

compiled Project

Seth Billiau

4/13/2020

```
source("na-convert.R")
library(tidyverse)

## -- Attaching packages ----- tidyverse 1.3.0 --
## v ggplot2 3.3.0      v purrr  0.3.3
## v tibble  3.0.0      v dplyr  0.8.5
## v tidyr   1.0.2      v stringr 1.4.0
## v readr   1.3.1      v forcats 0.5.0

## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()

library(ggpubr)

## Loading required package: magrittr
##
## Attaching package: 'magrittr'
##
## The following object is masked from 'package:purrr':
##
##   set_names
##
## The following object is masked from 'package:tidyr':
##
##   extract

library(xtable)
```

EDA

Read in the data and make note of missing values

```
data_raw = read.csv("data/chd_risk.csv")
summary(data_raw)
```

	age	education	cigsPerDay	totChol
## Min.	:32.00	College or Higher : 473	Min. : 0.000	Min. :107.0
## 1st Qu.:	:42.00	High School or GED:1253	1st Qu.: 0.000	1st Qu.:206.0
## Median :	:49.00	Some College : 687	Median : 0.000	Median :234.0
## Mean :	:49.58	Some High School :1720	Mean : 9.003	Mean :236.7
## 3rd Qu.:	:56.00	NA's : 105	3rd Qu.:20.000	3rd Qu.:263.0
## Max.	:70.00		Max. :70.000	Max. :696.0
##			NA's :29	NA's :50

```
##      sysBP      diaBP      BMI      heartRate
## Min.   : 83.5   Min.   : 48.00   Min.   :15.54   Min.   : 44.00
## 1st Qu.:117.0   1st Qu.: 75.00   1st Qu.:23.07   1st Qu.: 68.00
## Median :128.0   Median : 82.00   Median :25.40   Median : 75.00
## Mean   :132.4   Mean   : 82.89   Mean   :25.80   Mean   : 75.88
## 3rd Qu.:144.0   3rd Qu.: 89.88   3rd Qu.:28.04   3rd Qu.: 83.00
## Max.   :295.0   Max.   :142.50   Max.   :56.80   Max.   :143.00
##                                     NA's   :19   NA's   :1
##      glucose      sex      smoker      OnBPMeds      PrevStroke
## Min.   : 40.00   female:2419   Nonsmoker:2144   No :4061   No :4213
## 1st Qu.: 71.00   male :1819   Smoker :2094   Yes : 124   Yes: 25
## Median : 78.00
## Mean   : 81.97
## 3rd Qu.: 87.00
## Max.   :394.00
## NA's   :388
##      Hyp      Diab      CHD_Risk
## No :2922   No :4129   No :3594
## Yes:1316   Yes: 109   Yes: 644
##
##
##
##
##
```

Count number in on missingness:

```
# Generate the number of missing values for each predictor
apply(is.na(data_raw), 2, sum)
```

```
##      age  education  cigsPerDay  totChol  sysBP  diaBP  BMI
##      0      105      29      50      0      0      19
## heartRate  glucose      sex  smoker  OnBPMeds PrevStroke  Hyp
##      1      388      0      0      53      0      0
##      Diab  CHD_Risk
##      0      0
```

```
missing_preds = c("education", "cigsPerDay", "totChol", "BMI",
                  "heartRate", "glucose", "OnBPMeds")
```

Visualize distribution of quantitative predictors conditional on the CHD outcome:

```
quant_preds = c("age", "cigsPerDay", "totChol", "sysBP",
                "diaBP", "BMI", "heartRate", "glucose")
```

```
make_cond_hist = function(varname) {
  p1 = ggplot(data_raw, aes_string(x=varname)) +
    geom_histogram(aes(y = ..density..),
                  fill = "red", alpha = 0.5) +
    labs(title=paste(varname, "given CHD_Risk")) +
    xlab(varname) +
    ylab("Density") +
    facet_grid(. ~ CHD_Risk) +
    theme_bw()
  return(p1)
}
```

```

graphs = lapply(quant_preds, make_cond_hist)
figure1 = ggarrange(graphs[[1]], graphs[[2]], graphs[[3]], graphs[[4]],
  graphs[[5]], graphs[[6]], graphs[[7]], graphs[[8]],
  ncol = 2, nrow = 4)

## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## Warning: Removed 29 rows containing non-finite values (stat_bin).
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## Warning: Removed 50 rows containing non-finite values (stat_bin).
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## Warning: Removed 19 rows containing non-finite values (stat_bin).
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## Warning: Removed 1 rows containing non-finite values (stat_bin).
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## Warning: Removed 388 rows containing non-finite values (stat_bin).
annotate_figure(figure1,
  top = text_grob("Visualizing Quantitative Predictors given CHD_Risk (prevalence = 0.152)"
)

```

Visualizing Quantitative Predictors given CHD_Risk (prevalence = 0.152)

