonomiczne. Na podstawie tych danych należy zbudować model prognozujący szansę, że klient w wyniku prowadzonej kampanii zaakceptuje ofertę lokat terminowych.

Spis treści

1	Wprowadzenie	2
1.1	Opis problemu	2
1.2	Opis danych	2
1.3	Wstępna eksploracja danych	2
2	Analiza eksploracyjna	4
2.1	Age	4
2.2	Job	5
2.3	Marital status	6
2.4	Education	7
2.5	Has credit in default?	8
2.6	Has personal loan?	9
2.7	Has credit in default?	10
2.8	Has credit in default?	11
2.9	Contact communication type	12
2.10	Last contact month of year	13
2.11	Last contact day of the week	14
2.12	Last contact duration, in seconds	14
2.13	15
2.14	16
2.15	17
2.16	18
2.17	19
2.18	Contact communication type?	20

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 otworzyć 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'
```

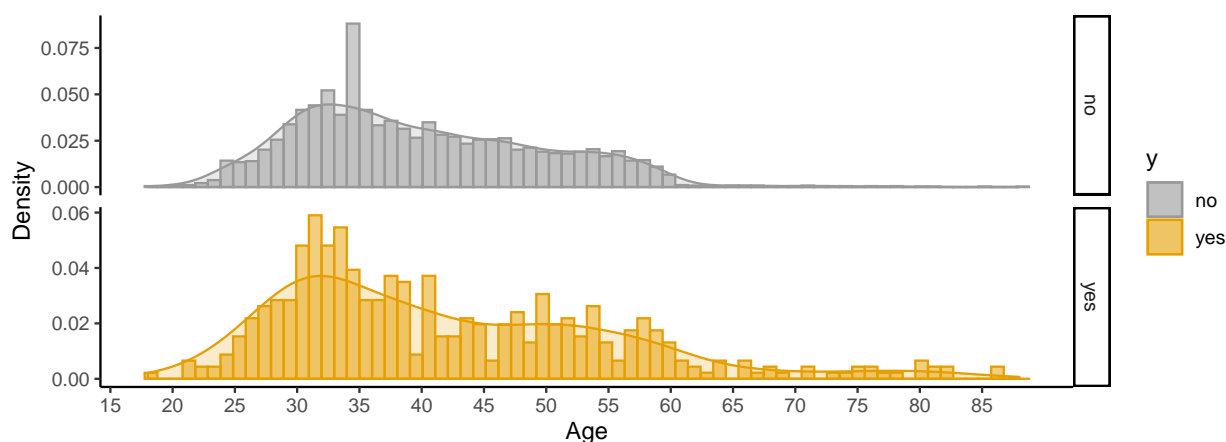
2 Analiza eksploracyjna

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

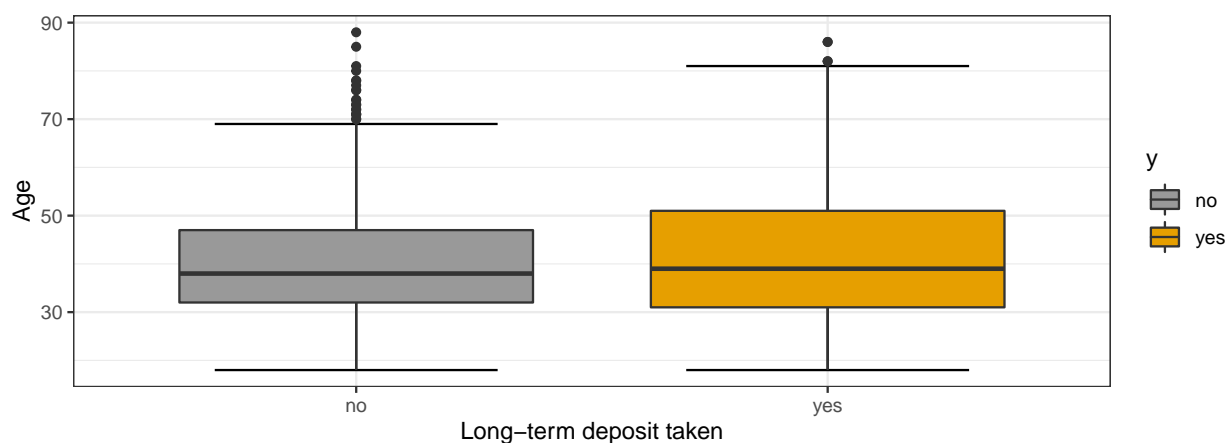
2.1 Age

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. Młodymi [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



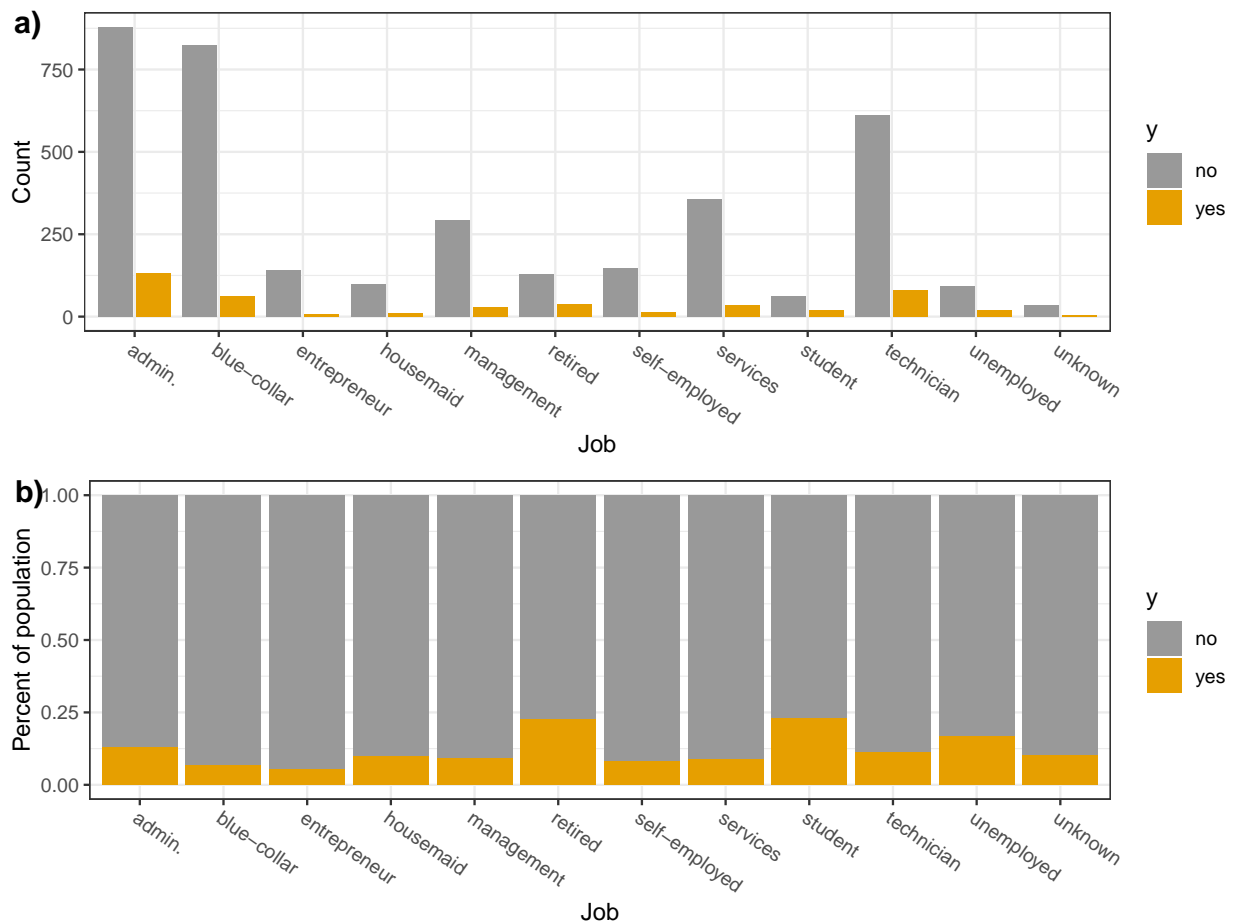
Rysunek 1: Histogram wieku klientów w zależności od wzięcia lokaty długoterminowej.



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

2.2 Job

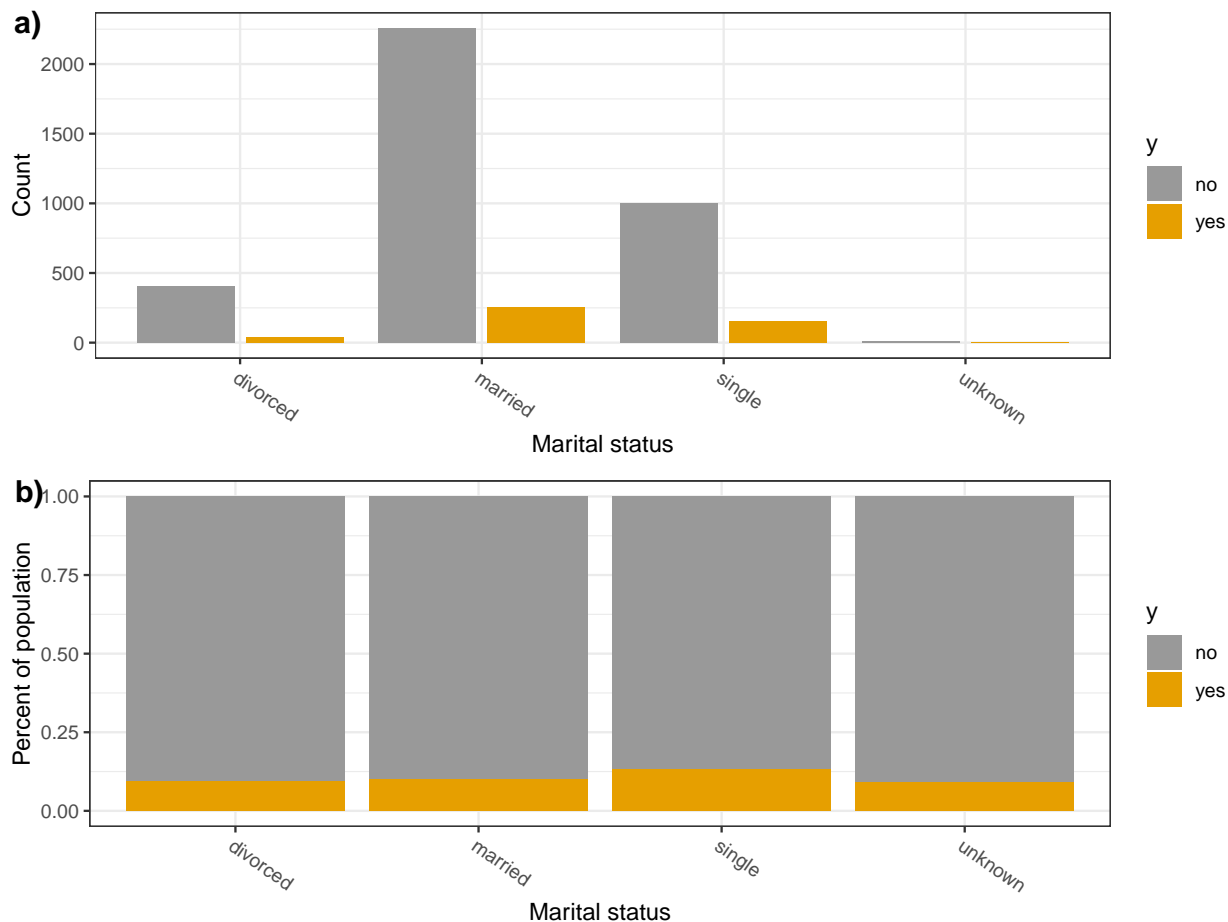
W tej kolumnie mamy 39 wartości nieznanych, co stanowi ledwo 1% całego zbioru, wiozbywamy sierszy, ktawieraj tformacj



Rysunek 3: Barplot typu a) przedstawia, jak wiele osób za daną pracę uważa, że jest to praca, b) przedstawia, jaki stosunek procentowy osób uważających, że praca jest, do osób uważających, że praca nie jest.

2.3 Marital status

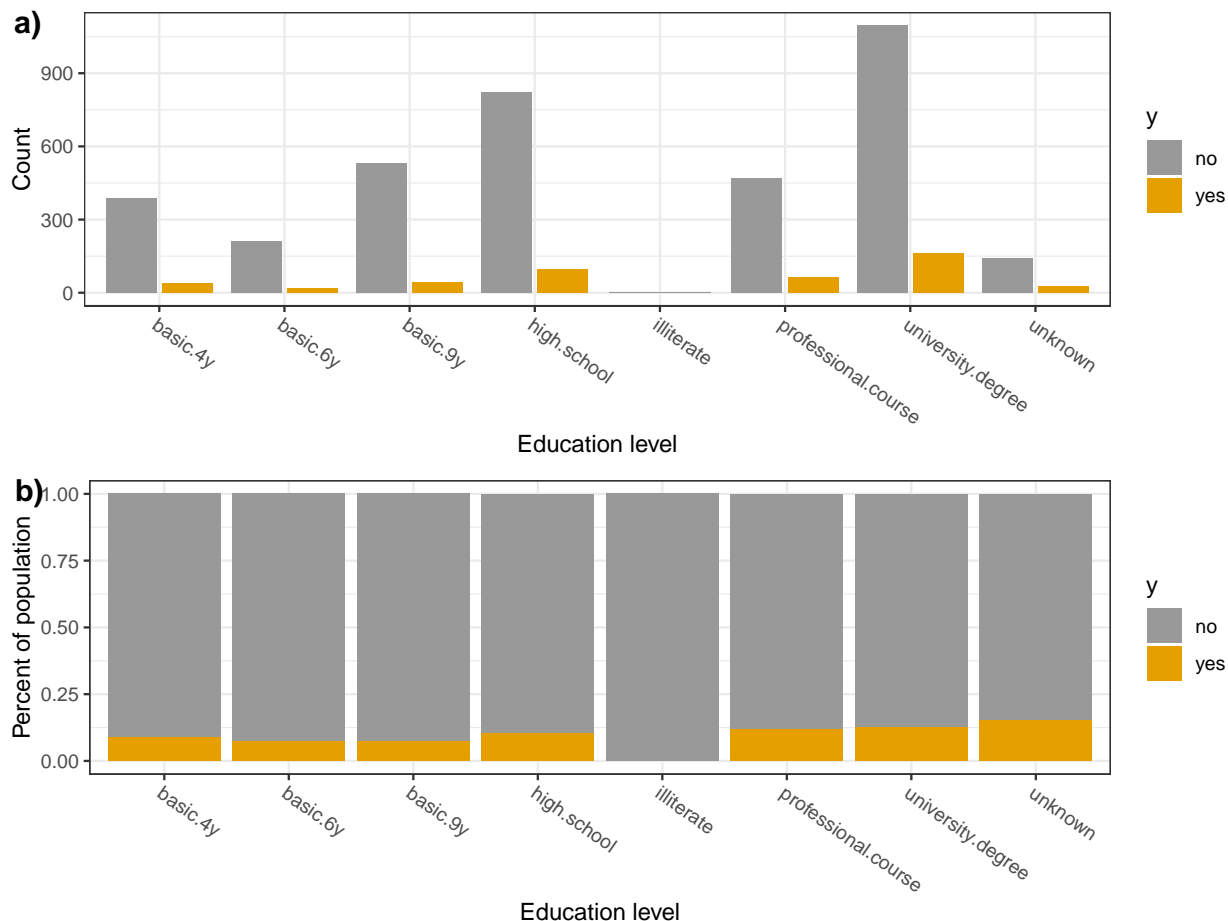
Sytuacja taka sama jak przy kolumnie 'job'. Mamy tutaj nieznane wartoci, ale stanowi one tylko 0.3% wszystkich danych, wie usuwamy te wiersze.



Rysunek 4: Barplot typu a) przedstawiaj $\langle U+623C \rangle \langle U+3E39 \rangle$ cy jak wiele os $\langle U+663C \rangle \langle U+3E33 \rangle$ b za $\langle U+623C \rangle \langle U+3E33 \rangle$ o $\langle U+623C \rangle \langle U+3E66 \rangle$ y $\langle U+623C \rangle \langle U+3E33 \rangle$ o lokat $\langle U+653C \rangle \langle U+3E61 \rangle$ w zale $\langle U+623C \rangle \langle U+3E66 \rangle$ no $\langle U+393C \rangle \langle U+3E63 \rangle$ ci od stanu cywilnego; b) przedstawiaj $\langle U+623C \rangle \langle U+3E39 \rangle$ cy stosunek procentowy os $\langle U+663C \rangle \langle U+3E33 \rangle$ b, kt $\langle U+663C \rangle \langle U+3E33 \rangle$ re za $\langle U+623C \rangle \langle U+3E33 \rangle$ o $\langle U+623C \rangle \langle U+3E66 \rangle$ y $\langle U+623C \rangle \langle U+3E33 \rangle$ y lokat $\langle U+653C \rangle \langle U+3E61 \rangle$ z zale $\langle U+623C \rangle \langle U+3E66 \rangle$ no $\langle U+393C \rangle \langle U+3E63 \rangle$ ci od stanu cywilnego.

2.4 Education

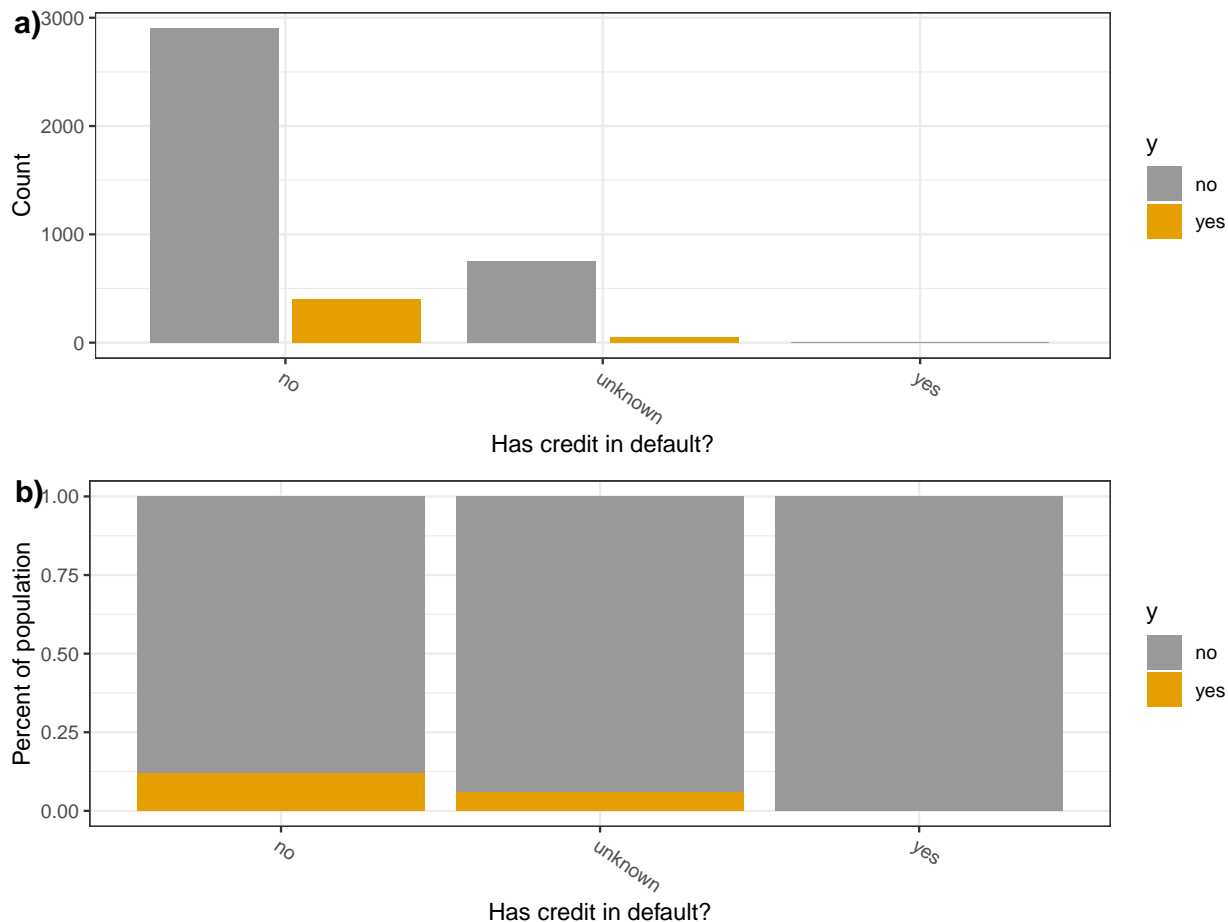
Pula klientów zmienną 'illiterate' zawiera tylko jedna obserwacja, nie ma więc sensu dodawać jej do tej puli klientów. Natomiast w tym przypadku mamy problem z nieznanymi wartościami. Po pierwsze stanowią one 4.1% wszystkich badanych. Najbardziej podobne proporcje danych między 'yes' i 'no' ma kategoria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

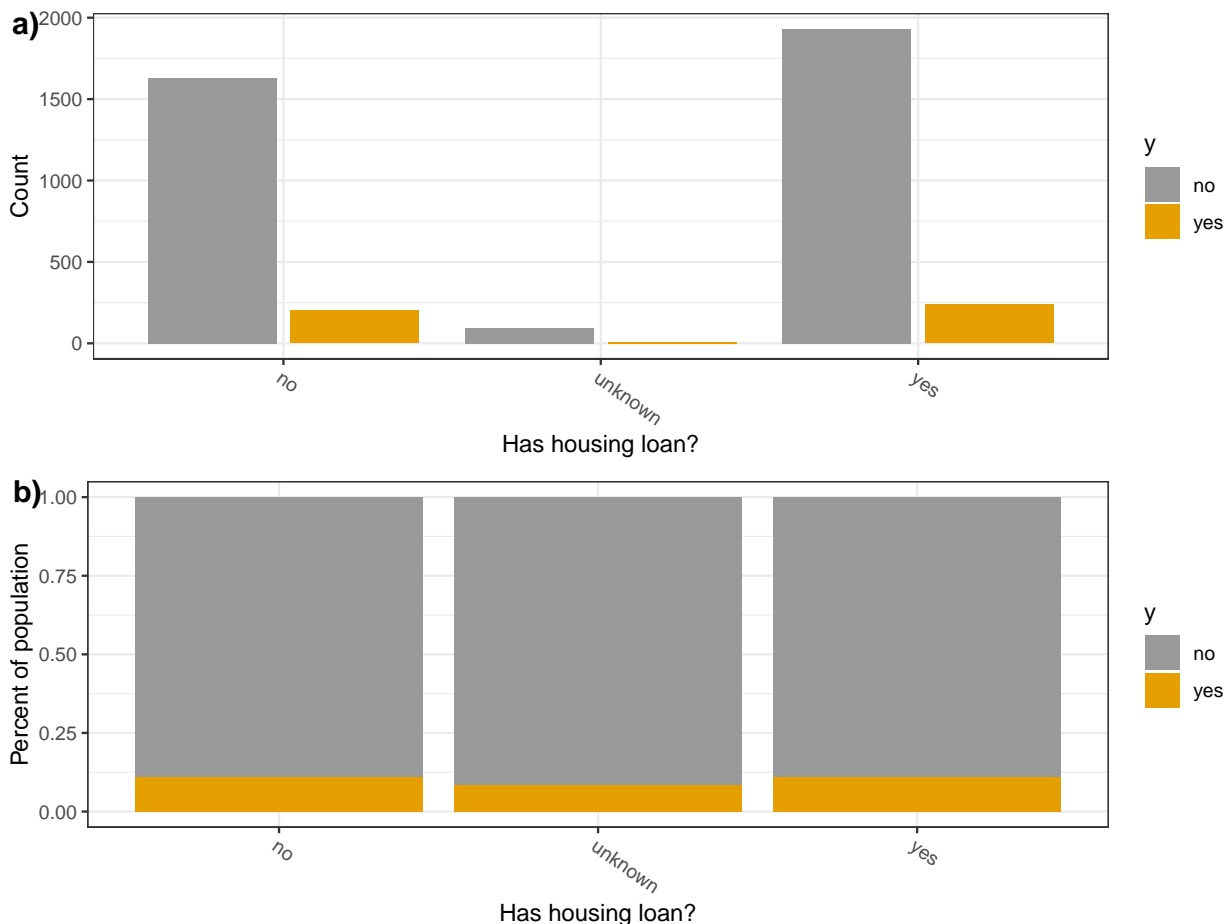
Ta kolumna nie przedstawia wystarczającej ilości danych o osobach, które wzięły ten kredyt. Z tego powodu ta kolumna nie będzie miała żadnego wpływu na nasze modele, dlatego ją usuwamy.



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

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.



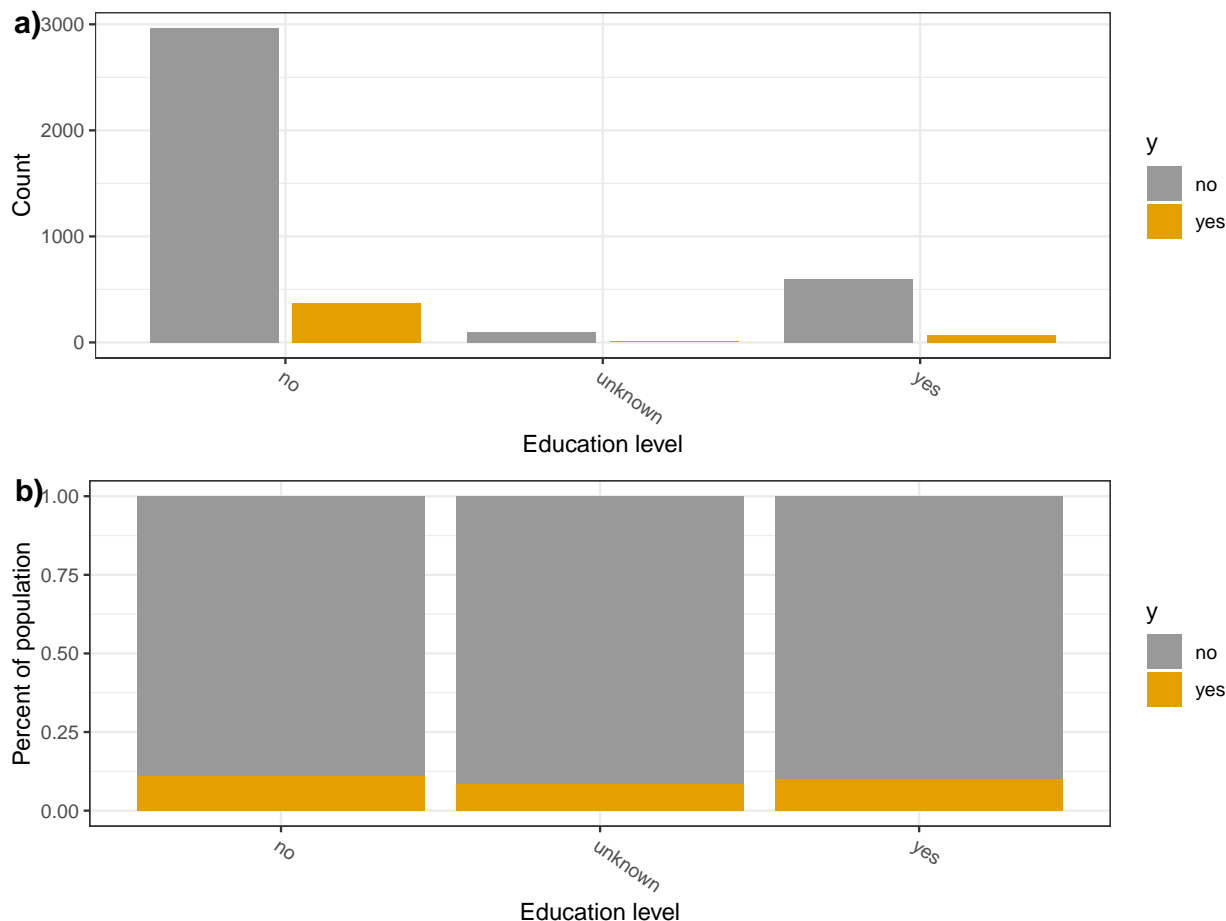
Rysunek 7: Barplot typu a) przedstawiajcy jak wiele os<U+663C><U+3E33>b za<U+623C><U+3E33>o<U+623C><U+3E66>y<U+623C><U+3E33>o lokat<U+653C><U+3E61> w zale<U+623C><U+3E66>no<U+393C><U+3E63>ci od ju<U+623C><U+3E66> posiadanej po<U+623C><U+3E66>yczki; b) przedstawiajcy stosunek procentowy os<U+663C><U+3E33>b,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

2.8 Has credit in default?

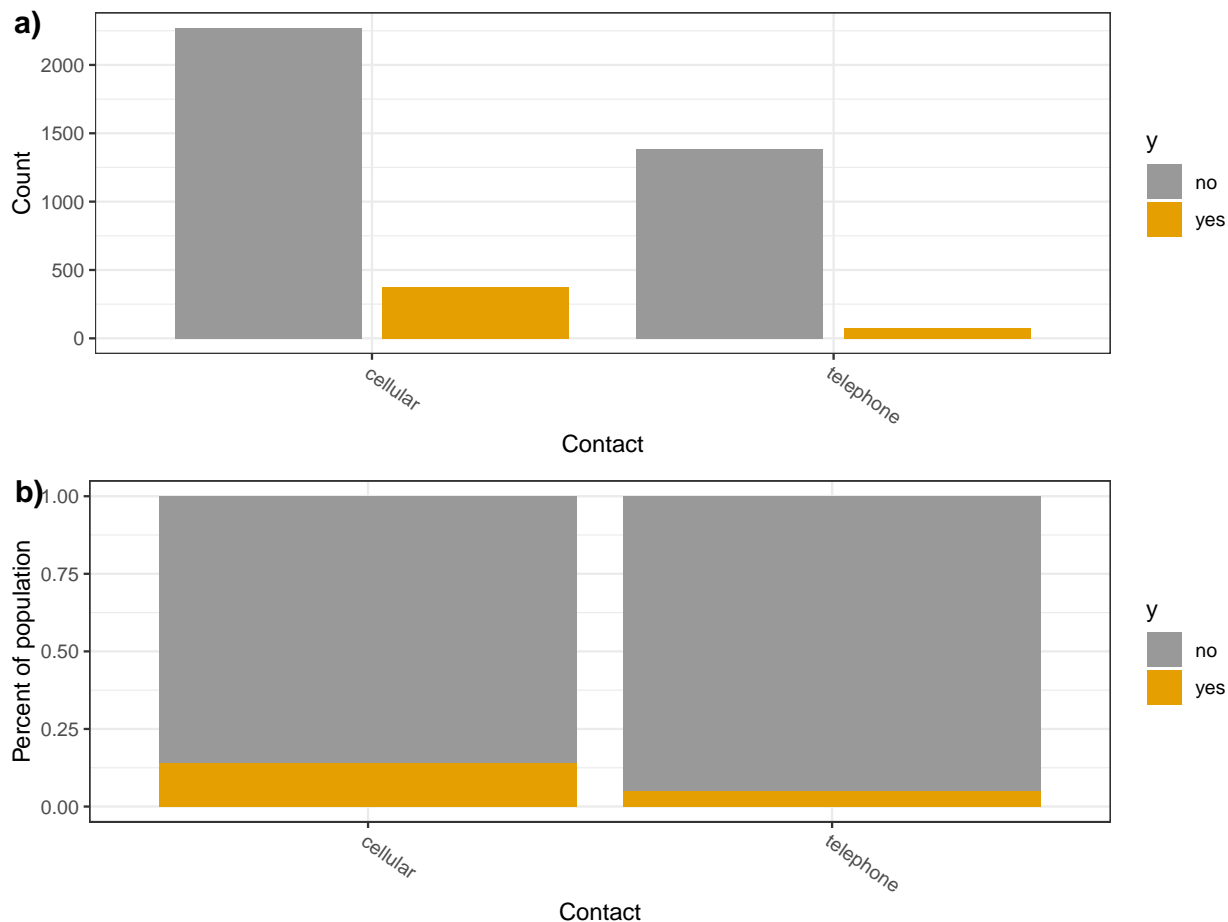
W tej kolumnie znajduj si informacje na temat posiadania kredytu. Sytuacja jest ta sama co w przypadku kredytu hipotecznego. Nie moemy pozwolobie na usunie tak duzej liczny wierszy (2.5%), a podaczenie do jakie innej opcji nie wchodzi w grykonamy test chisq w zelu zbadania niezalenoci miy 2 zmiennymi.

```
## Error in FUN(X[[i]], ...): nie znaleziono obiektu 'loan'
```

```
## Error in chisq.test(df_bank$loan, df_bank$y): 'x' and 'y' must have the same length
```

2.9 Contact communication type

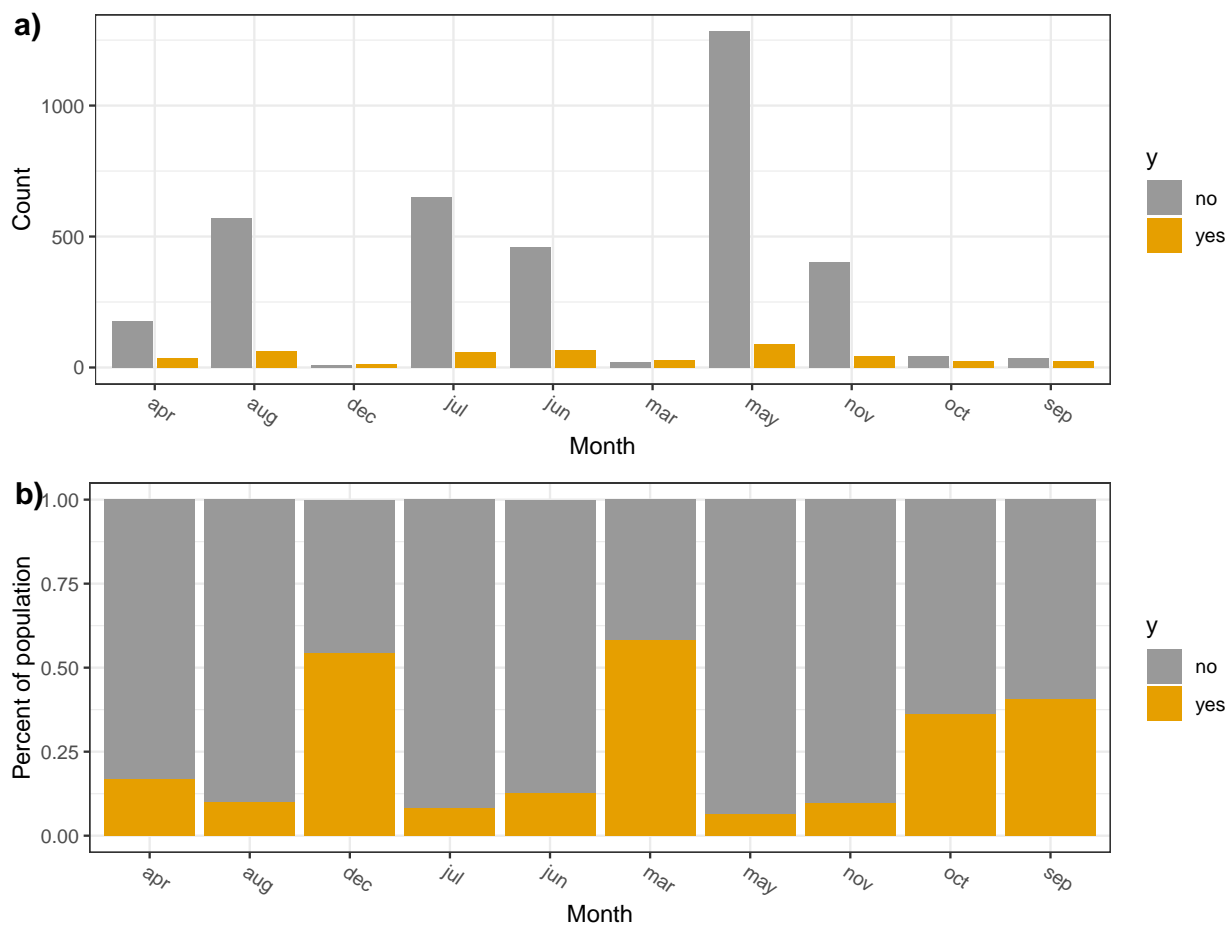
Osoby, z kti prano siontaktowa telefon komwy stanowi 64.4% caej badanej spoeczności i co 6 osoba z nich wziera lokatugoterminow.



Rysunek 9: 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.

2.10 Last contact month of year

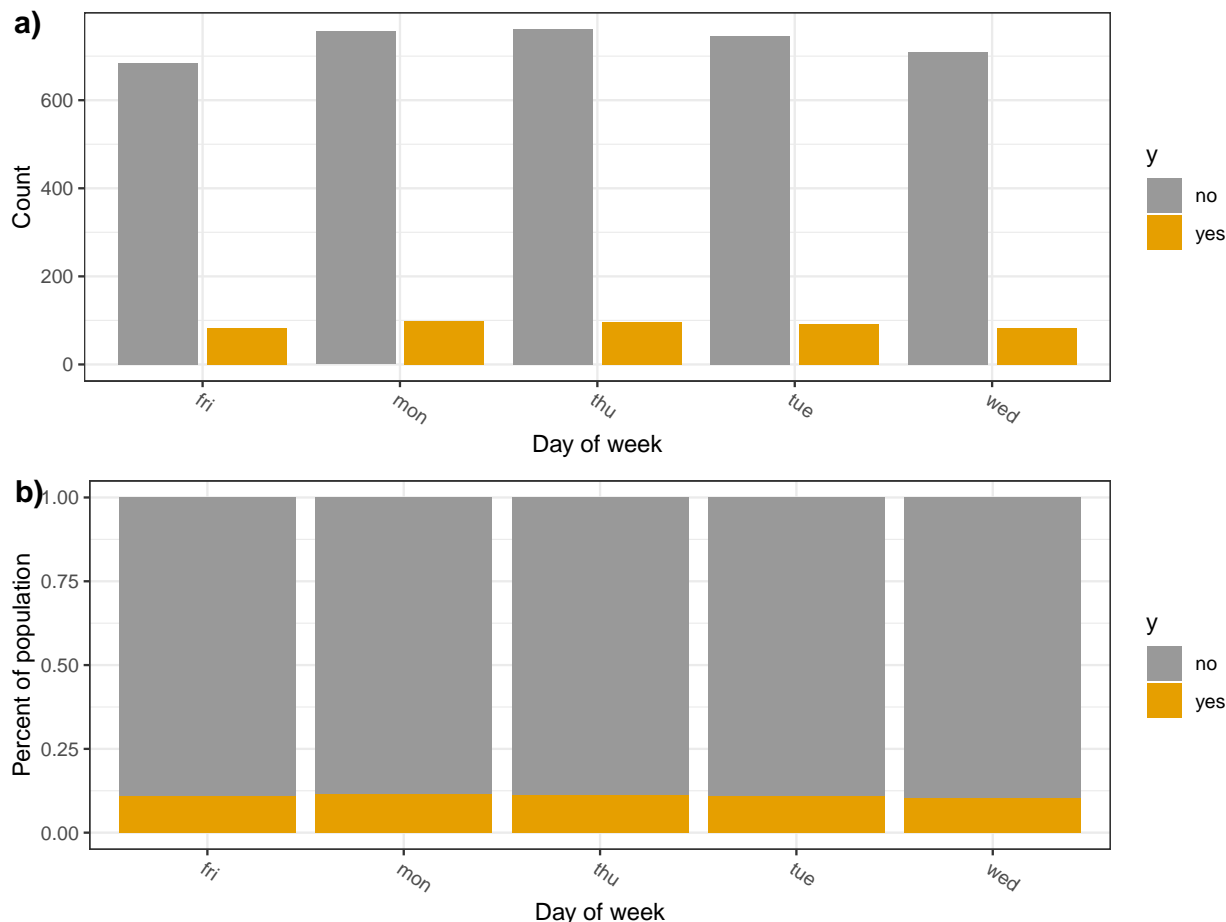
Ciekaw sytuacją jest fakt, że w zestawieniu w ogóle mamy transakcje ze stycznia oraz z lutego.



Rysunek 10: Barplot typu a) przedstawiający, jak wiele osób założyło dom w zależności od ostatniego kontaktu; b) przedstawiający stosunek procentowy osób, które założyły dom w zależności od ostatniego kontaktu.

2.11 Last contact day of the week

W każdym roboczym dniu tygodnia jest wykonywane mniej więcej tyle samo poczeintami, więc jesteśmy w stanie wygenerować wnioski obserwacji samych wykresów.



Rysunek 11: Barplot typu a) przedstawia, jak wiele osób założyło mieszkanie w zależności od dnia tygodnia, w którym zostało wykonane poczeintami; b) przedstawia stosunek procentowy osób, które założyły mieszkanie w zależności od dnia tygodnia, w którym zostało wykonane poczeintami.

2.12 Last contact duration, in seconds

Czas trwania ostatniego kontaktu jest atrybutem, który ma duży wpływ na cel wyjściowy (y). Wskazujemy uwagę: ten atrybut ma duży wpływ na cel wyjściowy. Jednak czas trwania nie jest znany przed wykonaniem poczeintami. Ponadto po wykonaniu poczeintami "y" jest oczywiście znane. W związku z tym należy odrzucić chęć stworzenia realistycznego modelu predykcyjnego.

2.13

2.14

2.15

2.16

2.17

2.18 Contact communication type?