

# Research log

Dennis Scheper

2020-10-18

## 1 Original dataset

Date: 11 September, 2020 (week 1)

### 1.1 Introduction

The article, “Impact of HbA1c Measurement on Hospital Readmission Rates: Analysis of 70,000 Clinical Database Patient Records” states that hyperglycemia management has a significant impact on the outcome and readmission rates of hospitalized patients. The authors based this conclusion on the comprehensive assessment of 70,000 diabetes patient records retrieved from 140 US hospitals. The results that depict the relationship between readmission rates and the measurement of HbA1c levels can further enlargement of already existing diabetes tactics to reduce readmission rates.

### 1.2 Dataset and Attributes Information

All information used in this article comes from a database, consisting of 41 tables and totals 117 features, such as demographics (like gender, race, and age), inpatient or outpatient, and (in-hospital) mortality. Data came from 130 hospitals in the USA for over ten years (1998-2008) and contained around 74 million unique visits by 18 million unique patients. This research used information that needs to accede to the following specifications:

1. Is a hospital admission;
2. The encounter is a diagnosed with 'diabetes', any kind will satisfy;
3. The length of admission was at least one day up to eighteen days;
4. Laboratory test results are available; and
5. Medications were administered.

Of these five criteria points, 101.000 encounters fulfilled all specifications. After some considerations with removing encounters based on incomplete (weight and medical specialty) or biased data (discharge to a hospice or death), 69.984 encounters remained in the final dataset.

The initial dataset consists of 55 attributes, with the class attribute being an encounter of one patient. As there are way too many attributes to describe, please refer to the codebook for all descriptions; we only look at some essential attributes and their type and possible valuations. The age of a patient is nominal and is grouped into ten-year intervals. The admission type or for what specific reason a patient was hospitalized and comes in 9 distinct values while the type is nominal. Some attributes are numeric and count, such as the number of lab procedures, the number of medications, and the number of emergency visits. The database consists of three diagnosis attributes, which can have 848, 923, and 954 distinct values, respectively. The values are based on ICD9 three-letter codes and are of the nominal type. Other vital attributes are whether a patient changed medications (with the values ‘no change’ and ‘change’; nominal type) or had diabetes medication (‘yes’ or ‘no’ values and nominal typing). Twenty-four other attributes depict whether medicine is prescribed or not, if prescribed, then if the dosage was increased (‘up’), decreased (‘down’), or stayed the same (‘steady’) during the encounter. Readmission rates were calculated by looking at a nominal type (‘Readmitted’) with the possible valuations of ‘<30’ for a patient that was readmitted within 30 days, ‘>30’ for a patient that was readmitted after 30 days, and ‘No’ for patients that were not readmitted. The authors’ goal was to determine whether a relationship between readmission rates and HbA1c measurement exists; therefore, they introduced a new attribute ‘HbA1c’ with four different valuations, based on the information from the database: 1) no HbA1c test performed; 2) HbA1c performed and in the normal range; 3) HbA1c performed and the result is greater than 8% with no change in diabetic medication; and 4) Hb1Ac performed, the result is more significant than 8%, and diabetic medication was changed.

### 1.3 Research Question

Is it possible, using machine learning techniques, to predict a patient’s time in hospital linked to the results of an HbA1c measurement?

## 2 EDA (Exploratory Data Analysis)

### 2.1 Date: 14 September, 2020

This section will perform an exploratory data analysis (EDA) on the dataset described above. With an EDA, we can explore our dataset and determine if any correlations exist between attributes and undermine any missing values. If any exist, we deal with them accordingly, looking to repair these values or remove them entirely. Additionally, we will look at variations between and in attributes for determining which ones are of most importance and interest to our research goals.

First, we load in our used packages, mostly used for visualization of data, and the data itself. Besides that, we also load in a codebook, which contains, i.e., a description for every attribute for the initial dataset. The codebook was not retrieved from any outside source but was constructed on the interpretation of the data.

```
# Load used packages
library(formattable)
library(plyr)
library(dplyr, warn.conflicts = FALSE)
library(tidyr)
library(tidyverse)
library(ggplot2)
library(grid)
library(gridExtra)
library(hexbin)
library(foreign)
library(xtable)

# Load in used data
encounter.data <- read.table(file = "dataset_diabetes/retrieved_data/diabetic_data.csv",
                             sep = ",", header = TRUE)

codebook <- read.csv(file = "codebook.csv", sep = ";", header = T)
```

To get a better picture of our dataset's distribution, we called upon the function `glimpse()`. We notice that the data contains 50 columns/ attributes. Additionally, a `summary()` function gives an outline of every attribute, showing each level with a count of its valuation.

```
# Give a glimpse of how the data is distributed
glimpse(encounter.data)
```

```
## Rows: 101,766
## Columns: 50
```

```

## $ encounter_id <int> 2278392, 149190, 64410, 500364, 16680, 357...
## $ patient_nbr <int> 8222157, 55629189, 86047875, 82442376, 425...
## $ race <fct> Caucasian, Caucasian, AfricanAmerican, Cau...
## $ gender <fct> Female, Female, Female, Male, Male, Male, ...
## $ age <fct> [0-10), [10-20), [20-30), [30-40), [40-50)...
## $ weight <fct> ?, ?, ?, ?, ?, ?, ?, ?, ?, ?, ?, ?, ?, ?, ...
## $ admission_type_id <int> 6, 1, 1, 1, 1, 2, 3, 1, 2, 3, 1, 2, 1, 1, ...
## $ discharge_disposition_id <int> 25, 1, 1, 1, 1, 1, 1, 1, 3, 1, 1, 3, 6, ...
## $ admission_source_id <int> 1, 7, 7, 7, 7, 2, 2, 7, 4, 4, 7, 4, 7, 7, ...
## $ time_in_hospital <int> 1, 3, 2, 2, 1, 3, 4, 5, 13, 12, 9, 7, 7, 1...
## $ payer_code <fct> ?, ?, ?, ?, ?, ?, ?, ?, ?, ?, ?, ?, ?, ?, ...
## $ medical_specialty <fct> Pediatrics-Endocrinology, ?, ?, ?, ?, ?, ?, ...
## $ num_lab_procedures <int> 41, 59, 11, 44, 51, 31, 70, 73, 68, 33, 47...
## $ num_procedures <int> 0, 0, 5, 1, 0, 6, 1, 0, 2, 3, 2, 0, 0, 1, ...
## $ num_medications <int> 1, 18, 13, 16, 8, 16, 21, 12, 28, 18, 17, ...
## $ number_outpatient <int> 0, 0, 2, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ...
## $ number_emergency <int> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, ...
## $ number_inpatient <int> 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ...
## $ diag_1 <fct> 250.83, 276, 648, 8, 197, 414, 414, 428, 3...
## $ diag_2 <fct> ?, 250.01, 250, 250.43, 157, 411, 411, 492...
## $ diag_3 <fct> ?, 255, V27, 403, 250, 250, V45, 250, 38, ...
## $ number_diagnoses <int> 1, 9, 6, 7, 5, 9, 7, 8, 8, 8, 9, 7, 8, 8, ...
## $ max_glu_serum <fct> None, None, None, None, None, None, ...
## $ A1Cresult <fct> None, None, None, None, None, None, None, ...
## $ metformin <fct> No, No, No, No, No, Steady, No, No, No...
## $ repaglinide <fct> No, No, No, No, No, No, No, No, No, ...
## $ nateglinide <fct> No, No, No, No, No, No, No, No, No, ...
## $ chlorpropamide <fct> No, No, No, No, No, No, No, No, No, ...
## $ glimepiride <fct> No, No, No, No, No, No, Steady, No, No, ...
## $ acetohexamide <fct> No, No, No, No, No, No, No, No, No, ...
## $ glipizide <fct> No, No, Steady, No, Steady, No, No, No, St...
## $ glyburide <fct> No, No, No, No, No, No, Steady, No, No, ...
## $ tolbutamide <fct> No, No, No, No, No, No, No, No, No, ...
## $ pioglitazone <fct> No, No, No, No, No, No, No, No, No, ...
## $ rosiglitazone <fct> No, No, No, No, No, No, No, No, No, ...
## $ acarbose <fct> No, No, No, No, No, No, No, No, No, ...
## $ miglitol <fct> No, No, No, No, No, No, No, No, No, ...
## $ troglitazone <fct> No, No, No, No, No, No, No, No, No, ...
## $ tolazamide <fct> No, No, No, No, No, No, No, No, No, ...
## $ examide <fct> No, No, No, No, No, No, No, No, No, ...
## $ citoglipton <fct> No, No, No, No, No, No, No, No, No, ...
## $ insulin <fct> No, Up, No, Up, Steady, Steady, Steady, No...
## $ glyburide.metformin <fct> No, No, No, No, No, No, No, No, ...

```

```

## $ glipizide.metformin      <fct> No, No, No, No, No, No, No, No, No, No...
## $ glimepiride.pioglitazone <fct> No, No, No, No, No, No, No, No, No, No...
## $ metformin.rosiglitazone   <fct> No, No, No, No, No, No, No, No, No, No...
## $ metformin.pioglitazone   <fct> No, No, No, No, No, No, No, No, No, No...
## $ change                   <fct> No, Ch, No, Ch, Ch, No, Ch, No, Ch, Ch, No...
## $ diabetesMed              <fct> No, Yes, Yes, Yes, Yes, Yes, Yes, Yes, Yes...
## $ readmitted               <fct> NO, >30, NO, NO, NO, >30, NO, >30, NO, NO, ...

```

*# Give a quick summary of the data, and give information such as min and max  
# value, median and mean*

```
summary(encounter.data)
```

```

##   encounter_id      patient_nbr          race
##   Min.    : 12522   Min.    : 135     ?: 2273
##   1st Qu.: 84961194 1st Qu.: 23413221 AfricanAmerican:19210
##   Median  :152388987 Median  : 45505143 Asian    : 641
##   Mean    :165201646 Mean    : 54330401 Caucasian:76099
##   3rd Qu.:230270888 3rd Qu.: 87545950 Hispanic  : 2037
##   Max.    :443867222 Max.    :189502619 Other    : 1506
##
##           gender       age        weight admission_type_id
##           Female     :54708 [70-80):26068 ?: 98569 Min.  :1.000
##           Male       :47055 [60-70):22483 [75-100) : 1336 1st Qu.:1.000
##           Unknown/Invalid: 3 [50-60):17256 [50-75)  : 897 Median :1.000
##                           [80-90):17197 [100-125): 625 Mean   :2.024
##                           [40-50): 9685 [125-150): 145 3rd Qu.:3.000
##                           [30-40): 3775 [25-50)   :  97 Max.   :8.000
##                           (Other): 5302 (Other)  :  97
##
##   discharge_disposition_id admission_source_id time_in_hospital   payer_code
##   Min.    : 1.000           Min.    : 1.000     Min.    : 1.000 ?  :40256
##   1st Qu.: 1.000           1st Qu.: 1.000     1st Qu.: 2.000 MC :32439
##   Median  : 1.000           Median : 7.000     Median : 4.000 HM : 6274
##   Mean    : 3.716           Mean   : 5.754     Mean   : 4.396 SP : 5007
##   3rd Qu.: 4.000           3rd Qu.: 7.000     3rd Qu.: 6.000 BC : 4655
##   Max.    :28.000           Max.   :25.000     Max.   :14.000 MD : 3532
##                           (Other): 9603
##
##           medical_specialty num_lab_procedures num_procedures
##           ?                  :49949   Min.    : 1.0     Min.    :0.00
##           InternalMedicine :14635   1st Qu.: 31.0    1st Qu.:0.00
##           Emergency/Trauma : 7565   Median  : 44.0    Median :1.00
##           Family/GeneralPractice: 7440   Mean   : 43.1    Mean   :1.34
##           Cardiology         : 5352   3rd Qu.: 57.0    3rd Qu.:2.00
##           Surgery-General   : 3099   Max.   :132.0    Max.   :6.00

```

```

## (Other) :13726
## num_medications number_outpatient number_emergency number_inpatient
## Min.   : 1.00   Min.   : 0.0000   Min.   : 0.0000   Min.   : 0.0000
## 1st Qu.:10.00   1st Qu.: 0.0000   1st Qu.: 0.0000   1st Qu.: 0.0000
## Median :15.00   Median : 0.0000   Median : 0.0000   Median : 0.0000
## Mean    :16.02   Mean    : 0.3694   Mean    : 0.1978   Mean    : 0.6356
## 3rd Qu.:20.00   3rd Qu.: 0.0000   3rd Qu.: 0.0000   3rd Qu.: 1.0000
## Max.    :81.00   Max.    :42.0000   Max.    :76.0000   Max.    :21.0000
##
## diag_1          diag_2          diag_3          number_diagnoses max_glu_serum
## 428   : 6862   276   : 6752   250   :11555   Min.   : 1.000 >200: 1485
## 414   : 6581   428   : 6662   401   : 8289   1st Qu.: 6.000 >300: 1264
## 786   : 4016   250   : 6071   276   : 5175   Median  : 8.000 None:96420
## 410   : 3614   427   : 5036   428   : 4577   Mean    : 7.423 Norm: 2597
## 486   : 3508   401   : 3736   427   : 3955   3rd Qu.: 9.000
## 427   : 2766   496   : 3305   414   : 3664   Max.   :16.000
## (Other):74419 (Other):70204 (Other):64551
## A1Cresult      metformin      repaglinide      nateglinide      chlorpropamide
## >7   : 3812   Down   : 575   Down   : 45   Down   : 11   Down   :     1
## >8   : 8216   No     :81778   No     :100227   No     :101063   No     :101680
## None:84748   Steady:18346   Steady: 1384   Steady:  668   Steady:    79
## Norm: 4990   Up     : 1067   Up     : 110   Up     :  24   Up     :     6
##
## glimepiride      acetohexamide      glipizide      glyburide      tolbutamide
## Down   : 194   No     :101765   Down   : 560   Down   : 564   No     :101743
## No     :96575   Steady:     1   No     :89080   No     :91116   Steady:    23
## Steady: 4670           Steady:11356   Steady: 9274
## Up     : 327           Up     : 770   Up     : 812
##
## pioglitazone      rosiglitazone      acarbose      miglitol      troglitazone
## Down   : 118   Down   :  87   Down   :    3   Down   :    5   No     :101763
## No     :94438   No     :95401   No     :101458   No     :101728   Steady:    3
## Steady: 6976   Steady: 6100   Steady: 295   Steady:  31
## Up     : 234   Up     : 178   Up     :   10   Up     :    2
##
## tolazamide      examide      citoglipton      insulin      glyburide.metformin
## No     :101727   No:101766   No:101766   Down   :12218   Down   :     6

```

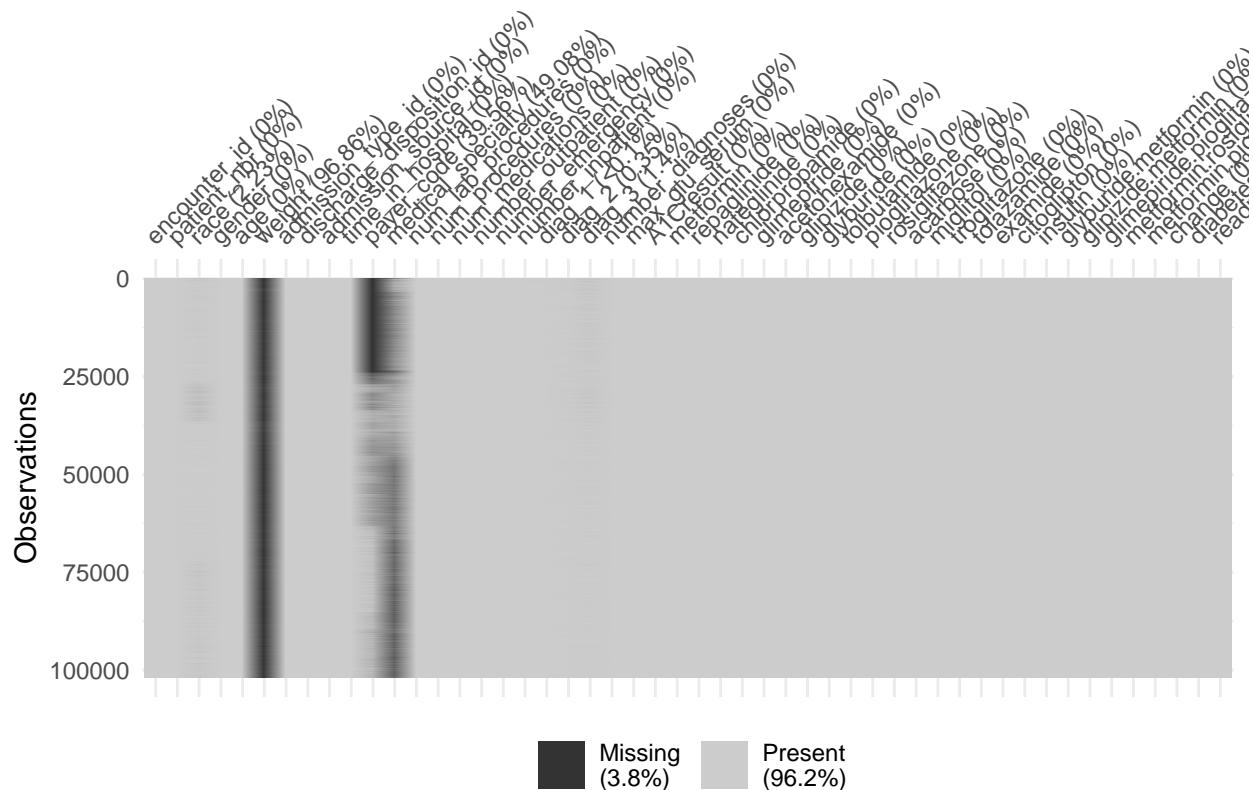
```
## Steady: 38 No :47383 No :101060
## Up : 1 Steady:30849 Steady: 692
## Up :11316 Up : 8
##
##
##
## glipizide.metformin glimepiride.pioglitazone metformin.rosiglitazone
## No :101753 No :101765 No :101764
## Steady: 13 Steady: 1 Steady: 2
##
##
##
##
## metformin.pioglitazone change diabetesMed readmitted
## No :101765 Ch:47011 No :23403 <30:11357
## Steady: 1 No:54755 Yes:78363 >30:35545
## NO :54864
##
##
##
##
```

## 2.2 Missing values

As we see in the result determined by the glimpse() function, the columns race, weight, payer\_code, gender, diag\_3, and medical specialty have some or significant amounts of missing values. To get a better picture, we will zoom in on these attributes and determine the amount and the percentage of missing values.

```
myvars <- c("race", "weight", "payer_code", "medical_specialty", "gender", "diag_3")
amountRecords <- length(rownames(episode.data))
tableMissingValues <- episode.data %>% summarise_each_(funz(TotalMissing = sum(. %in%
  c("?", "Unknown/Invalid"), na.rm = TRUE), MissingValuesPercentage = round(sum(. %in%
  c("?", "Unknown/Invalid"), na.rm = TRUE)/amountRecords * 100, 2)), myvars) %>%
  t()
tableMissingValues <- data.frame(Missing.Values = tableMissingValues[1:6, ], Missing.Val
  ])
rownames(tableMissingValues) <- myvars

# Added 27-09
library(naniar)
episode.data.missing <- replace(episode.data, episode.data == "?", NA)
vis_miss(episode.data.missing, warn_large_data = F)
```



	Missing.Values	Missing.Values.Percentage
race	2273.00	2.23
weight	98569.00	96.86
payer_code	40256.00	39.56
medical_specialty	49949.00	49.08
gender	3.00	0.00
diag_3	1423.00	1.40

As we notice in table 1, attribute weight is almost entirely made-up of missing values. For that reason, it is a candidate for removal. The same can be said for payer\_code and medical\_specialty, where 39.56% and 49.08% of values are missing, respectively. Payer\_code has no significant value to our research goals and questions; therefore, it can be removed. Medical\_specialty can be of value as it contains a range of useful information. Missing values can be set to ‘Missing’ or reevaluated based on other valuations; doing this is not without risk, as it is almost half of the attributes’ values. Deleting all missing values is not doable as it removes half of all records! Setting it to ‘Missing’ seems to be the most logical option. Attributes race and diag\_3 seem to have little amounts of missing data, but removing these values can alter our outcomes; therefore, we will use the same approach as to medical\_specialty. Gender only has three missing values. Considering the amount does not account to even one percent, removal does not seem to harm, and gender is, in most cases, considered binary data.

In the next section, we will introduce a new attribute (HbA1c measurement) based on an A1C test and the response to that result, which is defined as a change in diabetic medication. This test was performed at the time of hospital admission. We consider four groups of encounters:

1. no HbA1c test performed;
2. HbA1c performed and in normal range;
3. HbA1c performed and the result is greater than 8
4. HbA1c performed, result is greater than 8

```
encounter.data <- encounter.data %>% mutate(hba1c_res = ifelse(A1Cresult == "None"),
  1, NA)) %>% mutate(hba1c_res = ifelse((A1Cresult %in% c("Norm", ">7")) & is.na(hba1c_res),
  2, hba1c_res)) %>% mutate(hba1c_res = ifelse((A1Cresult == ">8") & (change ==
  "No") & is.na(hba1c_res), 3, hba1c_res)) %>% mutate(hba1c_res = ifelse((A1Cresult ==
  ">8") & (change == "Ch") & is.na(hba1c_res), 4, hba1c_res))
encounter.data$hba1c_res <- as.factor(encounter.data$hba1c_res)
```

(Date: 18-20 September, 2020)

Before we start analyzing our dataset, it is crucial to check on duplicates as our dataset contains multiple inpatient visits for some patients. These observations cannot be statistically independent and would create noise. We thus declare that only one encounter per patient is optimal.

```

detach("package:plyr", unload = TRUE)
duplicateCheck <- encounter.data %>% group_by(patient_nbr) %>% summarise(count = n())
  filter(count > 1)
table.dup <- data.frame(`Initial number of records` = nrow(encounter.data), `Number of Duplicates` = nrow(duplicateCheck),
  `Difference in percent` = round((nrow(encounter.data) - nrow(duplicateCheck)) / nrow(encounter.data) * 100, 2))

```

	Initial.number.of.records	Number.of.Duplicates	Difference.in.percent
1	101766	16773	-16.48

The initial dataset started with NA records, of which NA were duplicates. A potential removal of these records would leave us with 84993 records, a total loss of NA%. A small price to pay for a statistically independent dataset.

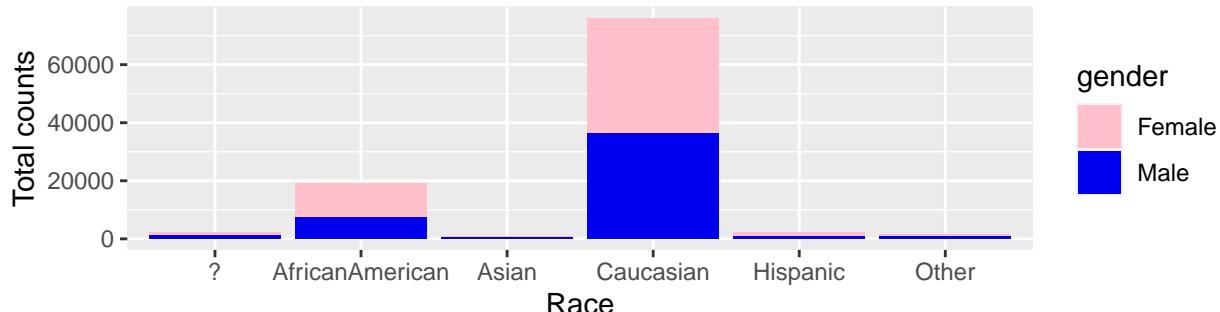
## 2.3 Categorial attributes

In this section, we take a look at variations in and between attributes. The attribute gender consists of three values ('female,' 'male,' and 'missing/unknown'). Since observe and use this attribute, it is essential to remove the abundant value.

```
# First we need to remove all records containing the third option of gender:  
# 'Missing/unknown', as this can be lethal in further analysis  
encounter.data <- encounter.data %>% select(everything()) %>% filter(gender != "Unknown"  
    droplevels()  
  
# Count specific variables of interest  
agg <- count(encounter.data, age, gender, A1Cresult, race)  
  
# Specify a color per value for visualization  
ecols <- c(Female = "pink", Male = "blue2")  
p1 <- ggplot(agg) + geom_col(aes(x = race, y = n, fill = gender)) + scale_fill_manual(v  
    labs(title = "Race distributed with gender", x = "Race", y = "Total counts",  
        subtitle = "1(a)")  
p2 <- ggplot(agg) + geom_col(aes(x = age, y = n, fill = gender)) + scale_fill_manual(va  
    labs(title = "Age distributed with gender", x = "Age", y = "Total counts", subtitle  
  
ecols <- c(`>7` = "blue", `>8` = "orange", None = "green", Norm = "yellow")  
p3 <- ggplot(agg) + geom_col(aes(x = age, y = n, fill = A1Cresult)) + scale_fill_manual  
    labs(title = "Age distributed with A1C test result", x = "Age", y = "Total counts",  
        subtitle = "2")  
  
# Arrange p1 and p2 in a multiplot and leave p3 on its own  
grid.arrange(p1, p2)
```

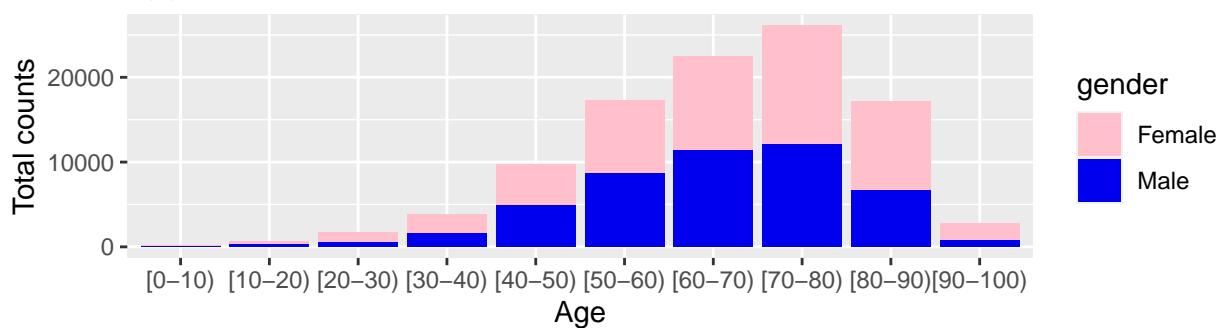
### Race distributed with gender

1(a)



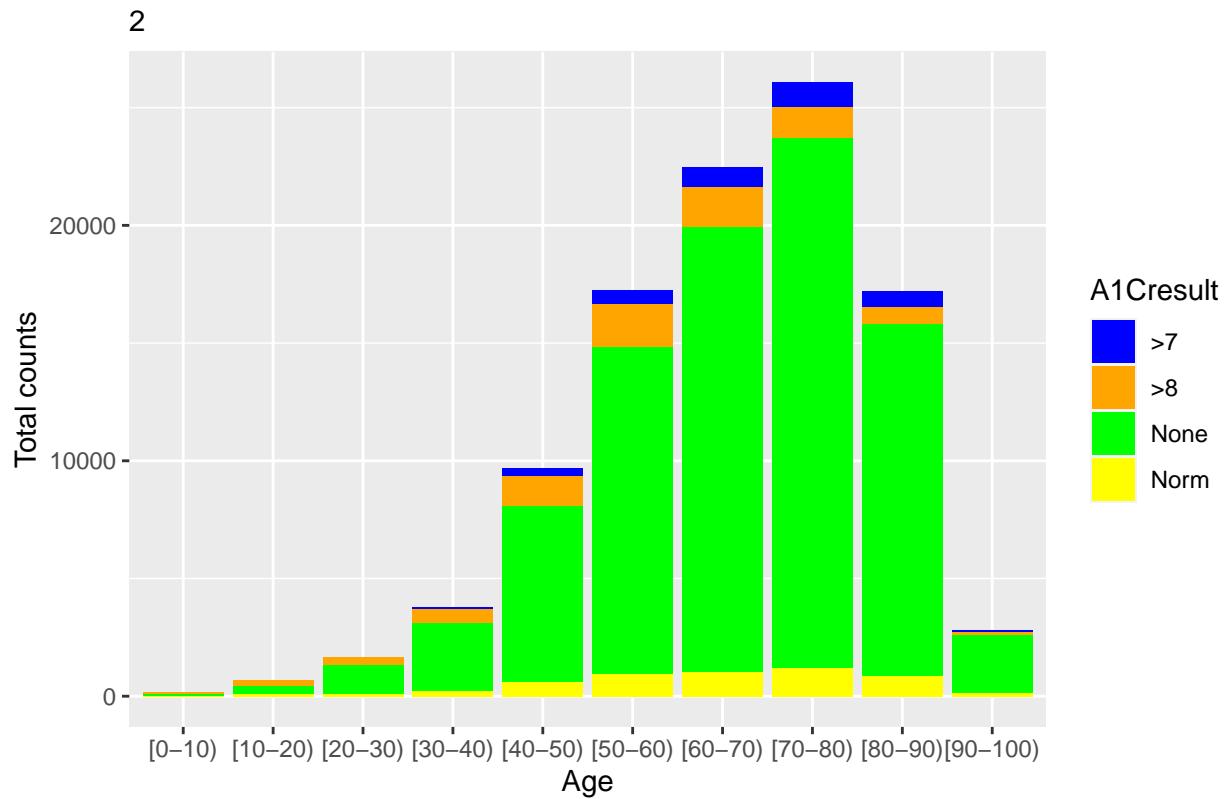
### Age distributed with gender

1(b)



```
grid.arrange(p3)
```

## Age distributed with A1C test result



Looking at figure 1a, we notice that most patients are from the Caucasian race with almost equivalent male and female ratios. We expect a rise in the number of diabetes patients when looking at older population groups. Figure 1b shows this exact prediction; the data is negatively skewed and increases in numbers per older age group. We observe the same results in figure 2, where an increase in A1C test is shown when comparing older population groups. Additionally, the figure gives an essential insight for many patients - in most cases, the test was not conducted and shows that strategies surrounding testing diabetes are not normalized in hospital protocols.

This graph does not, however, make a distinction between non-ICU and ICU patients. The authors of the original dataset analyses stated that ICU departments' protocols have a stricter policy surrounding testing for diabetes. When comparing non-ICU and ICU patient records, we need to construct a new attribute called 'icu\_or\_non,' comprising data from attributes admission\_type\_id, admission\_source\_id, and discharge\_disposition\_id. We distinguish between non-ICU and ICU patients.

```
encounter.data <- encounter.data %>% # Removing dead patients
filter(dischargeDisposition_id != 11) %>% # Distinction between ICU and non-ICU patient
mutate(admission_type_id = ifelse(admission_type_id %in% c(1, 2, 7), "ICU patient",
"Non-ICU patient")) %>% mutate(admission_source_id = ifelse(admission_source_id %in%
c(4, 7, 10, 12, 26), "ICU patient", "Non-ICU patient")) %>% mutate(discharge_dispositi
```

```

  c(13, 14, 19, 20, 21), "ICU patient", "Non-ICU patient"))
# Now that we changed the labels, let us discover how many ICU and non-ICU
# patients we have
encounter.data <- encounter.data %>% # All columns need to be equal to each other
mutate(icu_or_non = ifelse(admission_source_id == discharge_disposition_id & admission_t
  discharge_disposition_id, admission_source_id, ifelse(admission_source_id ==
  discharge_disposition_id, admission_source_id, admission_source_id)))
encounter.data$icu_or_non <- as.factor(encounter.data$icu_or_non)

agg <- count(encounter.data, gender, A1Cresult, icu_or_non, race)
ecols <- c(Female = "pink", Male = "blue2")
p1 <- ggplot(agg) + geom_col(aes(x = icu_or_non, y = n, fill = gender)) + scale_fill_ma
  labs(title = "ICU statistics distributed with gender", x = "ICU or non-ICU patient",
  y = "Total counts", subtitle = "3(a)")

p2 <- ggplot(agg) + geom_col(aes(x = icu_or_non, y = n, fill = race)) + labs(title = "IC
  x = "ICU or non-ICU patient", y = "Total counts", subtitle = "3(b)")

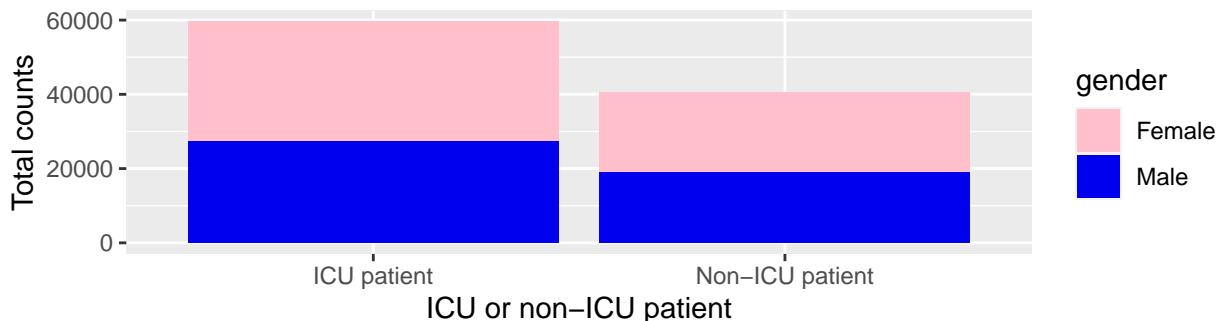
ecols <- c(`ICU patient` = "blue", `Non-ICU patient` = "red")
p3 <- ggplot(agg) + geom_col(aes(x = A1Cresult, y = n, fill = icu_or_non)) + scale_fill_
  labs(title = "ICU statistics distributed with A1C test result", x = "ICU or non-ICU
  y = "Total counts", subtitle = "4")

grid.arrange(p1, p2)

```

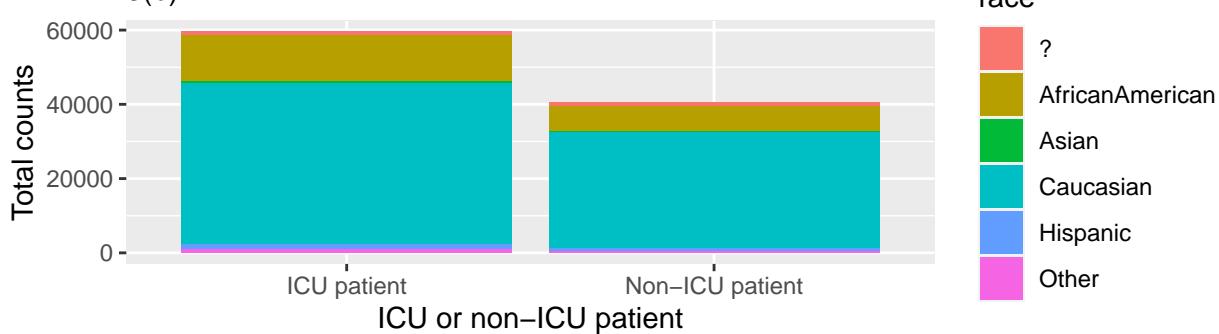
ICU statistics distributed with gender

3(a)



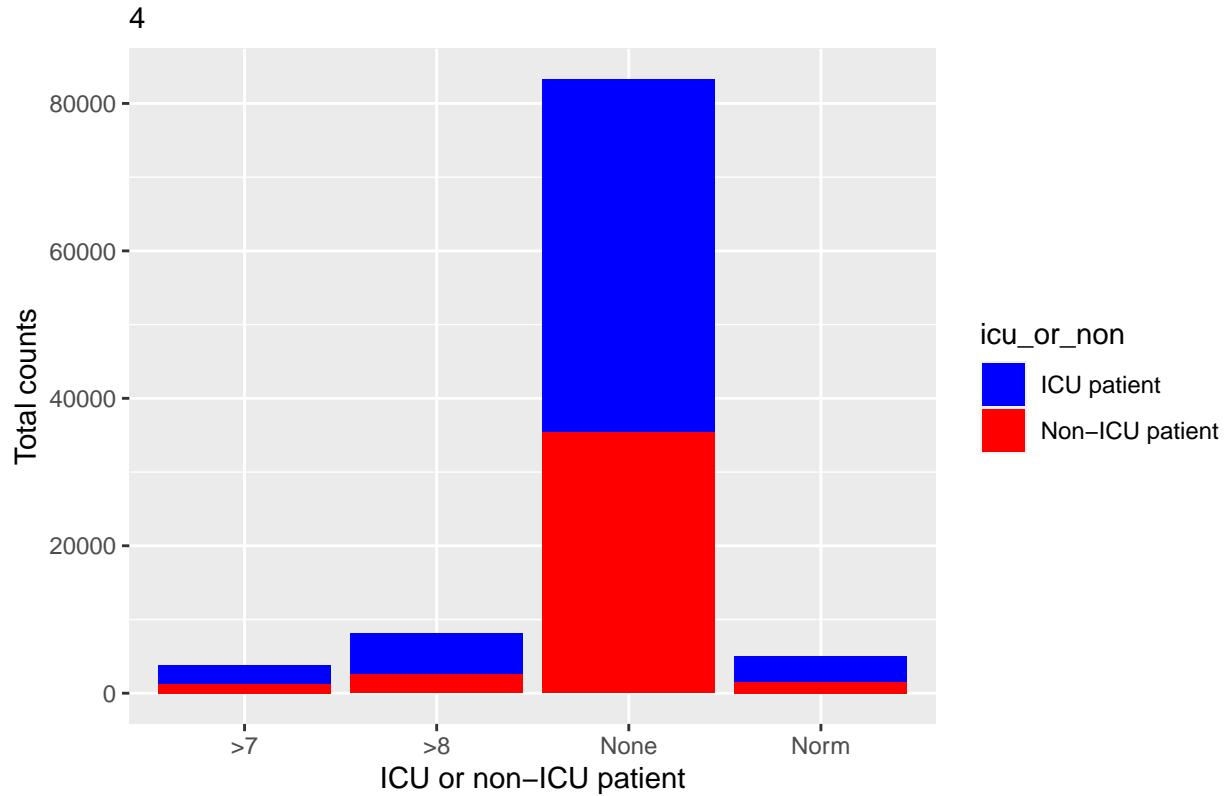
ICU statistics distributed with race

3(b)



```
grid.arrange(p3)
```

## ICU statistics distributed with A1C test result



We observe the result from the distinction between ICU and non-ICU in figures 3a, 3b, and 4. We notice a smaller population of non-ICU patients in all mentioned figures, whereas the ratio of not only gender but also seen in 3b, where race is depicted, ratios stay consequently the same. Figure 4 shows that no matter the patient type (ICU or non-ICU), protocols surrounding diabetes testing is not firmly conducted. Keeping the distinction between non-ICU and ICU might not be necessary. However, it can be of complimentary use when starting with machine learning; the data is binary, and it allows the removal of three attributes.

The attributes diag\_1, diag\_2, and diag\_3 consist of many three-digit ICD codes. Many of these codes belong together in a subgroup. In the following section, we will construct a better way of describing the exact diagnosis. ICD codes descriptions are retrieved from [2].

```
encounter.data <- encounter.data %>% mutate(diag_1 = case_when(diag_1 %in% c(390:459) ~
  "Circulatory", diag_1 %in% c(460:519) ~ "Respiratory", diag_1 %in% c(520:579) ~
  "Digestive", diag_1 %in% c(240:279) ~ "Diabetes", diag_1 %in% c(800:999) ~ "Injury",
  diag_1 %in% c(710:739) ~ "Musculoskeletal", diag_1 %in% c(580:629) ~ "Genitourinary",
  TRUE ~ "Others"))

encounter.data <- encounter.data %>% mutate(diag_2 = case_when(diag_2 %in% c(390:459) ~
  "Circulatory", diag_2 %in% c(460:519) ~ "Respiratory", diag_2 %in% c(520:579) ~
  "Digestive", diag_2 %in% c(240:279) ~ "Diabetes", diag_2 %in% c(800:999) ~ "Injury",
  diag_2 %in% c(710:739) ~ "Musculoskeletal", diag_2 %in% c(580:629) ~ "Genitourinary",
  TRUE ~ "Others")
```

```

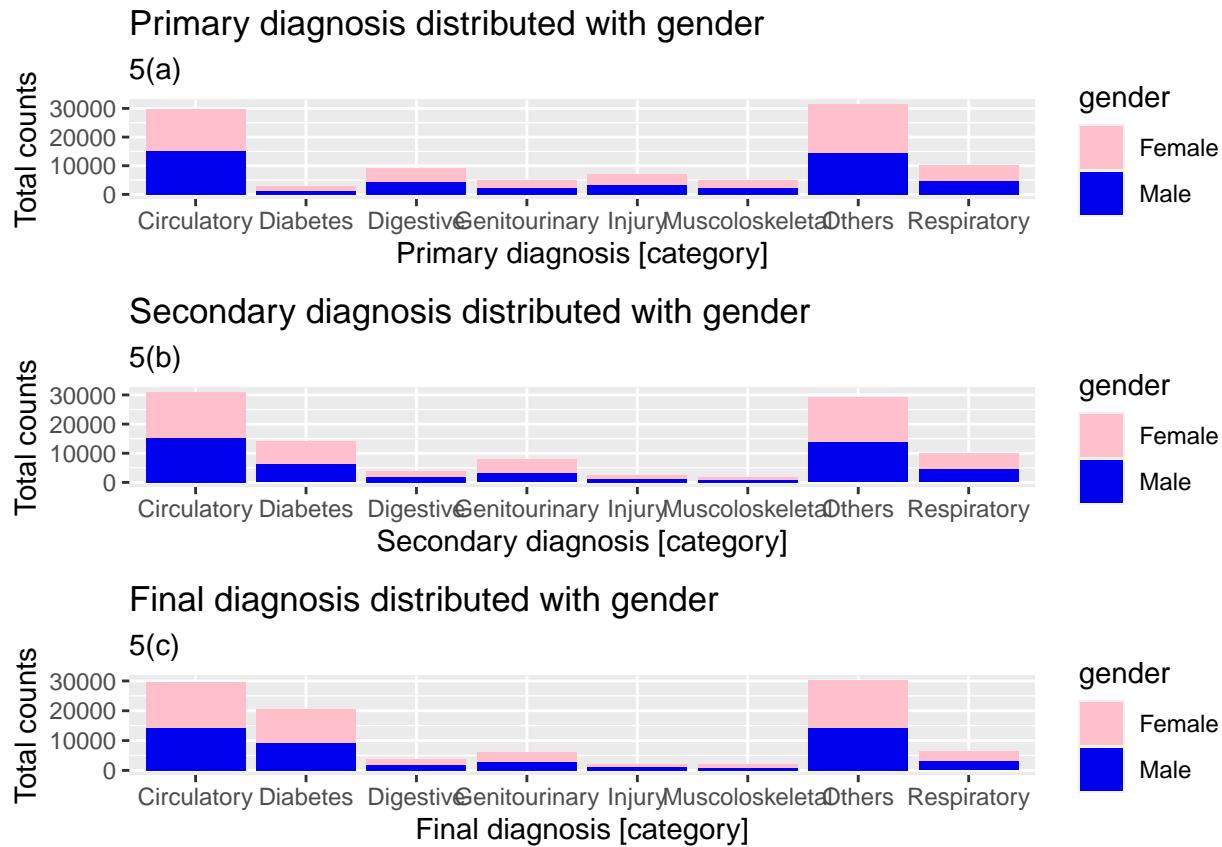
diag_2 %in% c(710:739) ~ "Musculoskeletal", diag_2 %in% c(580:629) ~ "Genitourinary"
TRUE ~ "Others"))
encounter.data <- encounter.data %>% mutate(diag_3 = case_when(diag_3 %in% c(390:459) ~
  "Circulatory", diag_3 %in% c(460:519) ~ "Respiratory", diag_3 %in% c(520:579) ~
  "Digestive", diag_3 %in% c(240:279) ~ "Diabetes", diag_3 %in% c(800:999) ~ "Injury",
  diag_3 %in% c(710:739) ~ "Musculoskeletal", diag_3 %in% c(580:629) ~ "Genitourinary",
  TRUE ~ "Others"))

# Convert to factors again
cols <- c("diag_1", "diag_2", "diag_3")
encounter.data[cols] <- lapply(encounter.data[cols], factor)
agg <- count(encounter.data, gender, diag_1, diag_2, diag_3)

# Specify a color per value for visualization
ecols <- c(Female = "pink", Male = "blue2")
p1 <- ggplot(agg) + geom_col(aes(x = diag_1, y = n, fill = gender)) + scale_fill_manual
  labs(title = "Primary diagnosis distributed with gender", x = "Primary diagnosis [category]",
       y = "Total counts", subtitle = "5(a)")
p2 <- ggplot(agg) + geom_col(aes(x = diag_2, y = n, fill = gender)) + scale_fill_manual
  labs(title = "Secondary diagnosis distributed with gender", x = "Secondary diagnosis",
       y = "Total counts", subtitle = "5(b)")
p3 <- ggplot(agg) + geom_col(aes(x = diag_3, y = n, fill = gender)) + scale_fill_manual
  labs(title = "Final diagnosis distributed with gender", x = "Final diagnosis [category]",
       y = "Total counts", subtitle = "5(c)")

grid.arrange(p1, p2, p3)

```



Looking at figure 5, we can see the results of the revaluation of attributes diag\_1, diag\_2, and diag\_3. These attributes depict the primary, secondary, and final diagnosis of a patient, respectively. Interestingly, the difference between the figures is the increase in the number of diabetes diagnoses. Due to not testing diabetes on the initial hospitalization, readmission rates increase as a patient's primary diagnosis is not sustainable. The data classification now shows a clearer picture and would undoubtedly be of fair use in the final dataset. The original authors did not keep diag\_2 and diag\_3, as it would make records too complex to achieve their goals. Removal of these two attributes is still up for discussion, but the analysis does not give - for now - a clear indication for potential removal.

In the next section, we look at attribute medical\_specialty and decide whether it is useful enough - considering the number of missing values - to be a candidate for the final dataset. We construct multiple plots that zoom in on the categories and valuations of this attribute.

```
library(pals)

# Count data we want to compare
agg <- count(encounter.data, age, medical_specialty, gender)

# Determine a color scheme for 71 values
ecols <- c(alphabet(26), cols25(25), glasbey(22))
```

```

# Order variables from high to low via '-n'
agg_ord <- mutate(agg, age = reorder(age, -n, sum), medical_specialty = reorder(medical_
  -n, sum))
# p1: bar plot combined, and p2: per gender
p1 <- ggplot(agg_ord) + geom_col(aes(x = age, y = n, fill = medical_specialty)) +
  scale_fill_manual(values = ecols) + theme(legend.position = "none") + labs(title =
    x = "Age [group]", y = "Total counts", subtitle = "6")

p1 <- p1 + facet_wrap(~gender) + theme(axis.text.x = element_text(angle = 90, vjust = 0,
  hjust = 1))

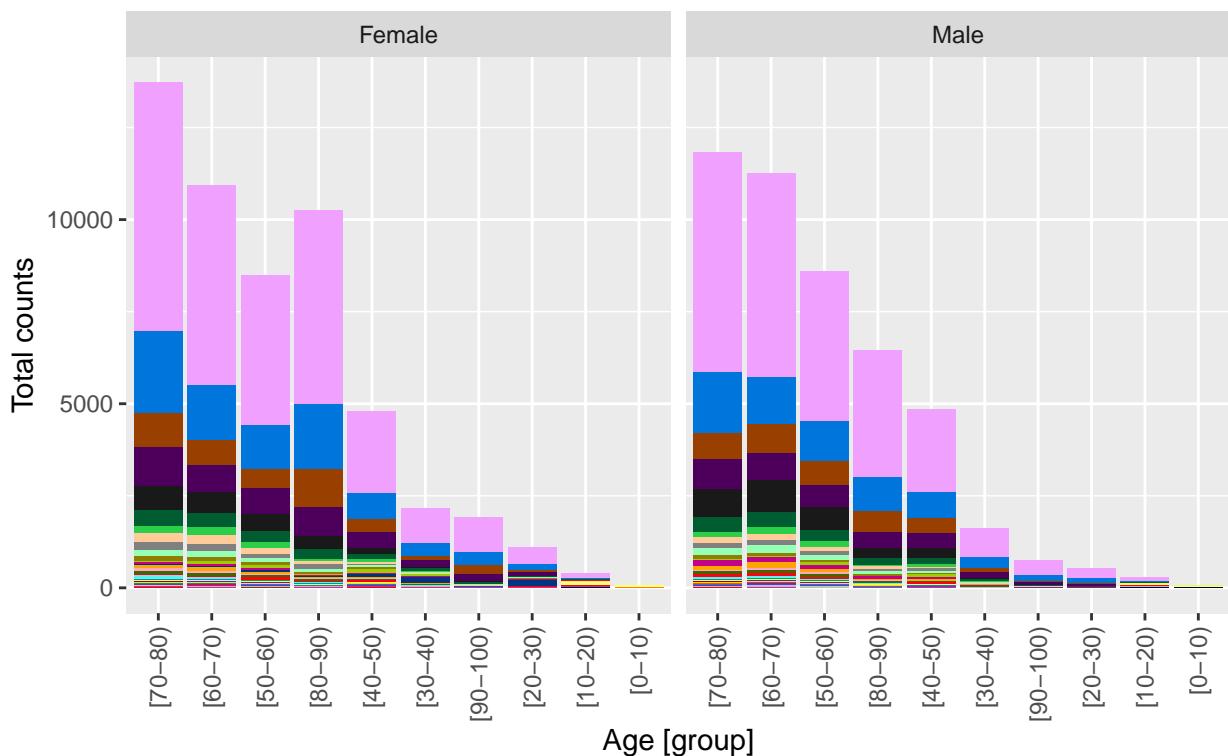
# Do the same as above, only now with pie charts
p2 <- ggplot(agg_ord) + geom_col(aes(x = 1, y = n, fill = medical_specialty), position =
  coord_polar(theta = "y")) + scale_fill_manual(values = ecols) + theme(legend.position =
  "none")
p2 <- p2 + facet_wrap(~gender) + theme_bw() + theme(legend.position = "none")

# Plot all results
grid.arrange(p1)

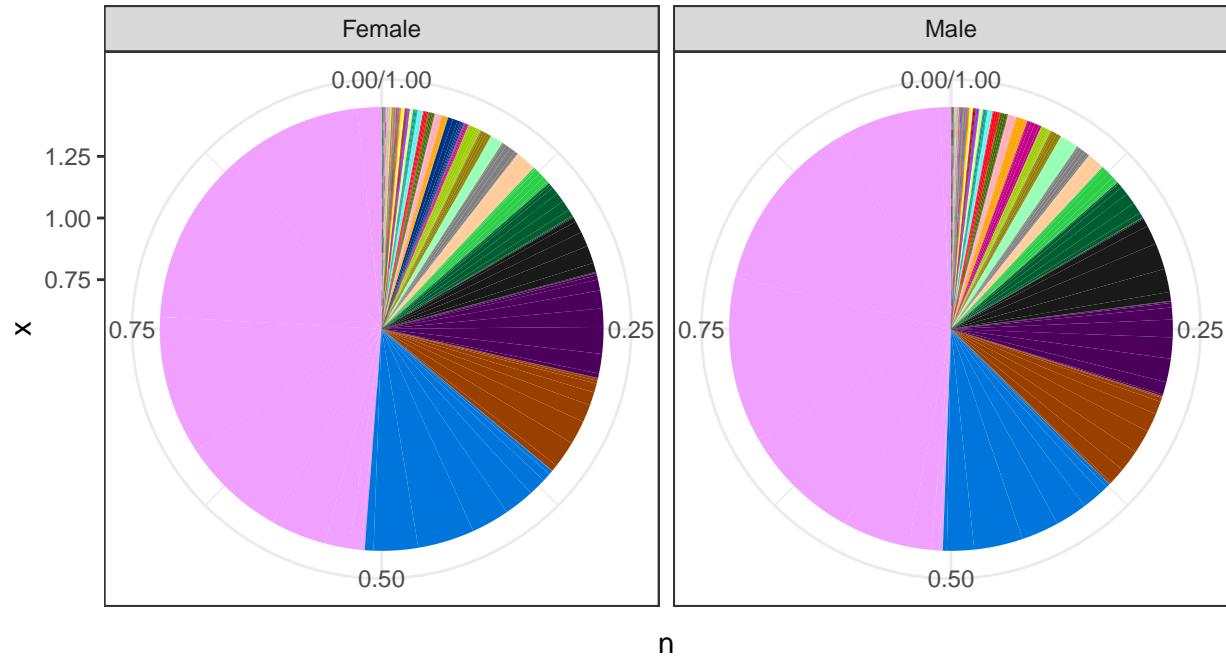
```

Medical specialty distributed with age groups and divided into gender

6



```
grid.arrange(p2)
```



Looking at both figures 6 and 7, we observe a large amount of valuation between all demographic attributes used. Depicted in purple are the missing values, as already established, make-up almost 50 percent of total records. Medical specialty is different from, for example, admission\_id, where we could shrink the attribute to a better format. This is an alternative to this attribute. We are considering this with the fact that medical specialty does not give better information than any diagnosis attribute, which also gives, maybe even more useful, guidance about an encounter's medical history. At this stage, removal of this attribute could not be of lethal harm.

For our final categorical attribute, and one of the most influential -according to the original authors- we will discuss readmitted. Readmitted is an attribute with three different valuations: '<30' for readmission within 30 days after release, '>30' for readmission after 30 days release, and 'NO' for no readmission.

```
agg <- count(encounter.data, readmitted, race, gender, age, diag_1, diag_2)
ecols <- c(`<30` = "red", `>30` = "blue", NO = "pink")

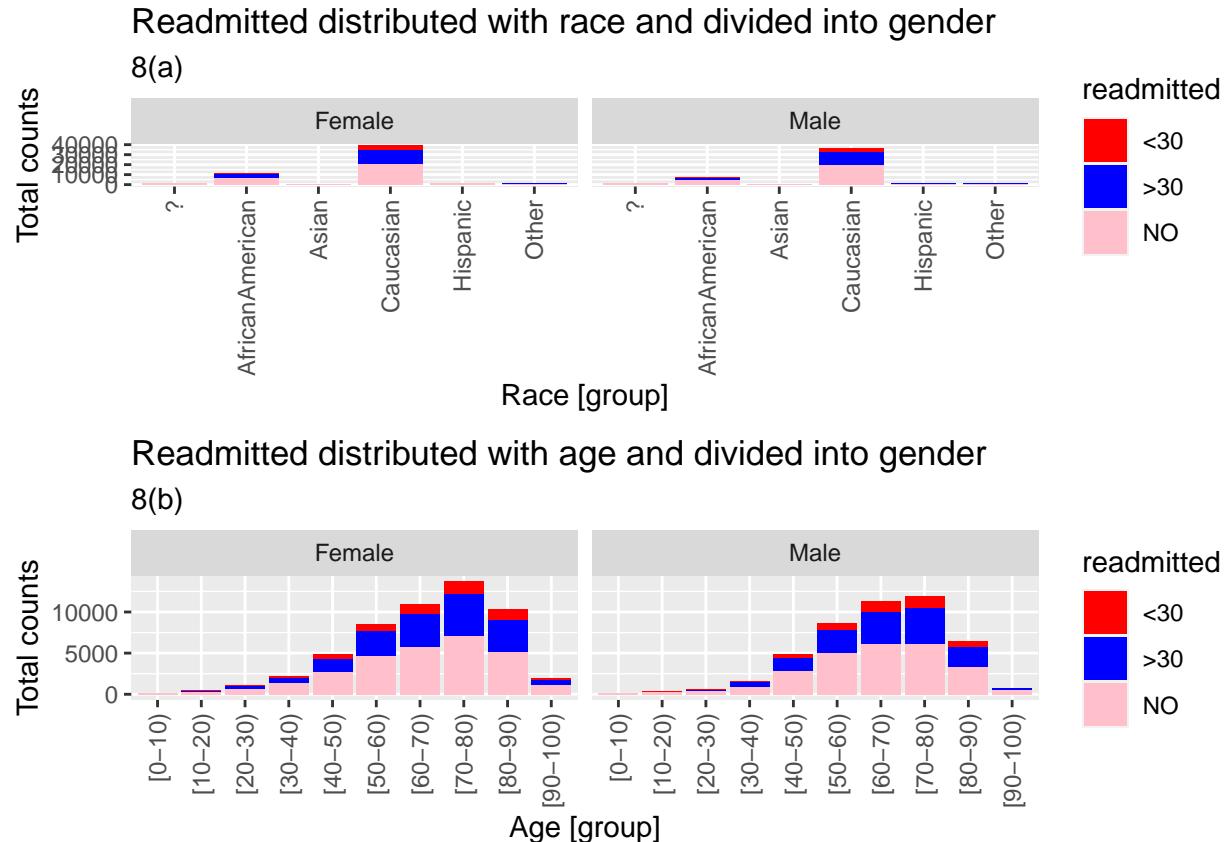
p1 <- ggplot(agg) + geom_col(aes(x = race, y = n, fill = readmitted)) + scale_fill_manual
```

```

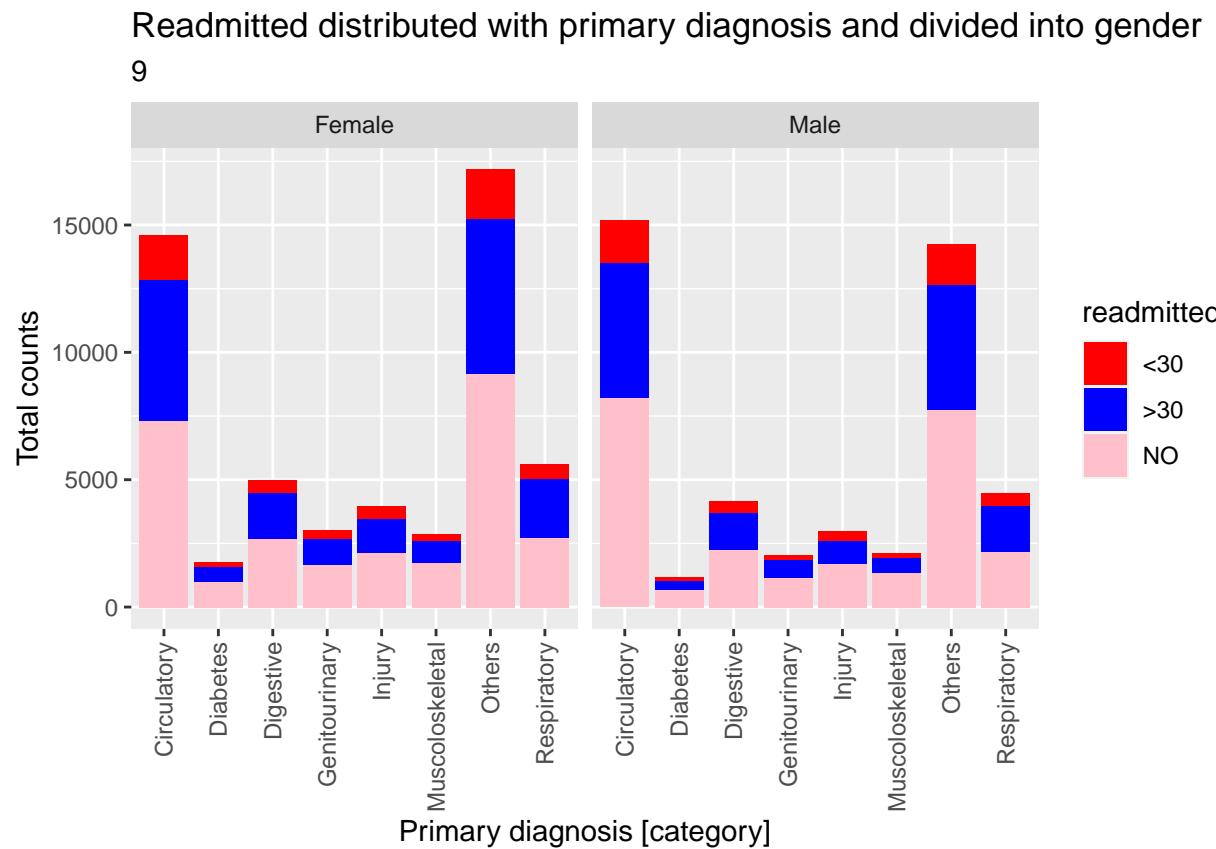
y = "Total counts", subtitle = "8(a)"
p1 <- p1 + facet_wrap(~gender) + theme(axis.text.x = element_text(angle = 90, vjust = 0,
hjust = 1))
p2 <- ggplot(agg) + geom_col(aes(x = age, y = n, fill = readmitted)) + scale_fill_manual
  labs(title = "Readmitted distributed with age and divided into gender", x = "Age [gr",
       y = "Total counts", subtitle = "8(b)")
p2 <- p2 + facet_wrap(~gender) + theme(axis.text.x = element_text(angle = 90, vjust = 0,
hjust = 1))
p3 <- ggplot(agg) + geom_col(aes(x = diag_1, y = n, fill = readmitted)) + scale_fill_manual
  labs(title = "Readmitted distributed with primary diagnosis and divided into gender",
       x = "Primary diagnosis [category]", y = "Total counts", subtitle = "9")
p3 <- p3 + facet_wrap(~gender) + theme(axis.text.x = element_text(angle = 90, vjust = 0,
hjust = 1))
p4 <- ggplot(agg) + geom_col(aes(x = diag_2, y = n, fill = readmitted)) + scale_fill_manual
  labs(title = "Readmitted distributed with secondary diagnosis and divided into gender",
       x = "Secondary diagnosis [category]", y = "Total counts", subtitle = "10")
p4 <- p4 + facet_wrap(~gender) + theme(axis.text.x = element_text(angle = 90, vjust = 0,
hjust = 1))

grid.arrange(p1, p2)

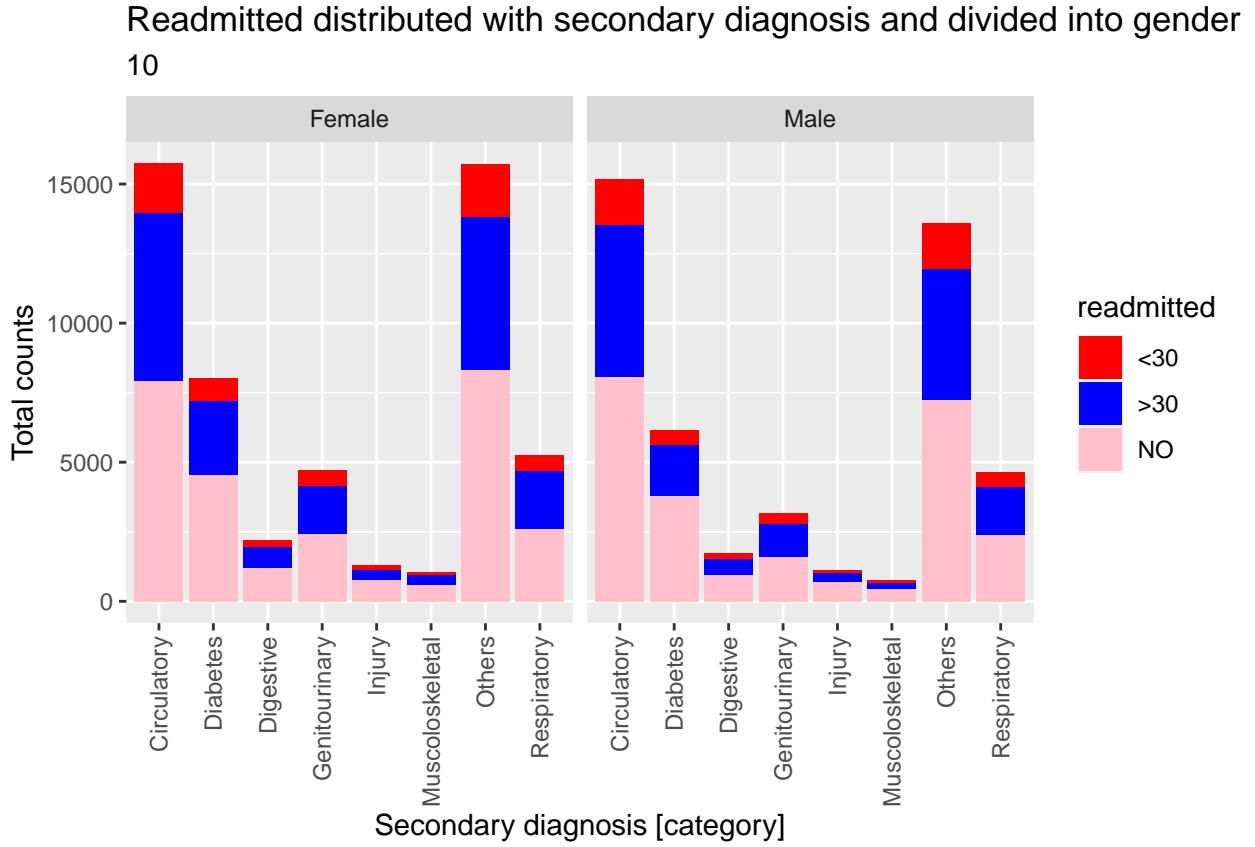
```



```
grid.arrange(p3)
```



```
grid.arrange(p4)
```



We observe the results of analyzing the readmitted attribute. We looked at readmitted with different demographics (gender, age, and race) and two of the three diagnosis attributes. Figure 8a shows that the Caucasian race has the most readmission rates in both genders. Observations made in figure 8b depict that as age increases, readmission rates also increase for females and males. This makes sense as immunity decreases with age, and the chance of recovery diminishes. Readmission rates decrease after the age of 80. The reasoning behind this might be death or transferring to a hospice or other facility.

Comparing the results in figures 9 and 10, we see a familiar trend: the increase of diabetes diagnoses between the primary and secondary diagnoses. This typically means that patients admitted were diagnosed differently, and over time, are getting a new diabetes diagnosis much quicker.

Seeing as we have multiple diagnosis categories but are somewhat interested in diabetes only, we can alter these attributes to have only two valuations ('diabetes' and 'other'). However, doing this can affect later analysis as we remove potentially valuable information.

## 2.4 Distribution - numeric attributes

We discussed many categorical data attributes, now turn to our few numeric data. It is important to have a good distribution between numeric data. If the range of distribution is

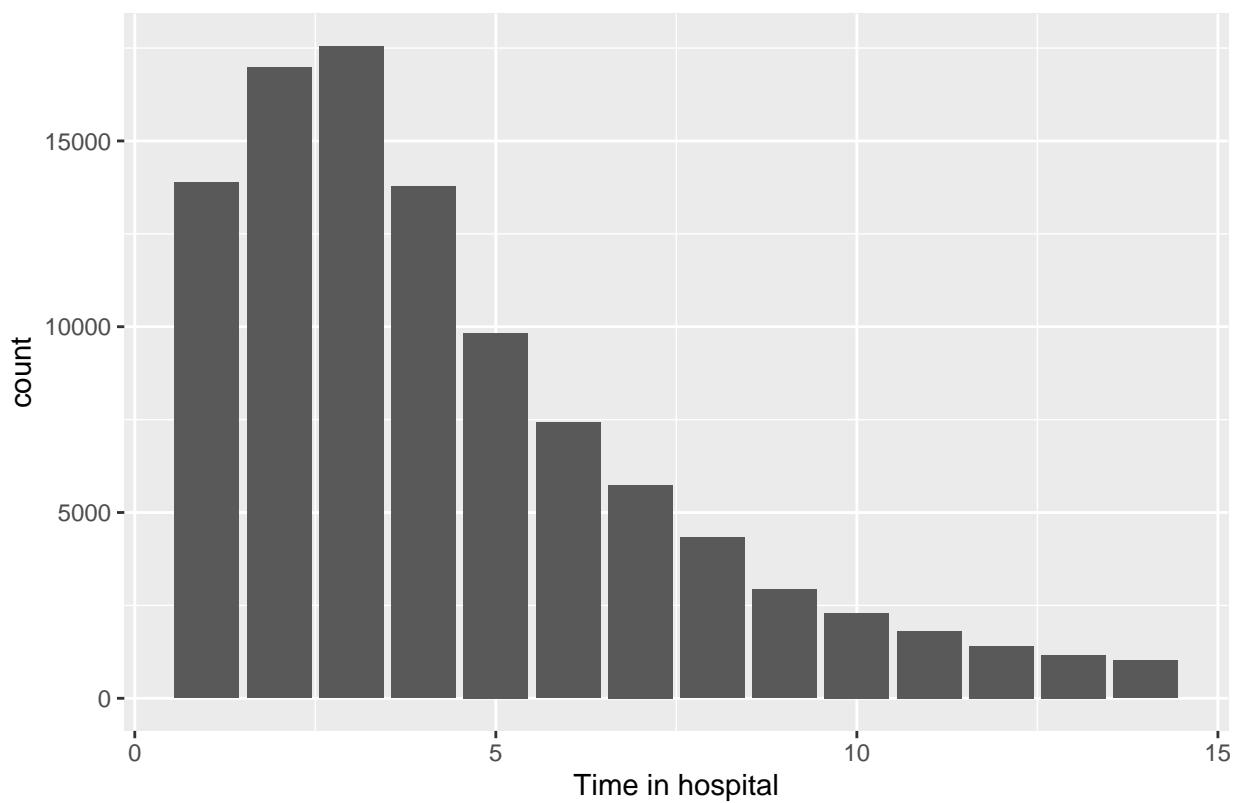
too large, then normalization is necessary. We construct multiple histograms, which we get through the histogram.plotter() function.

```
# Define a histogram plotter to construct multiple figures really quick
histogram.plotter <- function(attribute_1, number, attribute_name) {
  histogram <- ggplot() + geom_bar(mapping = aes(x = attribute_1)) + ggtitle(number) +
    xlab(attribute_name)
  print(histogram)
}

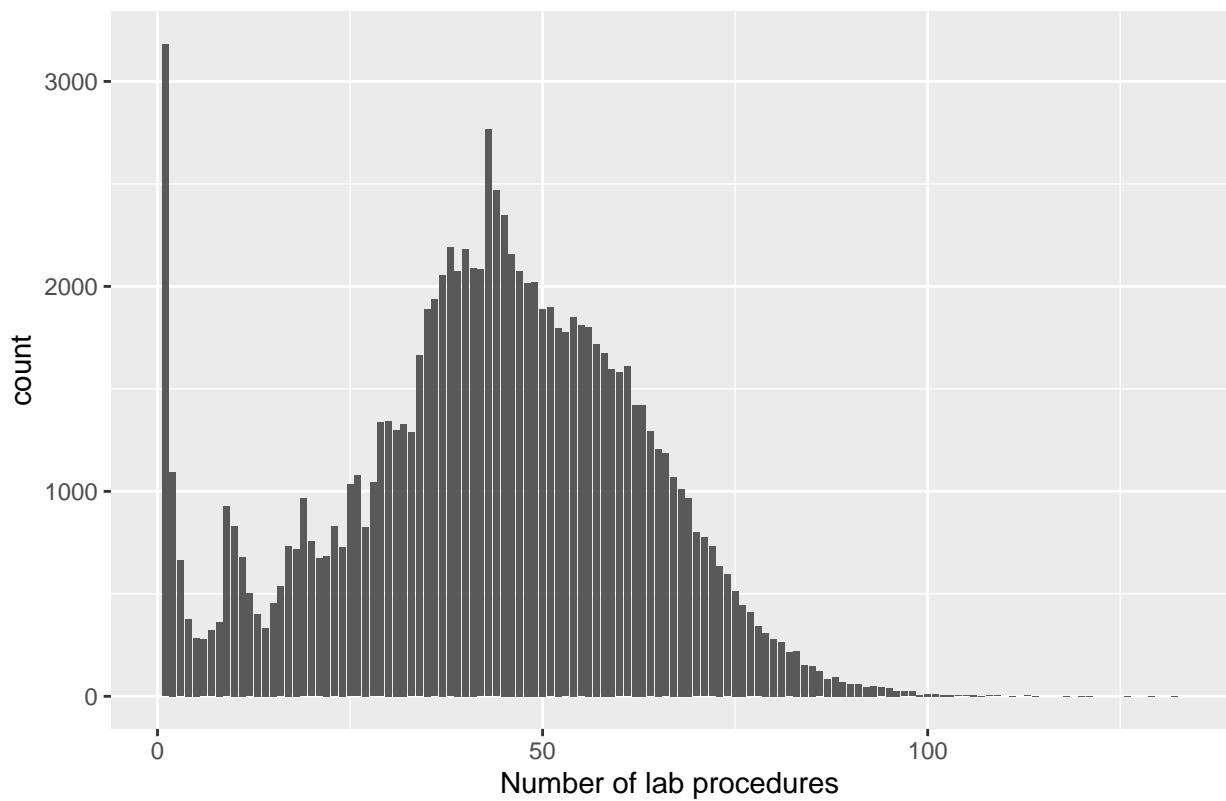
# Make a smaller dataset so that only the numeric data is collected
numeric.data <- select_if(encounter.data, is.numeric)
# Remove encounter and patient ID
numeric.data <- numeric.data[, -c(1, 2)]
smaller.codebook <- codebook[c(10, 13:18, 22), 2]
numeric.data.log <- log2(numeric.data + 0.1)

wrapper <- function(data) {
  p <- list()
  for (i in c(1:length(data))) {
    p[[i]] <- histogram.plotter(data[, i], LETTERS[i], smaller.codebook[i])
  }
  return(p)
}
numeric.normal <- wrapper(numeric.data)
```

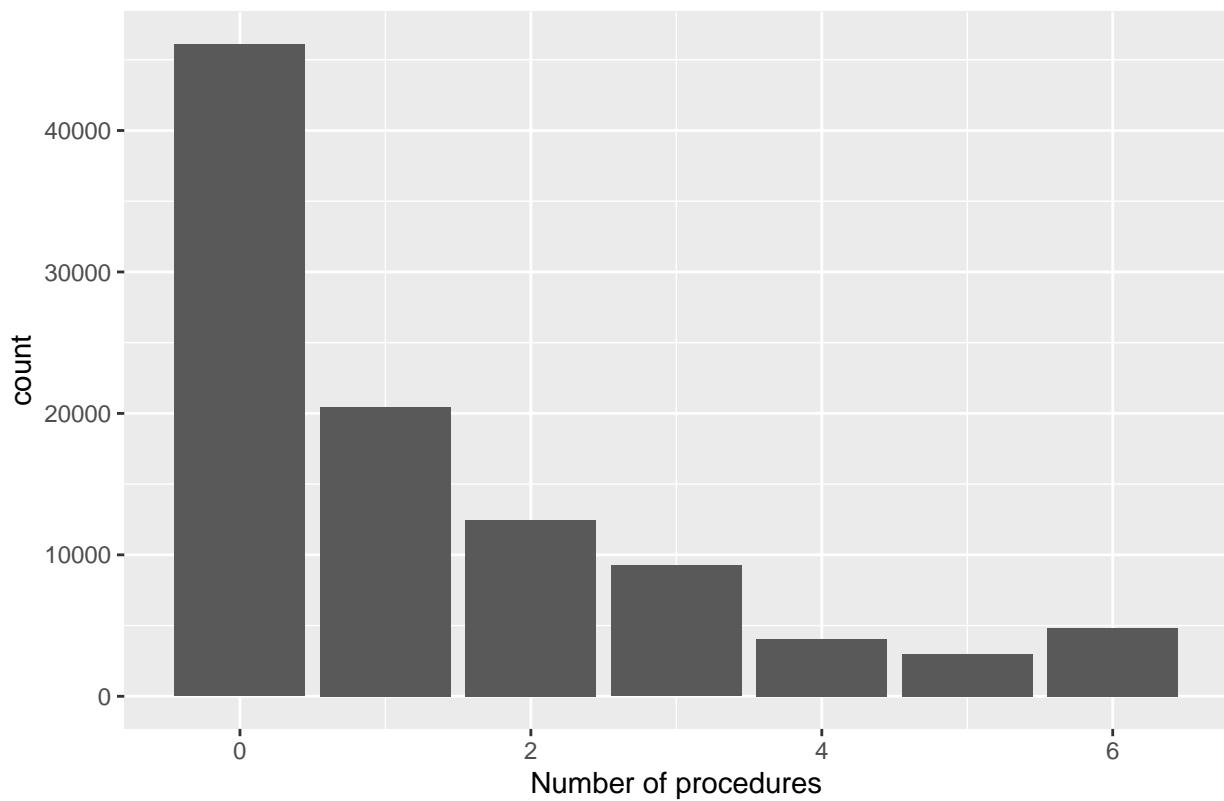
A

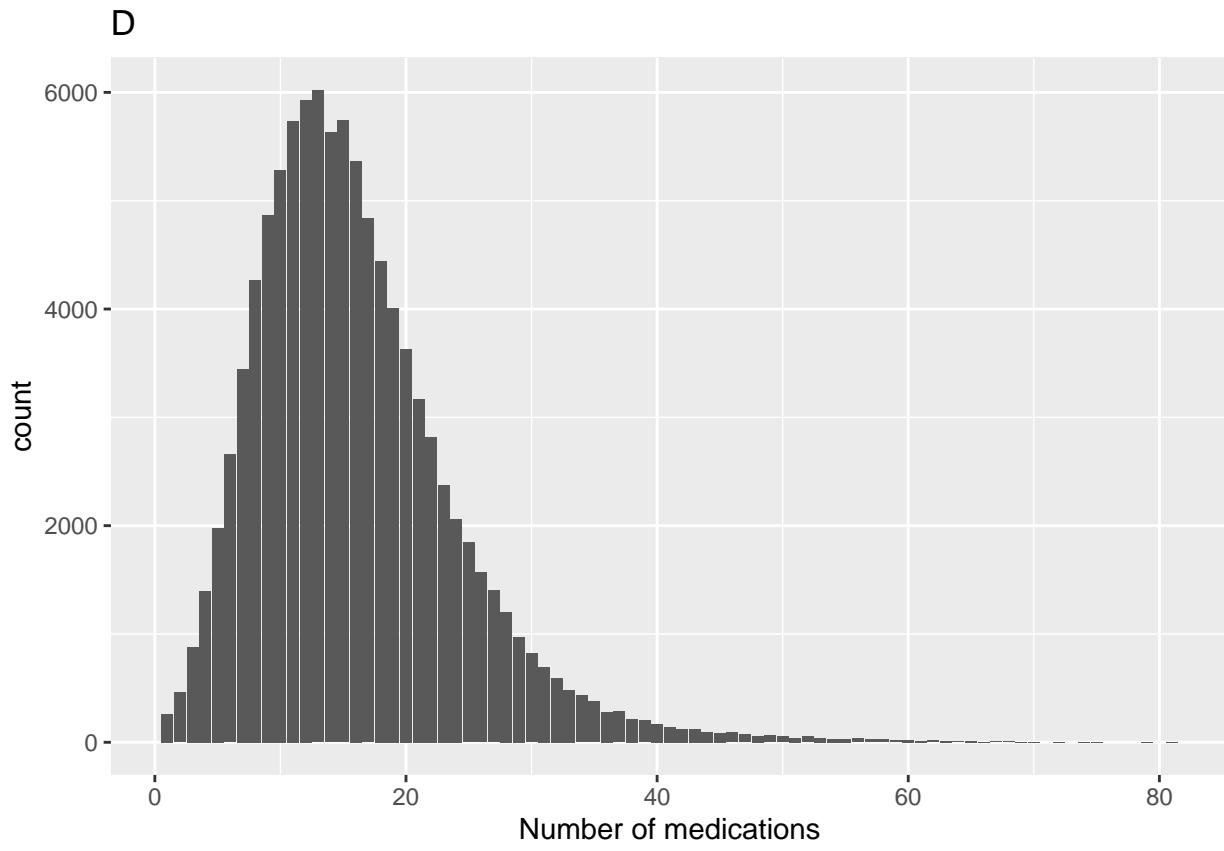


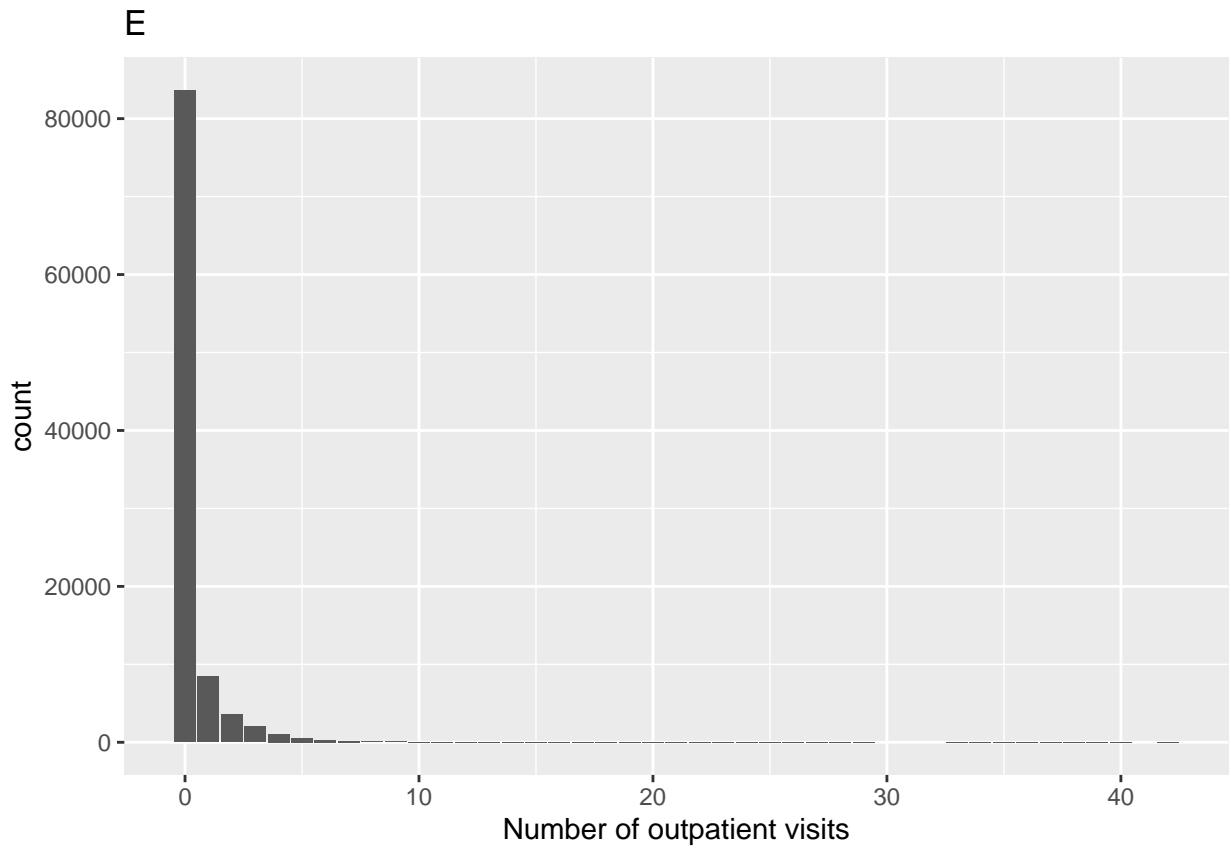
**B**



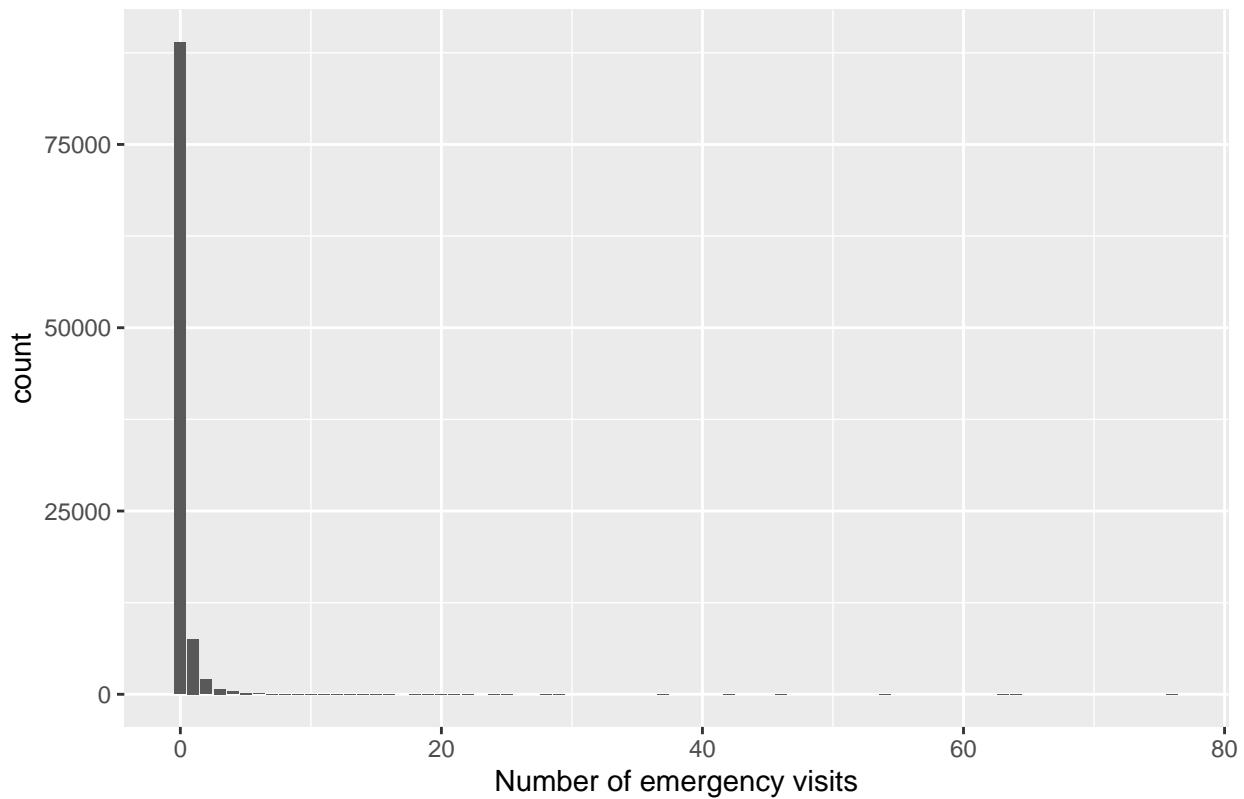
C

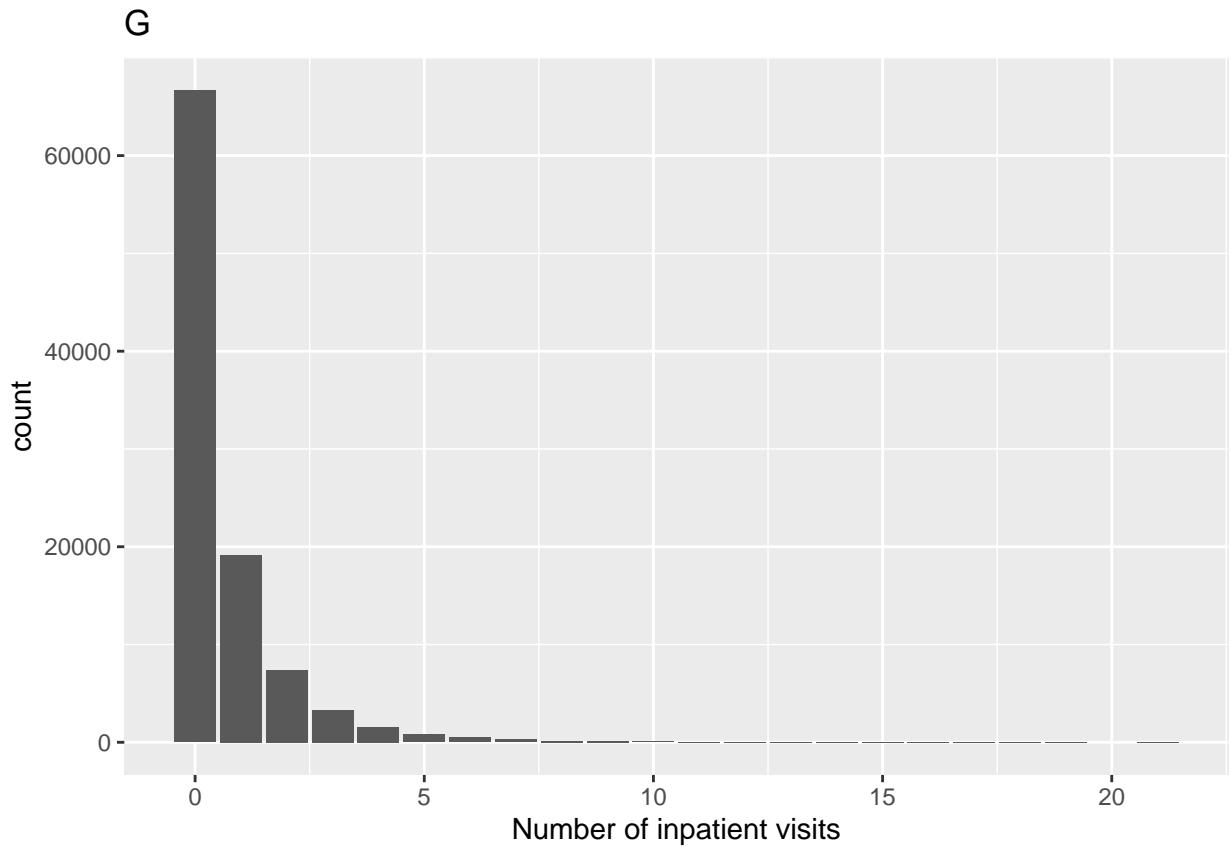


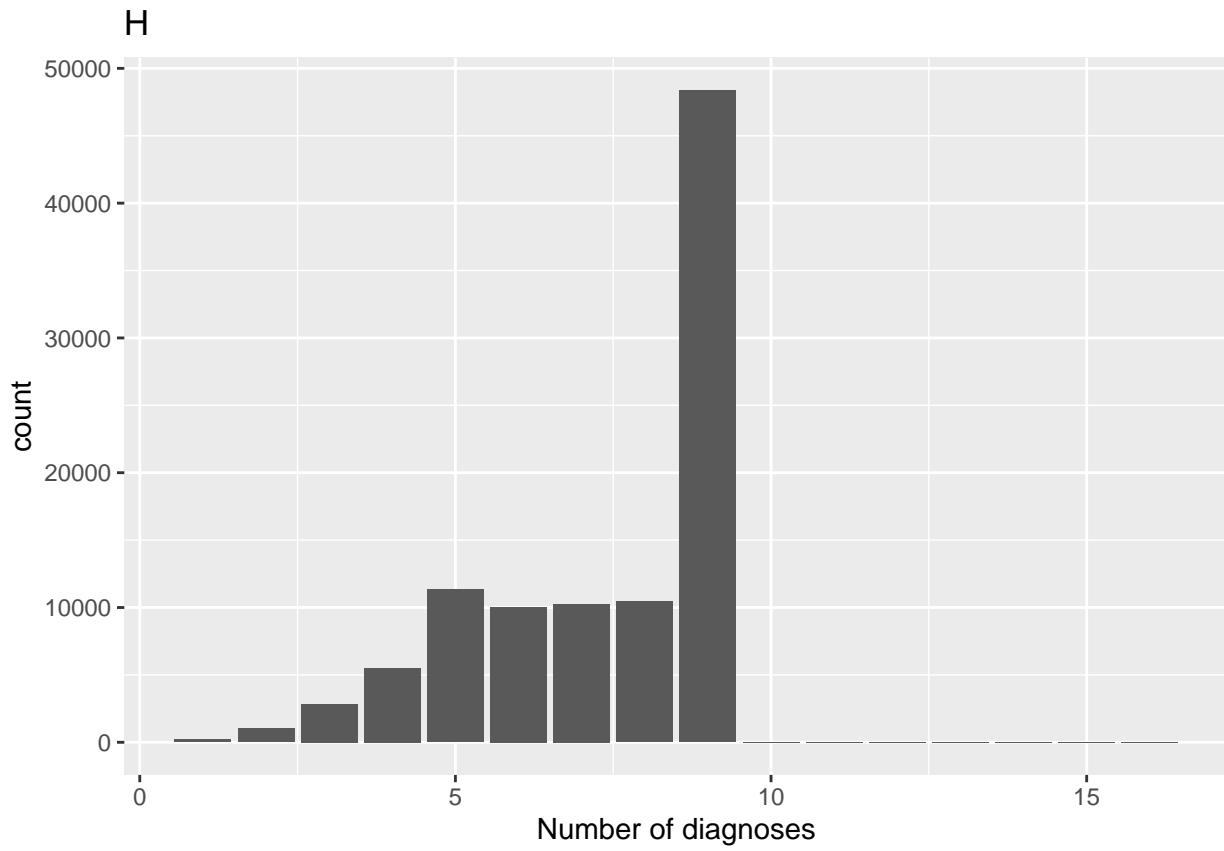




F

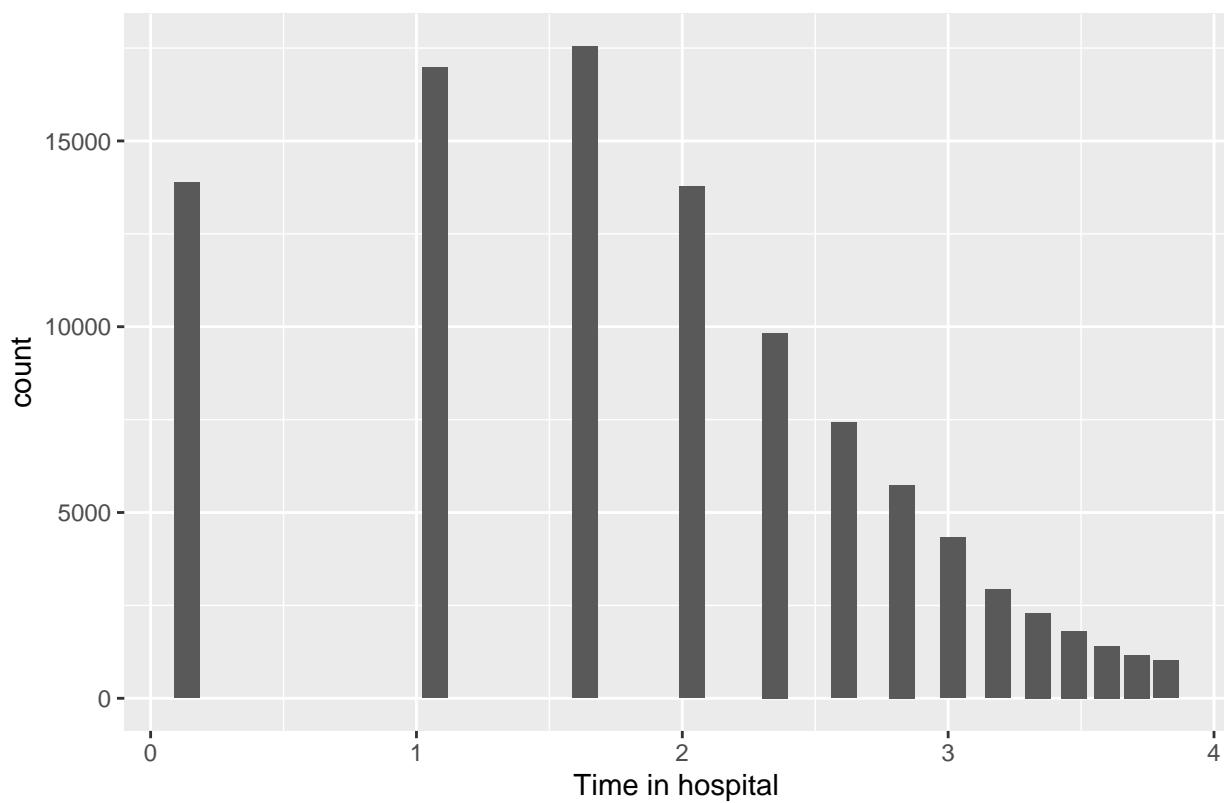




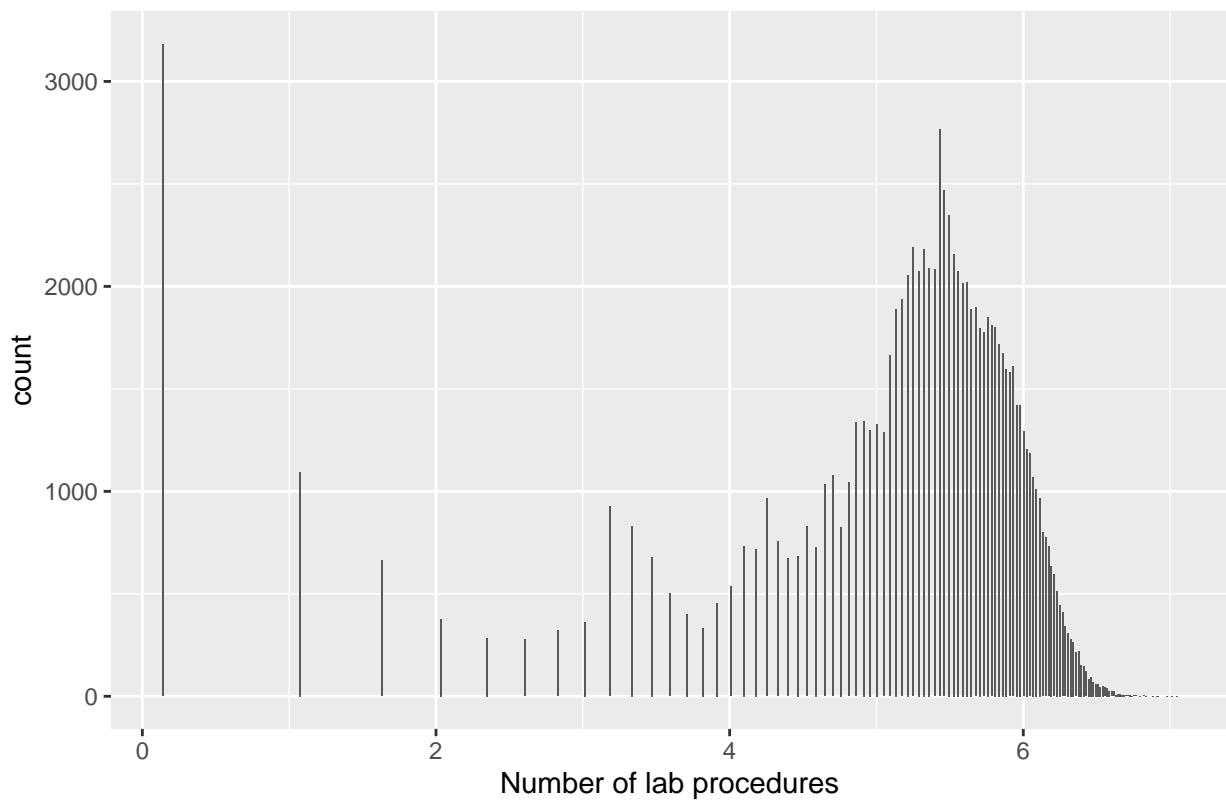


```
numeric.log <- wrapper(numeric.data.log)
```

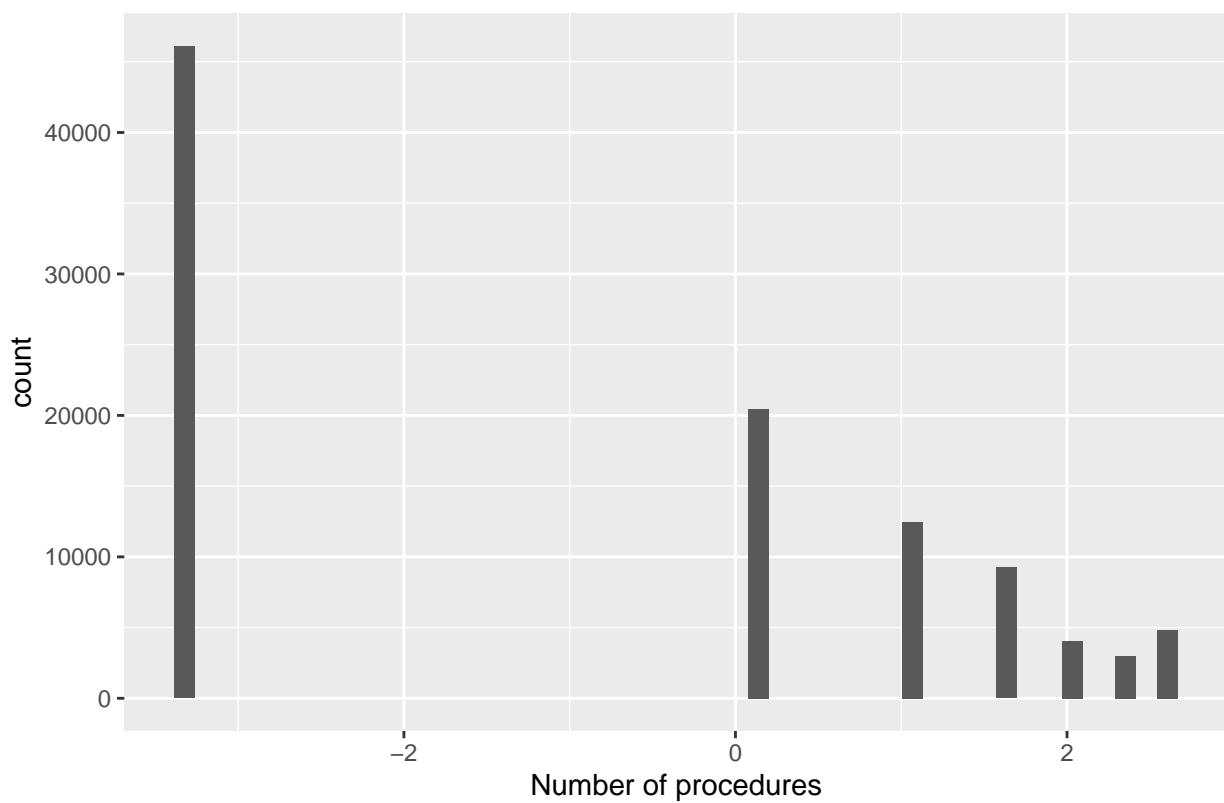
A

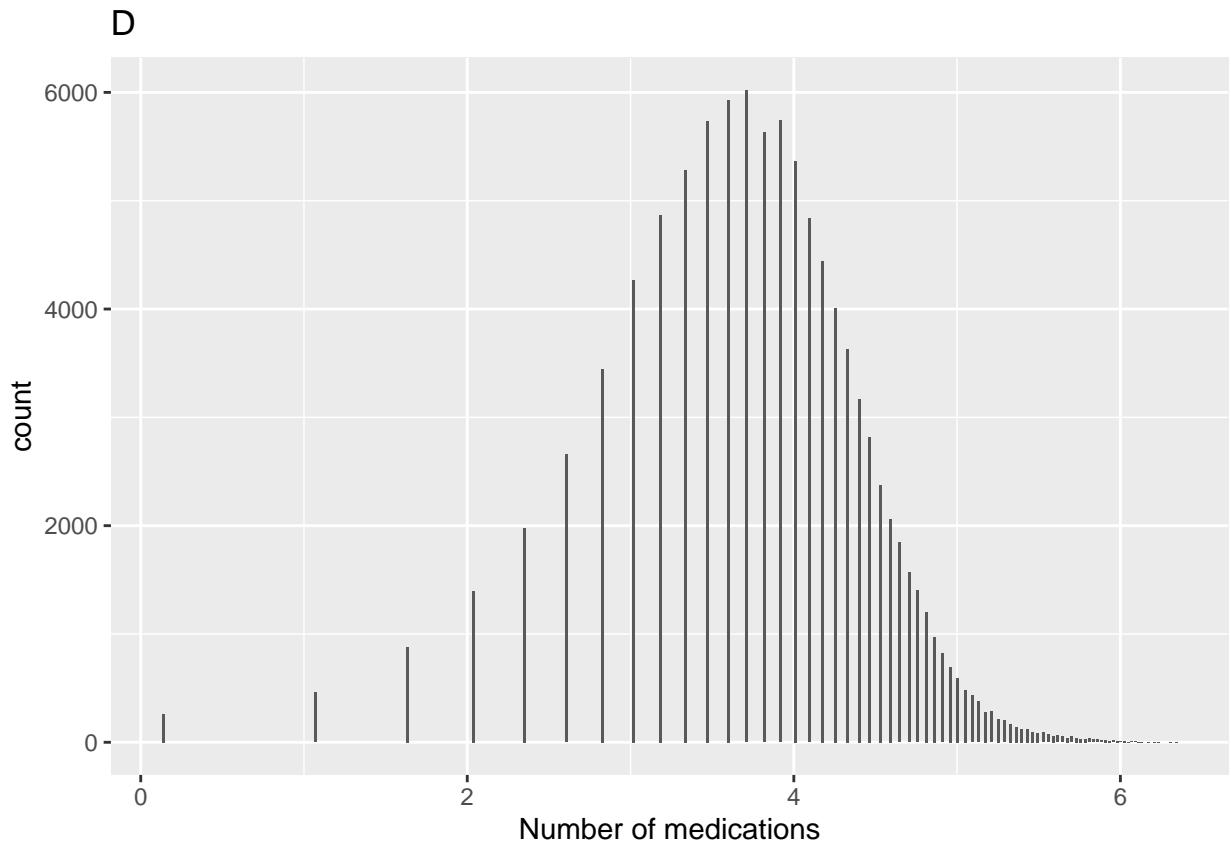


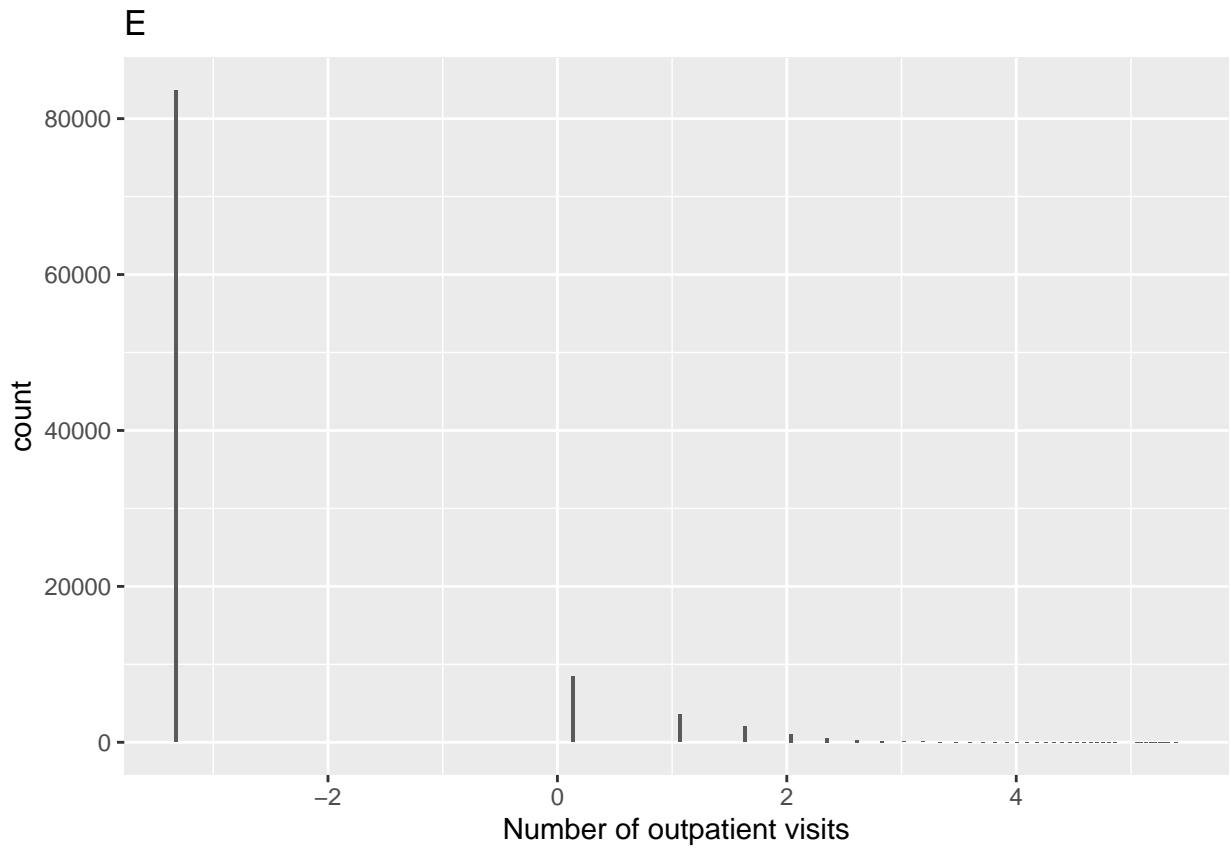
B



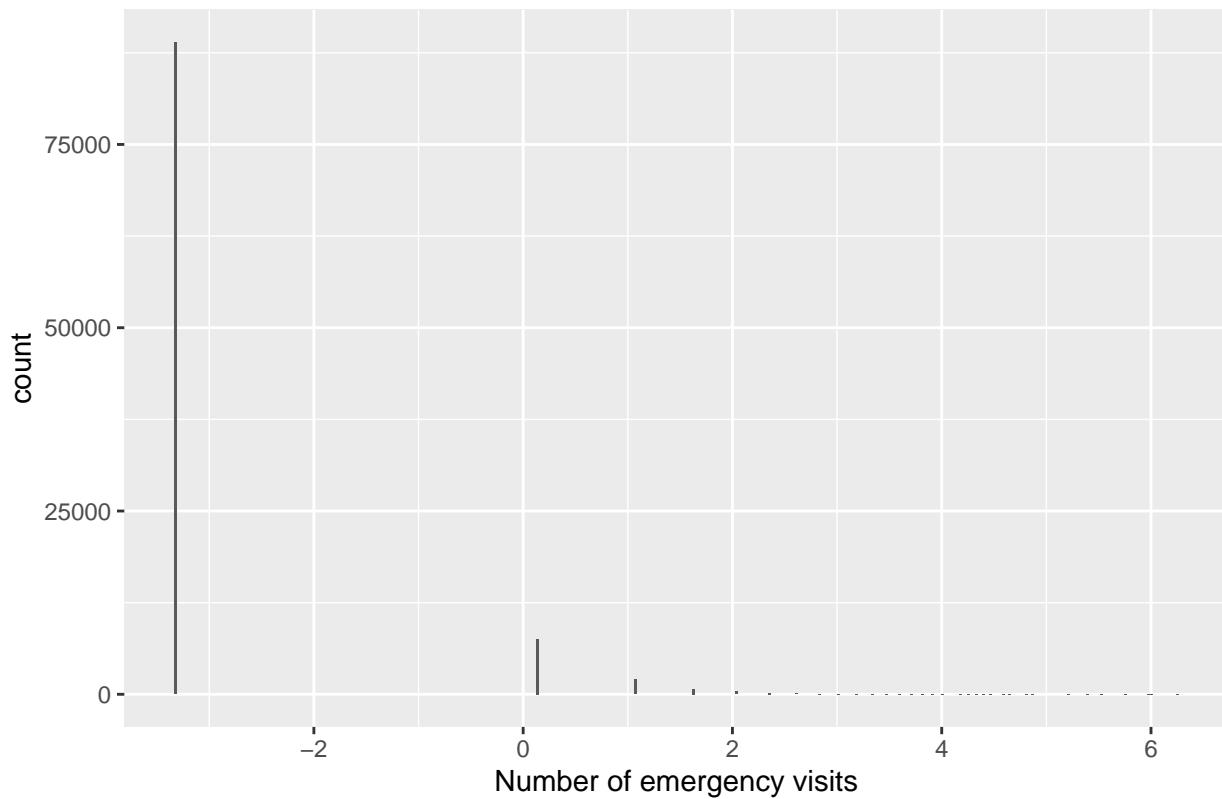
C

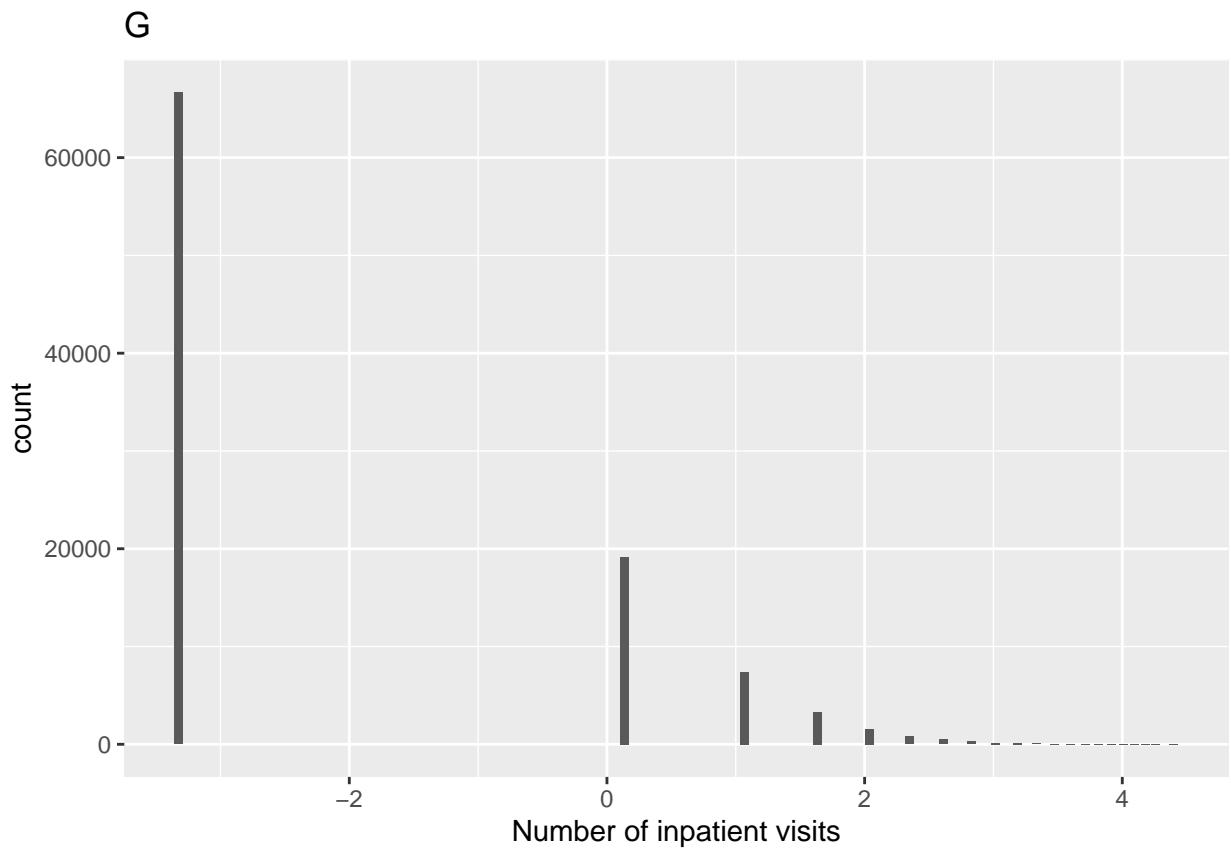


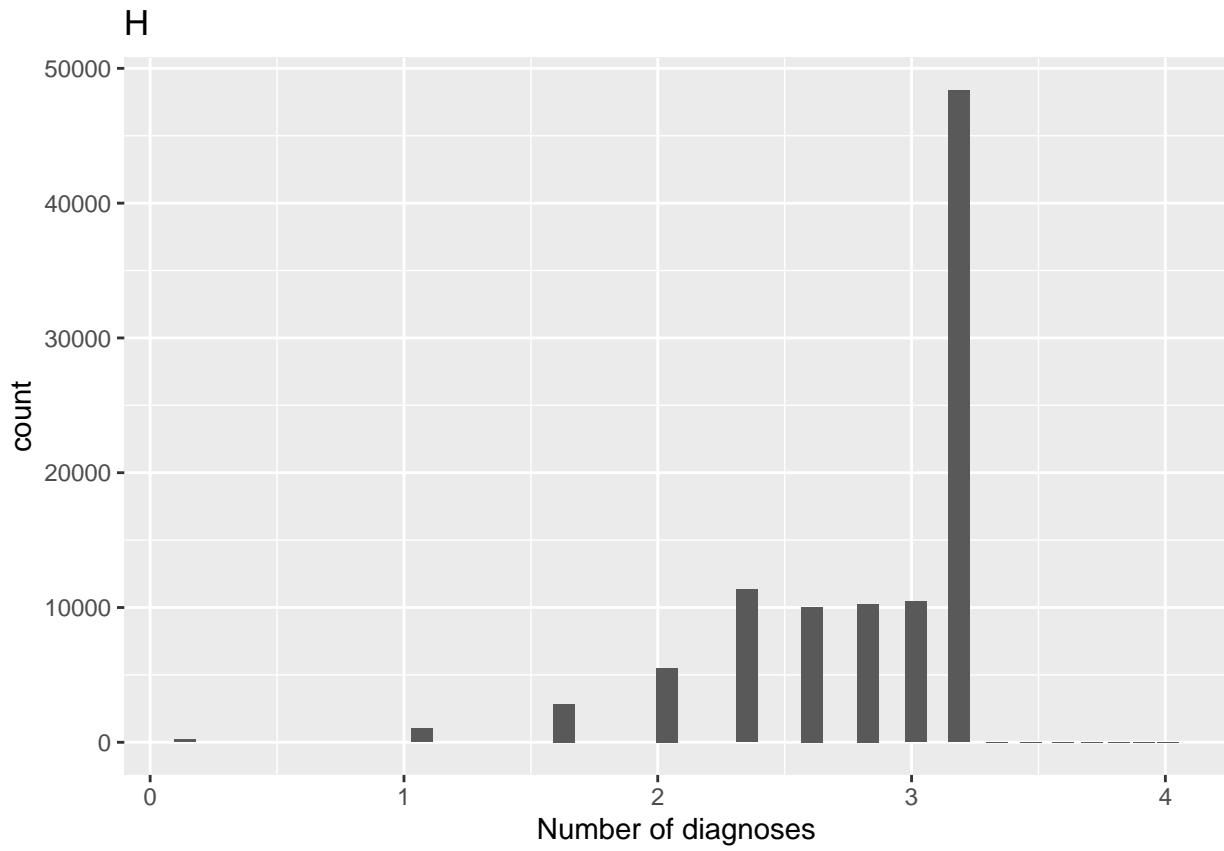




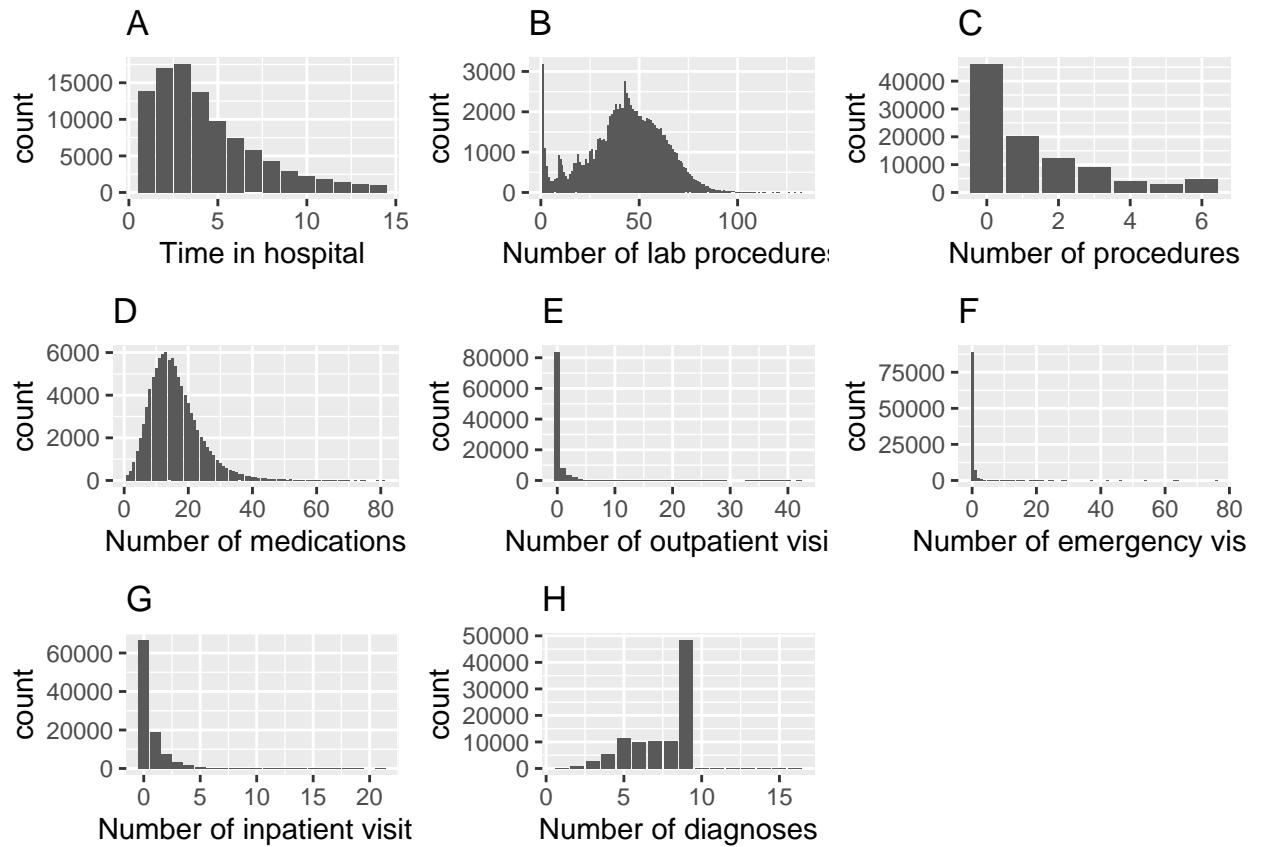
F



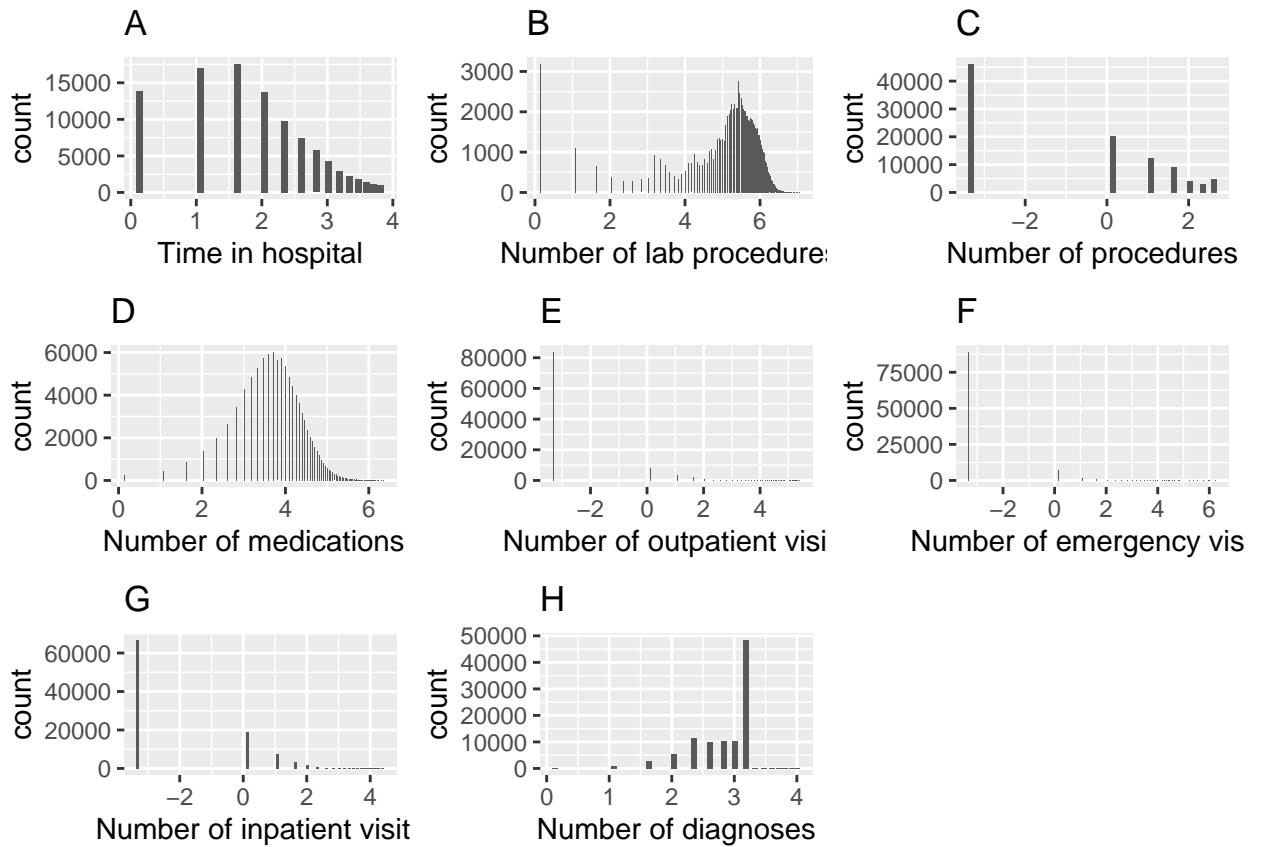




```
# Arrange in two multiplots
do.call(grid.arrange, numeric.normal[1:8])
```



```
ggsave(filename = "figures/distribution_his_1.pdf", width = 10, height = 8, plot = do.call(grid.arrange, numeric.log[1:8]))
```

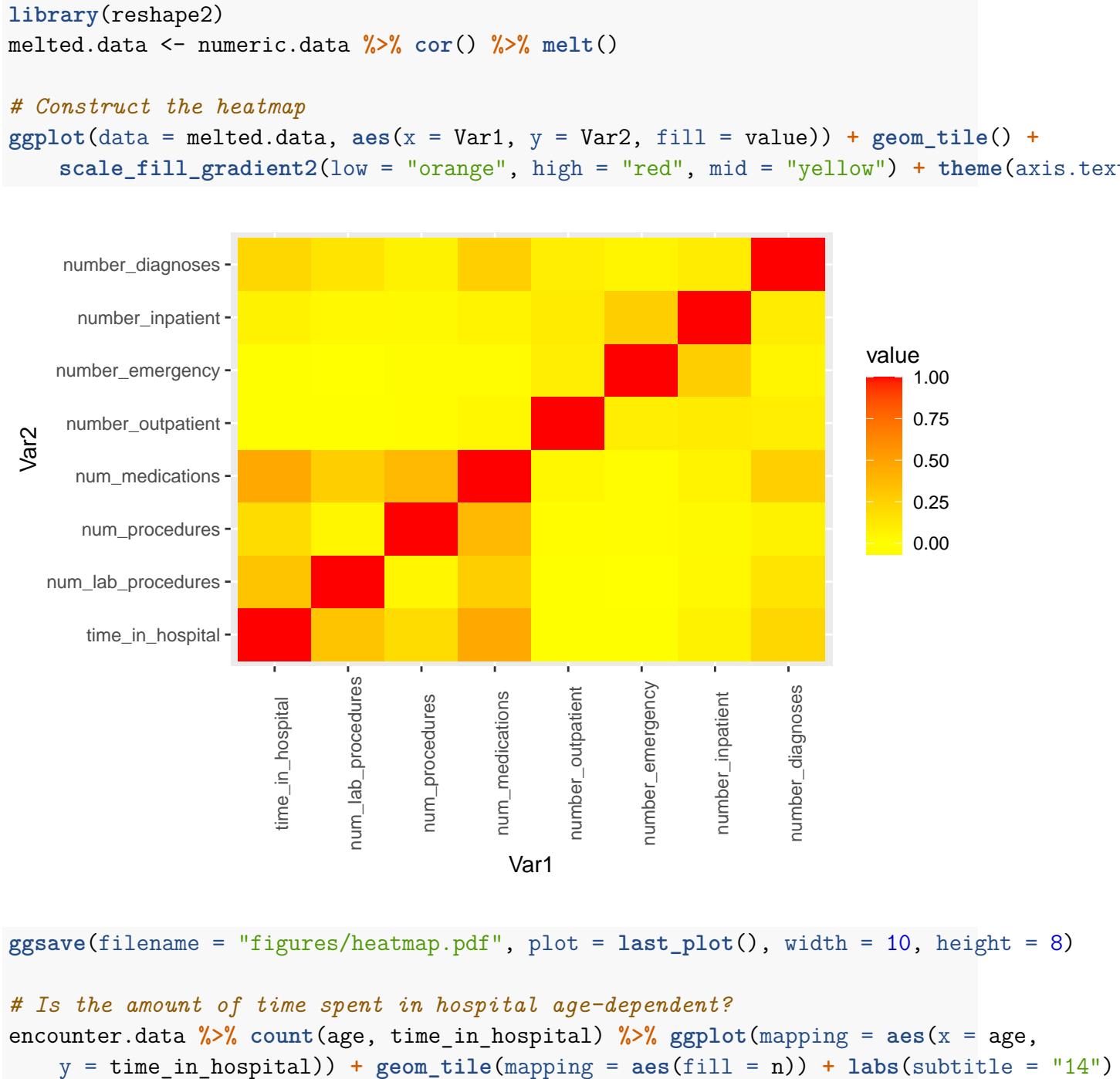


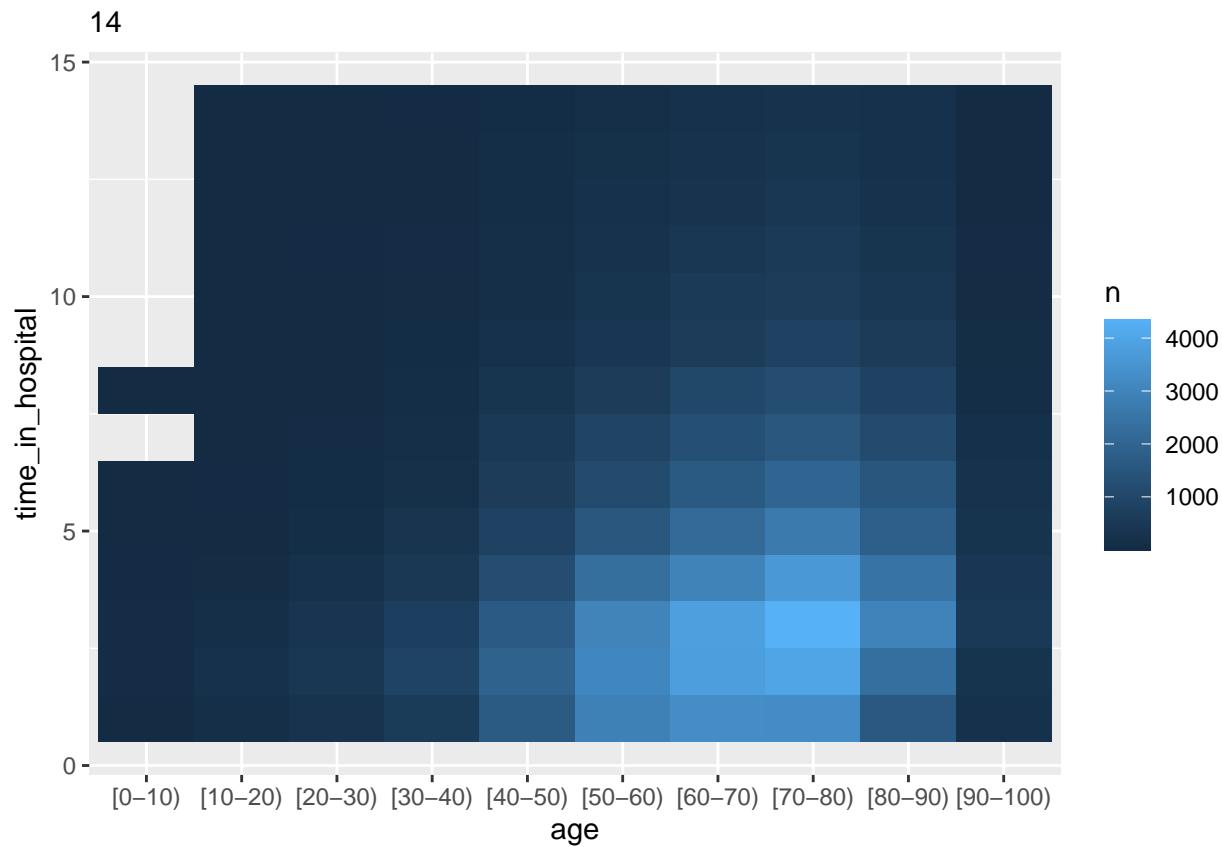
```
ggsave(filename = "figures/distribution_his_2.pdf", width = 10, height = 8, plot = do.call(ggplot, numeric.log[1:8]))
```

Looking at the results, depicted in figures 11 and 12, shows that not all of the attributes are evenly distributed. Normalization might be necessary to correct the distributions of namely B, which makes a couple of weird spikes, E, F, and G, which are all -but B- number of visits of an encounter. We might want to introduce a new attribute, total\_visits, which accounts for all visits made by one encounter. The method of normalization is still up for debate.

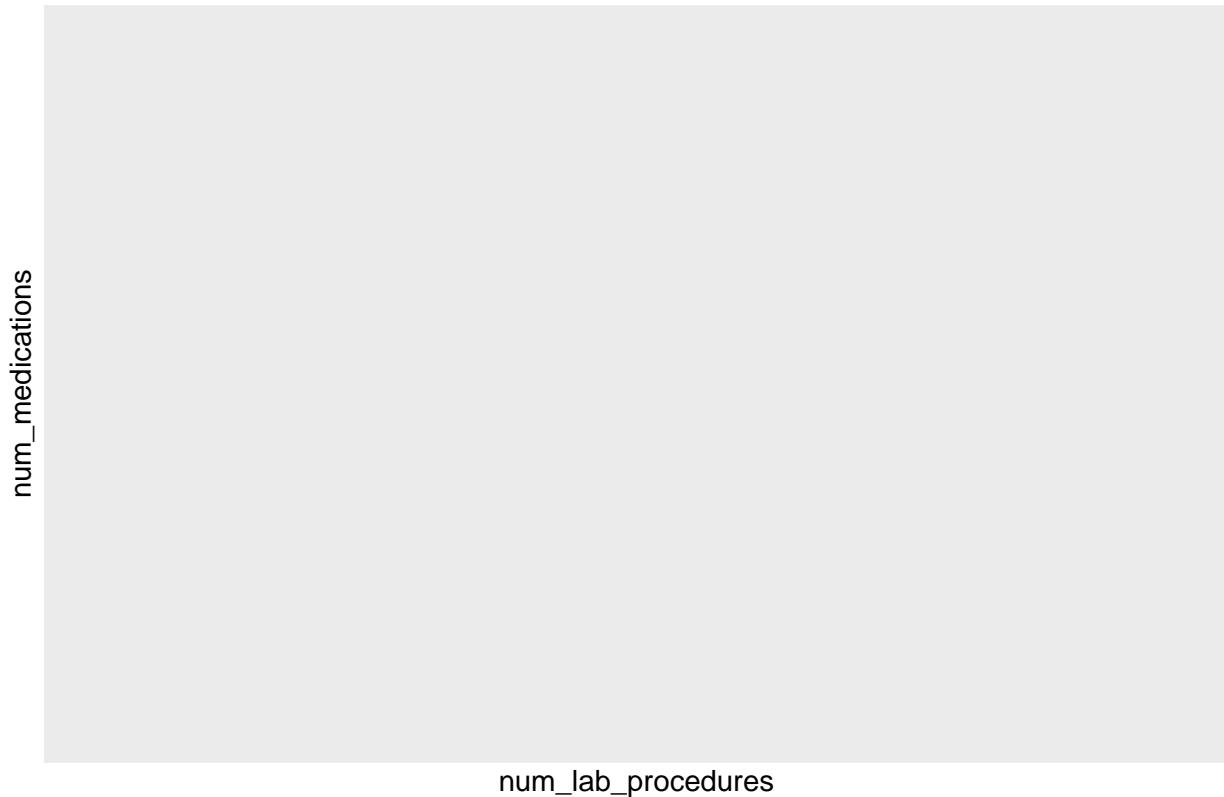
## 2.5 Correlation - numeric attributes

To discover any correlation between the numeric data in our dataset, we constructed a heatmap, which is depicted below. We discuss a possible correlation between age and time spent in the hospital and the potential relationship between the num\_lab\_procedures and num\_medications.





```
# Is there a correlation between num_lab_procedures and num_medications?  
encounter.data %>% ggplot() + geom_hex(mapping = aes(x = num_lab_procedures, y = num_medications), bins = 15, subtitle = "15")
```

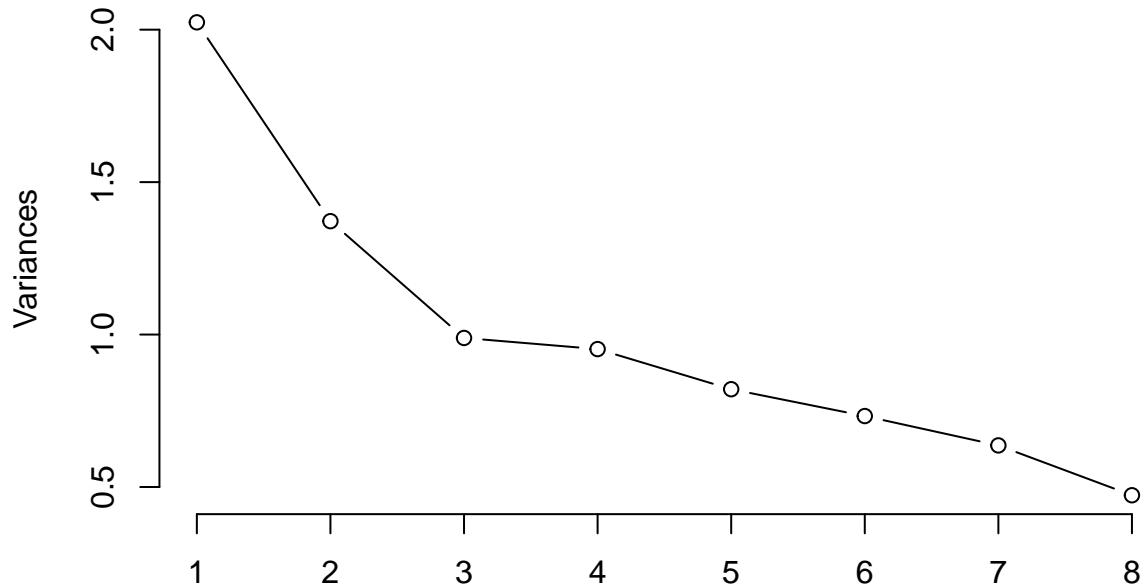


Looking at our heatmap (shown in figure 13), we determine a strong correlation as red, a small correlation as orange, and no correlation as yellow. The attributes that stand out are, for example, num\\_medications-time\\_in\\_hospital, num\\_procedures-num\\_medications, and hba1c\\_res-num\\_medications. As we can see, it is mostly the num\\_medication attribute that has some correlation. Others do not stand out that much. We can conclude that most of the attributes are non-correlated, at least for the numeric attributes.

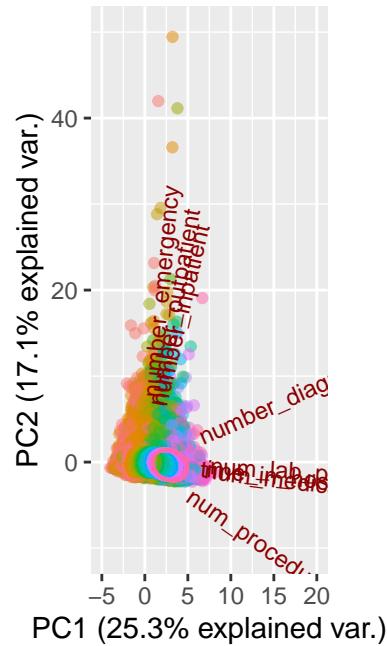
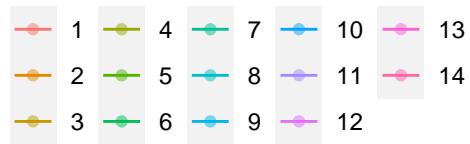
Looking at figure 14, we observe that the age groups above [40-50) increases in total counts; consequently, the total time spent in the hospital before release also increases. We can, therefore, conclude that the time\_in\_hospital attribute is dependent on a patient's age. Finally, we zoom in on two attributes that had a (small) correlation depicted in figure 13. Figure 15 shows this the result of zooming in on num\\_medications and num\\_lab\\_procedures, the result shows that most encounters have a valuation between 15-20 medications and 40-60 amount of lab procedures. The result also shows a lot of outliers, which normalization could solve.

```
cc.pca <- prcomp(numeric.data, center = TRUE, scale. = TRUE)
plot(cc.pca, type = "l", main = "Correlation PCA")
```

## Correlation PCA



```
library(ggbiplot)
ggbiplot(cc.pca, obs.scale = 1, var.scale = 1, groups = as.factor(encounter.data$time_in
  ellipse = TRUE, circle = TRUE, alpha = 0.5) + scale_color_discrete(name = "") +
  theme(legend.direction = "horizontal", legend.position = "top") + xlim(-5, 20) +
  ylim(-10, 50)
```



```
ggsave(filename = "figures/pca-plot.pdf", plot = last_plot(), width = 10, height = 8,
       dpi = 300)
```

### 3 A Clean Dataset

#### 3.1 28th of September

In this section, we will remove all unnecessary attributes and instances as discussed in the EDA. Many filters and mutations have already been performed throughout the EDA, but not all have made it and seems to be scattered around - to give a more clear picture of the final, clean dataset, we perform all filtering, mutations, and relabeling that made it to the final dataset, again.

```
# Load in the data again to perform only filtering that is needed, discussed in
# the EDA
final.data <- read.table(file = "dataset_diabetes/retrieved_data/diabetic_data.csv",
                         sep = ",", header = TRUE)

final.data <- final.data %>% # Filter out every duplicate, based on patient number
distinct(patient_nbr, .keep_all = TRUE) %>% # Relabel '?' to 'Missing'
mutate(race = recode(race, `?` = "Missing")) %>% # Remove instances with a missing label
filter(gender != "Unknown/Invalid" & discharge_disposition_id != 11) %>% # Introduce the gender
mutate(hba1c_res = ifelse((A1Cresult == "None"), "1", NA)) %>% mutate(hba1c_res = ifelse(
  c("Norm", ">7") & is.na(hba1c_res), "2", hba1c_res)) %>% mutate(hba1c_res = ifelse(
  ">8" & (change == "No") & is.na(hba1c_res), "3", hba1c_res)) %>% mutate(hba1c_res =
  ">8" & (change == "Ch") & is.na(hba1c_res), "4", hba1c_res)) %>% # Introduce the new hba1c_res
mutate(admission_type_id = ifelse(admission_type_id %in% c(1, 2, 7), "ICU patient",
                                    "Non-ICU patient")) %>% mutate(admission_source_id = ifelse(admission_source_id %in%
  c(4, 7, 10, 12, 26), "ICU patient", "Non-ICU patient")) %>% mutate(discharge_disposition_id =
  c(13, 14, 19, 20, 21), "ICU patient", "Non-ICU patient")) %>% mutate(icu_or_non =
  ifelse(admission_source_id == discharge_disposition_id, admission_source_id,
         admission_source_id)) %>% # Change the label system in attribute diag_1
mutate(diag_1 = case_when(diag_1 %in% c(390:459) ~ "Circulatory", diag_1 %in% c(460:519) ~
  "Respiratory", diag_1 %in% c(520:579) ~ "Digestive", diag_1 %in% c(240:279) ~
  "Diabetes", diag_1 %in% c(800:999) ~ "Injury", diag_1 %in% c(710:739) ~ "Musculoskeletal",
  diag_1 %in% c(580:629) ~ "Genitourinary", TRUE ~ "Others")) %>% # Set all character
mutate_at(vars("time_in_hospital", "hba1c_res", "icu_or_non"), funs(factor)) %>%
  # Introduce a log2 scale on all numeric attributes
  mutate_at(.vars = names(select_if(., is.numeric)), ~log2(. + 0.1)) %>% # Remove all redundant
  select(-c("encounter_id", "patient_nbr", "weight", "payer_code", "medical_specialty",
          "diag_2", "diag_3", "admission_source_id", "discharge_disposition_id", "admission_type_id"))
# Write final dataset to datafile
write.csv(final.data, "dataset_diabetes/own_data/final_dataset.csv", row.names = FALSE)
```

## 4 Determine quality metrics relevancy

For the next part of our research, we will now predict the time spent in the hospital. Weka, a Java-based data analysis and algorithm interface, is used in this research. Its powerful interface makes it easy to work with comprehensive data modeling techniques.

Weka comes with quite a few ways to evaluate results gathered from a classifier. When predicting a class variable, Weka reports of the number of true positives (TP), true negatives (TN), false positives (FP), and false negatives (FN). These are essential values for evaluating (predicted) results after running on a classifier. Knowing/Having these four crucial numbers, we can calculate quality metrics to determine how well a classifier works on the supplied dataset. The number of quality metrics to choose from is substantial, but we only need a couple to establish whether a classifier is worth-while. We look at the speed a classification is finished, the accuracy, precision, and recall of the results, the rates of true and false positives, F1-score measurement, determine a ROC area, and construct a confusion matrix.

The speed (in seconds) at which a classifier finishes its process is of interest for determining whether a classifier is useful or not. Some classifications may be slower than others. For example, a Naïve Bayes is much quicker than performing a J48 on this dataset due to the complexity of the dataset - 70.000 instances and 44 attributes to consider - and because both classifications have different goals. Therefore, a slower classification is not necessarily worse than a faster one; it depends on what we are interested in researching, and if classifications are similar in end goals. When the latter is true, and there is a considerable difference in construction speed, we can dismiss the slower one. In most cases, other metrics may have the final verdict.

A confusion (or error) matrix represents the performance of an algorithm. It is a table with rows for the predicted instances, and the columns represent the instances of the actual class. The matrix reports on the number of FP, FN, TP, and TN. This representation allows for a more detailed analysis than a metric that specifies a certain value like accuracy. However, it can complicate when too many class variables need to be depicted. For example, our class attribute has 14 valuations and creates a confusion matrix of 196 (14x14) elements. It contains a lot of information; therefore, we introduce other quantity metrics that work more specific and are less complex.

The number of instances correctly identified as either true positive or true negative is called the accuracy. This metric gives a peak in how “true” our classification is. It is essential to know how accurate our classification is as we want the result (prediction) to be as close to the ‘true’ value.

Recall (also called sensitivity or the true positive rate) reports on the number of instances correctly identified as positive out of the total actual positives, and precision (also called positive predictive value) is the number of instances correctly identified as positives of the total instances identified as positive. Both metrics give an insight into the relevance of a classification. The recall is related to the type II error rate and gives information about the completeness. In contrast, precision is more related to the type I error rate, which we

are more interested in as this is about false positives. Having a high precision means good quality, which - mostly - translates to good accuracy. Precision can be a good metric in deciding whether a classification is worth-while or not.

An F1-score is the mean of the precision and recall; it measures the effectiveness of identification done by the recall and precision metrics. It can range from 0 to 1 (perfect precision and recall). In other words, the F1-score conveys the balance between precision and recall. Since we want to balance between precision and recall, the F1-score is a perfect metric to do this with.

The last metric we introduce is the ROC area, which can be depicted in a ROC curve. The ROC curve is created by plotting the recall (or TPR) as a function of the false-positive rate. ROC analysis arranges a way to select optimal models and to discharge suboptimal ones.

$$\text{ACC} = \frac{\text{TP} + \text{TN}}{\text{P} + \text{N}}$$

$$\text{PPV} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

$$\text{TPR} = \frac{\text{TP}}{\text{P}}$$

$$\text{F}_1 = \frac{2\text{TP}}{2\text{TP} + \text{FP} + \text{FN}}$$

Figure 1: All calculation regarding ...

## 5 Machine learning algorithm performances

We will perform a multitude of algorithms. We are interested in performing with Naïve Bayes, Simple Logistics, SVM (SMO), Nearest Neighbour (IBk), a decision tree in the form of a J48/C4.5, and Random Forest. Besides, we also include ZeroR and OneR to measure baseline performance. We only report on the weighted average in evaluating classifiers. Weka gives outcome for every class valuation (so in our case 1 to 14); we take the average sum of these outcomes, as it is better for comparison with other classifiers. The results are depicted in table [NUMBER].

```
library(kableExtra)
normal.result <- read.csv("dataset_weka/normal.csv", sep = ";")
kbl(normal.result, booktabs = T) %>% kable_styling(full_width = F) %>% column_spec(1,
  bold = T)
```

Classifier	TP.Rate	FP.Rate	Accuracy	Precision	Recall	F.Measure	ROC.Area
<b>ZeroR</b>	0,178	0,178	17,82%	?	0,178	?	0,5
<b>OneR</b>	0,21	0,154	20,96%	?	0,21	?	0,528
<b>J48/C4.5</b>	0,198	0,124	19,81%	0,181	0,198	0,187	0,552
<b>Naïve Bayes</b>	0,209	0,118	20,90%	0,193	0,209	0,198	0,648
<b>IBk</b>	0,162	0,118	16,20%	0,161	0,162	0,162	0,522
<b>Simple.Logistics</b>	0,236	0,138	23,55%	0,204	0,236	0,205	0,678
<b>SMO</b>	0	0	0	0	0	0	0
<b>Random.Forest</b>	0	0	0	0	0	0	0

We now observe the results shown in table [NUMBER]. The results of all involved classifiers show a not-so-great picture; the accuracy is very low across all classifiers, ranging from a mere 16,20% for Nearest Neighbor (IBk) to a high-low result of 23,55% with Simple Logistics. This situation is not ideal as - on average - only one-fifth of the class variable is correctly predicted. This is because the class variable is positively skewed, which means low-valuations are more present than higher values. In our case, this means patients spent just a couple of days in hospital rather than, for example, more than a week.

There is no magic trick that would increase the accuracy without severely adjusting our data. Think about removing instances with higher valuation, reconstructing instances, or manually giving weight to high values to create a more normal-looking distribution. We already performed a cost-sensitive classifier to give more weight to values with fewer instances. Awarding even more weight could create a false image of what is really going on. Removing instances with higher values that are less covered is the last thing we want to do. Then we need to draw a line somewhere: what values can stay and which ones need to go? To illustrate the effects of removing high-value instances, we consider removing instances with a valuation ranging from 11 to 14; therefore, we would remove 3567 instances, which translates to a loss of around 6%. After removal, we are left with 66,871 instances. Now that we have a new, temporary dataset, we turn to Weka once more and see if the rate of correctly classified instances has improved. The short-hand answer to that question is yes, it has improved, but not substantially enough to consider removing these instances for our final model. We performed the J48 and Simple Logistics classifiers again and observed an improvement of 1,13% compared to our first run (19,81% versus 20,94%) for J48. We noticed a similar trend for Simple Logistics. Here, the improvement is just 1,33%. The results are depicted in table [NUMBER]. Removing 6% of our instances for such little improvement in accuracy overall is not worthwhile; therefore, we can safely conclude that this method is not what we are looking for in determining a better accuracy.

```
diff.instances <- read.csv("dataset_weka/diff_instances.csv", sep = ";")
kbl(diff.instances, booktabs = T) %>% kable_styling(full_width = F) %>% column_spec(1,
  bold = T)
```

Classifier	TP.Rate	FP.Rate	Accuracy	Precision	Recall	F.Measure	ROC.Area
<b>J48</b>	0,198	0,124	19,81%	0,181	0,198	0,187	0,552
<b>J48 (rstrd)</b>	0,209	0,135	20,94%	0,195	0,209	0,2	0,551
<b>Simple.Logistics</b>	0,236	0,138	23,55%	0,204	0,236	0,205	0,678
<b>Simple.Logistics (rstrd)</b>	0,249	0,146	24,88%	0,226	0,249	0,221	0,669

Another option to improve accuracy could be a relabeling of the attribute. This method has the upside of retaining all instances. The downside of this method is that we create bigger data groups, thereby removing the chance of exactly determining how many days were spent in the hospital. While we cannot determine exactly how many days a patient spent in the hospital, we can separate visits based on classification, such as “short” or “average”. So classification is still possible but is less specific. This is the price to pay if we want a more accurate model. We considered a relabeling with four different valuations:

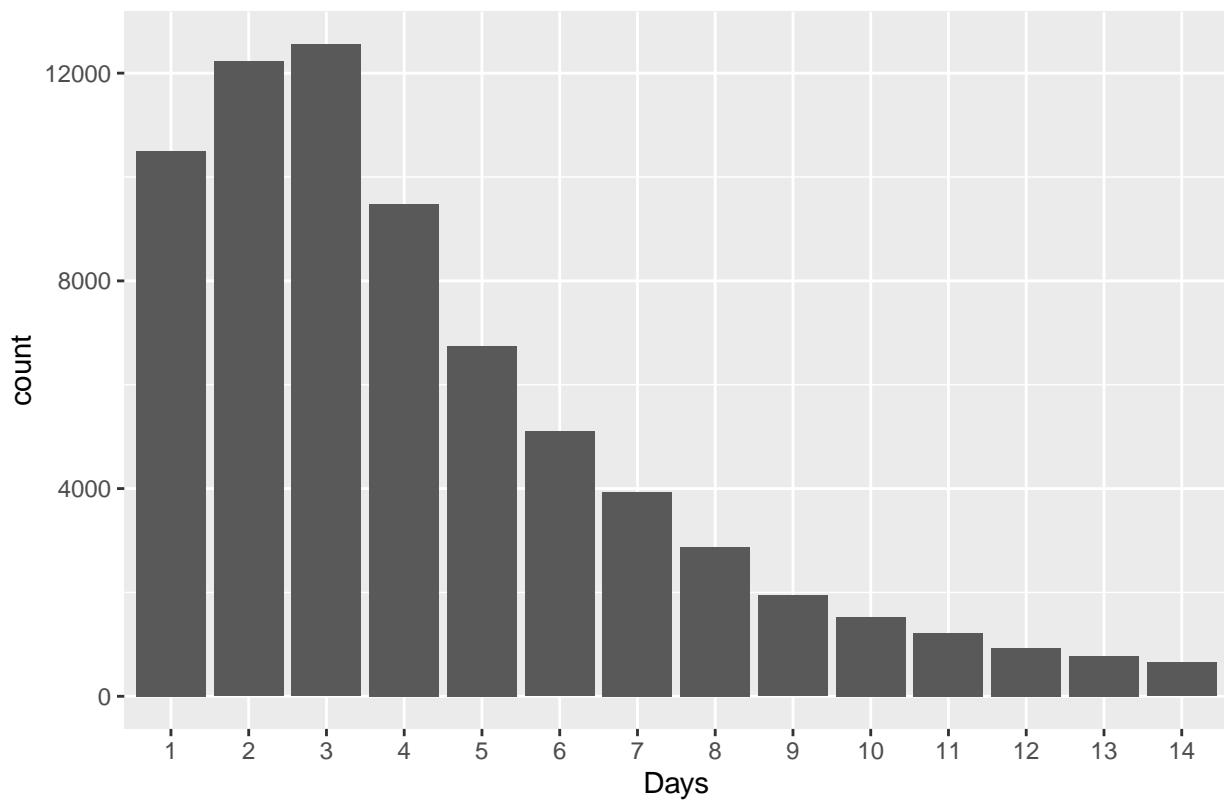
1. Time spent in hospital is between one and three days;
2. Time spent in hospital is between four and six days;
3. Time spent in hospital is between seven and nine days; and
4. Time spent in hospital is equal to or more than ten days.

In the next code section, we mutate the time\_in\_hospital attribute, as described above. We also show the implications of relabeling in figures ... and .... As we can see, the data is still very much positively skewed but has simply fewer classes. Therefore, some classes are relatively big (1 with > 30.000 instances) or small (4 with ~ 5.000 instances). (more?)

```
final.data.label <- final.data %>% mutate(time_in_hospital = case_when(time_in_hospital
  c("1", "2", "3") ~ "1", time_in_hospital %in% c("4", "5", "6") ~ "2", time_in_hospital
  c("7", "8", "9") ~ "3", TRUE ~ "4"))

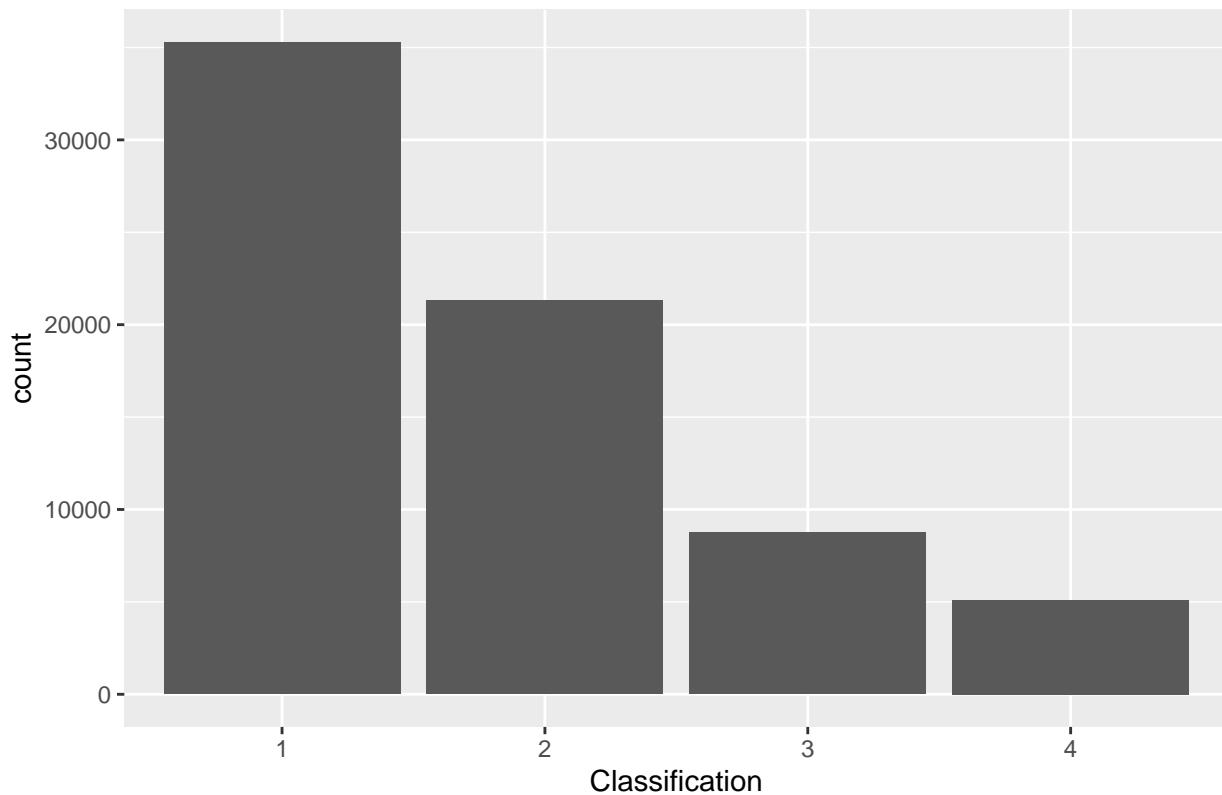
# Add legends and introduce a two-in-one plot
ggplot() + geom_bar(mapping = aes(x = final.data$time_in_hospital)) + ggtitle("Time in Days")
  xlab("Days")
```

Time in hospital (1–14 days)



```
ggplot() + geom_bar(mapping = aes(x = final.data.label$time_in_hospital)) + ggttitle("Ti  
xlab("Classification")
```

### Time in hospital (four classifications)



```
write.csv(final.data.label, "dataset_diabetes/final_dataset_2.csv", row.names = FALSE)
```

Now that we have a new dataset, we turn to Weka once more to analyze the new dataset. The results of this are shown in table [NUMBER]. Having obtained results regarding both datasets, we can determine the differences. As expected, accuracy rose across the broad of classifiers; most experienced an increase of at least 30 percentage points, except for IBk, which rose 26 percentage points. IBk already had the least sufficient outcome in previous Weka runs, so logically we would not have expected it to become best-performing. When considering taking a final dataset to the next stage in our research, we must weigh if a less specific attribute is worth having better accuracy (and for every other metric, for that matter). As we already discussed, having a sense of the duration spent in the hospital is still present. Additionally, we can now predict more accurately by a large margin. Even now, classifiers only predict correctly half of the time, and is even worse when we do not relabel. In contrast, if we come put with a way to produce more sufficient outcomes with the initial dataset, we can always switch back. It seems that the relabeled dataset looks the most promising right now, and therefore we will continue our analysis with this dataset. However, we will report on the initial dataset in the next section as well.

Next up is determining which classifiers are most worthwhile regarding our research goals. As stated earlier, accuracy and precision are the most important metrics. The best option seems to be Simple Logistics with an accuracy of 56,53% and a precision of 0,523 - both are

best-performing in their category. Deciding on a runner-up appears more challenging as J48 and Naïve Bayes both perform well. J48 has better accuracy (52,76%) than Naïve Bayes (51,28%), but Bayes outperforms J48 in every other aspect. This sometimes happens in far bigger margins than their difference in accuracy. For example, the dissimilarity for the ROC area is 0,071. To conclude, both classifiers seem as likely candidates for further analysis, and it is hard to decide which is the best. For that reason, in addition to taking Simple Logistics, both Naïve Bayes and J48 will be considered in further analysis.

```
label.result <- read.csv("dataset_weka/labelling.csv", sep = ";")
kbl(label.result, booktabs = T) %>% kable_styling(full_width = F) %>% column_spec(1,
  bold = T)
```

Classifier	TP.Rate	FP.Rate	Accuracy	Precision	Recall	F.Measure	ROC.Area
<b>ZeroR</b>	0,501	0,501	50,08%	?	0,501	?	0,5
<b>OneR</b>	0,532	0,365	53,15%	0,467	0,532	0,468	0,583
<b>J48/C4.5</b>	0,528	0,305	52,76%	0,489	0,528	0,501	0,624
<b>Naïve Bayes</b>	0,513	0,271	51,28%	0,505	0,513	0,506	0,695
<b>IBk</b>	0,426	0,321	42,60%	0,423	0,426	0,425	0,552
<b>Simple.Logistics</b>	0,565	0,314	56,53%	0,523	0,565	0,52	0,732
<b>SMO</b>	0	0	0	0	0	0	0
<b>Random.Forest</b>	0	0	0	0	0	0	0

## **6 References**

### **References**

- [1] Beata Strack, Jonathan P. DeShazo, Chris Gennings, Juan L. Olmo, Sebastian Ventura, Krzysztof J. Cios, and John N. Clore, "Impact of HbA1c Measurement on Hospital Readmission Rates: Analysis of 70,000 Clinical Database Patient Records," BioMed Research International, vol. 2014, Article ID 781670, 11 pages, 2014. DOI: <https://doi.org/10.1155/2014/781670>
- [2] Centers for Disease Control and Prevention, National Center for Health Statistics, ICD-9, <https://www.cdc.gov/nchs/icd/icd9.htm>, November 6, 2015.