

Predicting Diabetes in PIMA Women

edX Capstone Project Submission

Kirtimay Pendse

6/23/2020

Introduction

Diabetes is a metabolic disorder defined as when one's blood glucose is too high (known as hyperglycemia) for a prolonged period of time. Glucose is an essential simple sugar widely consumed daily, and the hormone insulin helps absorbing glucose from food and transform it into energy; however, sometimes one's body doesn't make enough insulin or is unable to use it well, resulting in glucose staying in the blood stream undigested and unable to reach the cells.¹. This can cause health problems, especially diabetes. Around 9.5% -almost 30.5 million- of the United States population had diabetes in 2015 ², and factors such as being overweight, being physically inactive, having a family history are linked with higher chances of developing diabetes. Due to several factors not discussed in this paper ³, diabetes is extremely prevalent in Native Americans, most notably within the Pima tribe- since the Pima tribe is a mostly homogenous group, Pima people have been the subject of several studies of diabetes.

This project is the final part of the HarvardX: PH125.9x Data Science: Capstone course⁴, the last course for the Data Science Professional Certificate. This project is centered around predicting the presence of diabetes in Pima Indian women using data on factors such as age, body mass index, blood pressure etc. compiled together in the Pima Indians Diabetes dataset.

The dataset, loaded as 'pima_diabetes', is split into a training set containing 80% of the data and a test set containing 20% of the data for validation. This report is split into four sections: first, the objective and motivation behind the project is highlighted, then exploratory data analysis is conducted, following which the modeling approach to develop the diabetes prediction algorithm is presented. Finally, the modeling results are presented along with a discussion on the algorithm's performance and its limitations.

Objective

The dataset⁵ is available on Kaggle and is originally sourced from the National Institute of Diabetes and Digestive and Kidney Diseases, a part of the Department of Health and Human Services. The objective of this analysis is to diagnostically predict whether or not a patient is diabetic, based on select diagnostic measurements included in the dataset (such as BMI, Age, Blood Pressure). There are 786 individuals in the dataset, all of whom are females of at least 21 years of age, and of Pima Indian heritage.

¹<https://www.niddk.nih.gov/health-information/diabetes/overview/what-is-diabetes#:~:text=Diabetes%20is%20a%20disease%20that,to%20be%20controlled%20with%20insulin,Diabetes%20is%20a%20disease%20that,to%20be%20controlled%20with%20insulin>

²Centers for Disease Control and Prevention. National diabetes statistics report, 2017. www.cdc.gov/diabetes/pdfs/data/statistics/national-diabetes-statistics-report.pdf

³more can be found at <https://care.diabetesjournals.org/content/29/8/1866>

⁴<https://courses.edx.org/courses/course-v1:HarvardX+PH125.9x+1T2020/course/>

⁵<https://www.kaggle.com/ksp585/pima-indian-diabetes-logistic-regression-with-r>

Methods and Analysis

Preparing the data

First, the dataset is downloaded and split into a train set and a test set. The train set is used to create the prediction algorithm, and then the algorithm is tested on the test set for a final validation.

```
#Loading required packages
library(lubridate)
if(!require(ggthemes))
  install.packages("ggthemes", repos = "http://cran.us.r-project.org")
if(!require(scales))
  install.packages("scales", repos = "http://cran.us.r-project.org")
if(!require(tidyverse)) install.packages("tidyverse", repos = "http://cran.us.r-project.org")
if(!require(caret)) install.packages("caret", repos = "http://cran.us.r-project.org")
if(!require(data.table)) install.packages("data.table", repos = "http://cran.us.r-project.org")
library(dplyr)
library(knitr)
library(ggplot2)
library(dslabs)
library(lubridate)
library(corrplot)
library(readr)

#Downloading the data
dl <- tempfile()
download.file("https://github.com/kirtimay/edX_Capstone/blob/master/cyo-diabetes/diabetes.csv", dl)
pima_diabetes <- read.csv("diabetes.csv", col.names=c("Pregnancies","Glucose","BP","SkinThickness","Insulin",
head(pima_diabetes) #check whether dataset downloaded properly
```

```
##   Pregnancies Glucose BP SkinThickness Insulin BMI   DPF Age Outcome
## 1           6    148 72           35      0 33.6 0.627  50         1
## 2           1     85 66           29      0 26.6 0.351  31         0
## 3           8    183 64            0      0 23.3 0.672  32         1
## 4           1     89 66           23     94 28.1 0.167  21         0
## 5           0    137 40           35    168 43.1 2.288  33         1
## 6           5    116 74            0      0 25.6 0.201  30         0
```

```
#convert outcome to factor
pima_diabetes$Outcome <- factor(pima_diabetes$Outcome)
```

Description of Variables:

```
v_type <- lapply(pima_diabetes, class)
v_desc <- c("No. of Pregnancies", "Plasma Glucose Concentration (2 Hrs after an oral test)", "Diastolic
v_name <- colnames(pima_diabetes)
desc_table <- as_data_frame(cbind(v_name, v_type, v_desc))
colnames(desc_table) <- c("Variable","Class","Description")
desc_table #%>% knitr::kable()
```

```
## # A tibble: 9 x 3
##   Variable Class      Description
##   <list>   <list>    <list>
## 1 <chr [1]> <chr [1]> <chr [1]>
## 2 <chr [1]> <chr [1]> <chr [1]>
## 3 <chr [1]> <chr [1]> <chr [1]>
## 4 <chr [1]> <chr [1]> <chr [1]>
```

```
## 5 <chr [1]> <chr [1]> <chr [1]>
## 6 <chr [1]> <chr [1]> <chr [1]>
## 7 <chr [1]> <chr [1]> <chr [1]>
## 8 <chr [1]> <chr [1]> <chr [1]>
## 9 <chr [1]> <chr [1]> <chr [1]>

set.seed(1, sample.kind="Rounding")
test_index <- createDataPartition(y = pima_diabetes$Outcome, times = 1, p = 0.2, list = FALSE)
train_set <- pima_diabetes[-test_index,]
test_set <- pima_diabetes[test_index,]
```

Exploratory Analysis

text

```
## Pregnancies Glucose BP SkinThickness Insulin BMI DPF Age Outcome
## 1          6    148 72         35         0 33.6 0.627 50         1
## 2          1     85 66         29         0 26.6 0.351 31         0
## 3          8    183 64          0         0 23.3 0.672 32         1
## 4          1     89 66         23        94 28.1 0.167 21         0
## 5          0    137 40         35       168 43.1 2.288 33         1
## 6          5    116 74          0         0 25.6 0.201 30         0
```

text

```
## Pregnancies      Glucose      BP      SkinThickness      Insulin
## Min.   : 0.00   Min.   : 0   Min.   : 0.0   Min.   : 0.0   Min.   : 0.0
## 1st Qu.: 1.00   1st Qu.: 99   1st Qu.: 62.0   1st Qu.: 0.0   1st Qu.: 0.0
## Median : 3.00   Median :117   Median : 72.0   Median :23.0   Median : 30.5
## Mean   : 3.85   Mean   :121   Mean   : 69.1   Mean   :20.5   Mean   : 79.8
## 3rd Qu.: 6.00   3rd Qu.:140   3rd Qu.: 80.0   3rd Qu.:32.0   3rd Qu.:127.2
## Max.   :17.00   Max.   :199   Max.   :122.0   Max.   :99.0   Max.   :846.0
## BMI      DPF      Age      Outcome
## Min.   : 0.0   Min.   :0.078   Min.   :21.0   0:500
## 1st Qu.:27.3   1st Qu.:0.244   1st Qu.:24.0   1:268
## Median :32.0   Median :0.372   Median :29.0
## Mean   :32.0   Mean   :0.472   Mean   :33.2
## 3rd Qu.:36.6   3rd Qu.:0.626   3rd Qu.:41.0
## Max.   :67.1   Max.   :2.420   Max.   :81.0
```

text

```
## Pregnancies      Glucose      BP SkinThickness      Insulin
##           0           0           0           0           0
## BMI      DPF      Age      Outcome
##           0           0           0           0
```

Plots

Outcome Variable

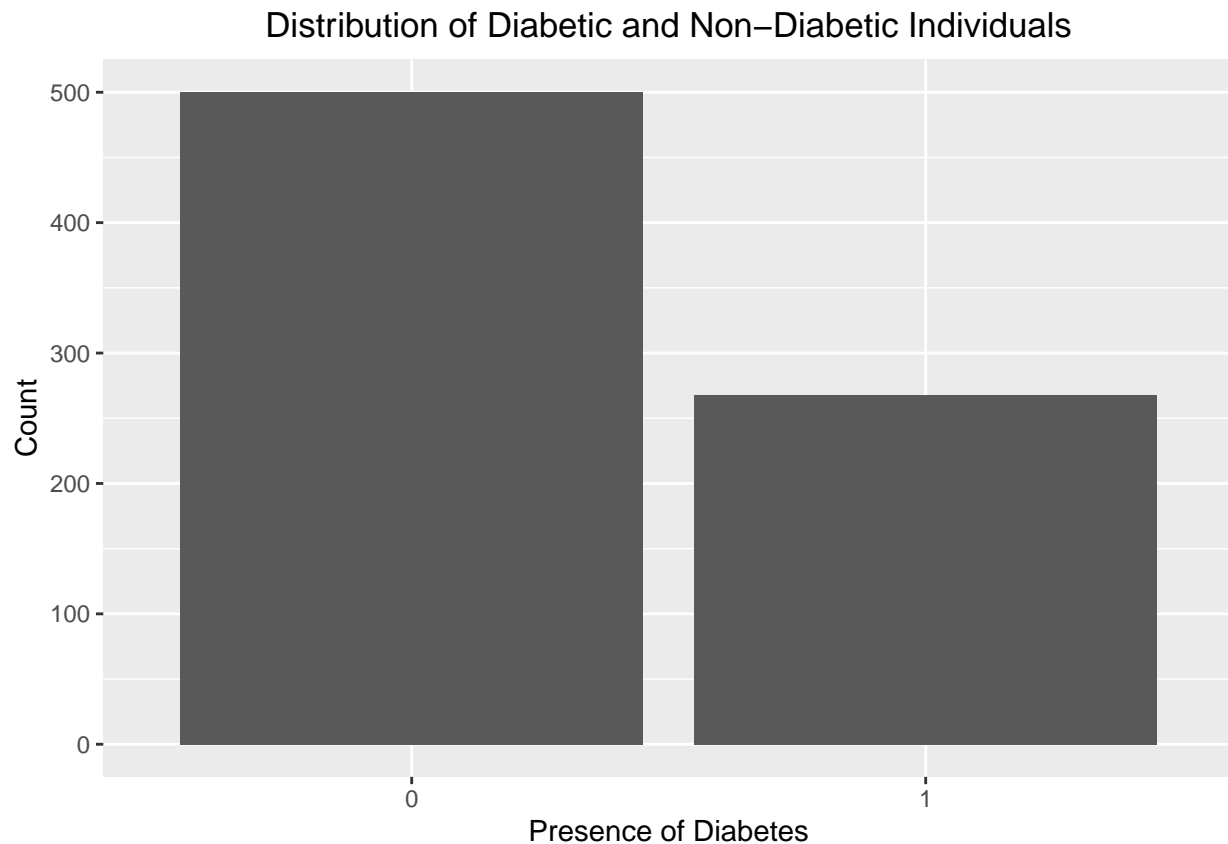
Diabetes

```
ggplot(pima_diabetes, aes(Outcome)) +
  geom_bar() +
  ggtitle("Distribution of Diabetic and Non-Diabetic Individuals") +
  theme(plot.title = element_text(hjust = 0.5)) +
```

```

xlab("Presence of Diabetes") +
ylab("Count")

```



Predictor Variables

Pregnancies

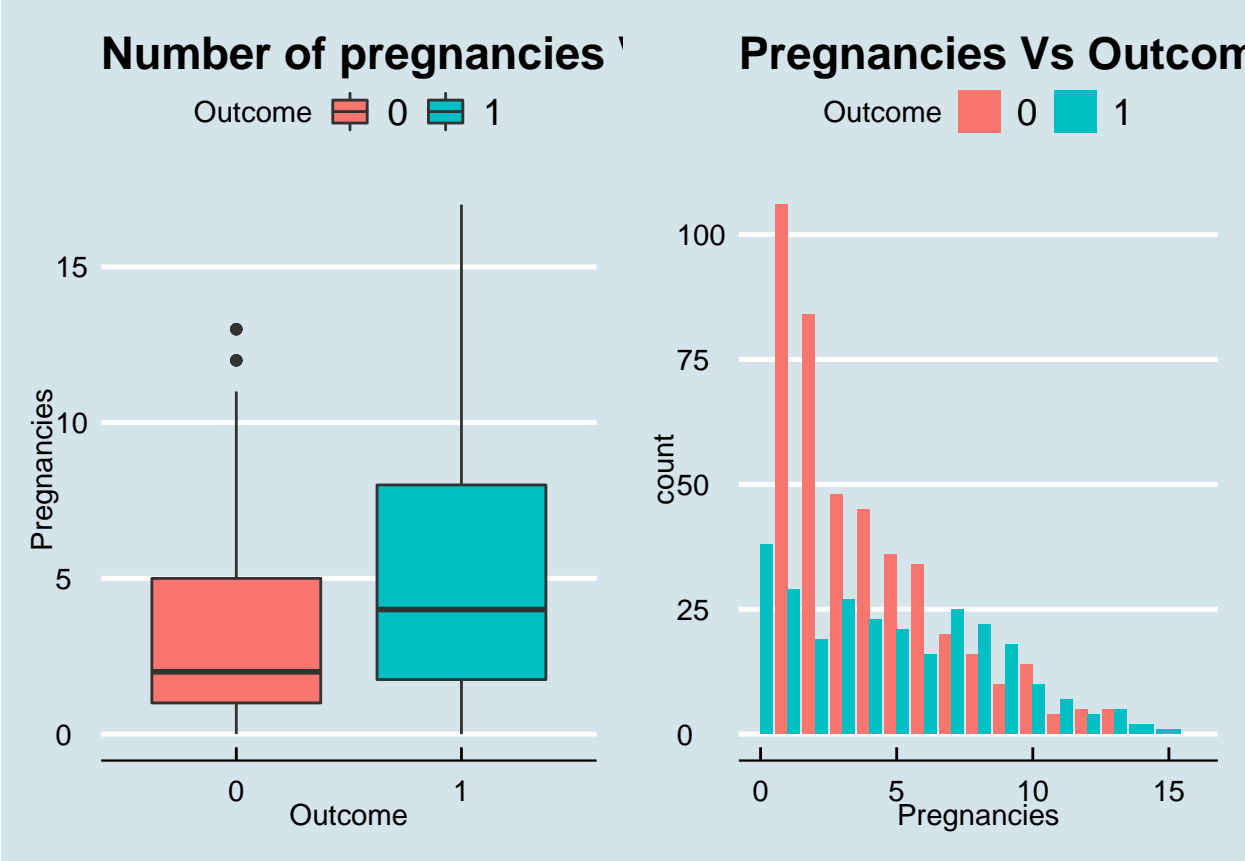
```

p1 <- ggplot(pima_diabetes, aes(x = Outcome, y = Pregnancies, fill=Outcome)) +
  geom_boxplot() +
  theme(legend.position = "bottom") +
  ggtitle("Number of pregnancies Vs Diabetes") +
  theme_economist()

p2 <- ggplot(pima_diabetes, aes(x = Pregnancies, fill=Outcome)) +
  geom_bar(position = "Dodge") +
  scale_x_continuous(limits = c(0,16)) +
  theme(legend.position = "bottom") +
  labs(title = "Pregnancies Vs Outcome")+
  theme(legend.position="right") +
  theme_economist()

gridExtra::grid.arrange(p1, p2, ncol = 2)

```



Glucose

Blood Pressure

Skin Thickness

Insulin

BMI

DPF

Age

Correlation Matrix

##	Pregnancies	Glucose	BP	SkinThickness	Insulin	BMI	DPF	Age
## Pregnancies	1.0	0.1	0.1	-0.1	-0.1	0.0	0.0	0.5
## Glucose	0.1	1.0	0.2	0.1	0.3	0.2	0.1	0.3
## BP	0.1	0.2	1.0	0.2	0.1	0.3	0.0	0.2
## SkinThickness	-0.1	0.1	0.2	1.0	0.4	0.4	0.2	-0.1
## Insulin	-0.1	0.3	0.1	0.4	1.0	0.2	0.2	0.0
## BMI	0.0	0.2	0.3	0.4	0.2	1.0	0.1	0.0
## DPF	0.0	0.1	0.0	0.2	0.2	0.1	1.0	0.0
## Age	0.5	0.3	0.2	-0.1	0.0	0.0	0.0	1.0

