

MACHINE LEARNING FOR PREDICTING DIABETES MELLITUS

A project report submitted in partial fulfilment of the
requirement for the degree of

Bachelor in Science

In

Computer Science

by

Alexander Roque Rodrigues

under the supervision of

Ashweta Fondekar

**PARVATIBAI CHOWGULE COLLEGE OF ARTS
& SCIENCE AUTONOMOUS**

June 2019

Declaration by Candidate

I declare that this project report has been prepared by me and to the best of my knowledge, it has not previously formed the basis for the award of any diploma or degree by any other University.

Certificate by Supervisor

Certified that the Project Report is a record of work done by the candidate himself/herself/themselves under my guidance during the period of study and that to the best of my knowledge, it has not previously formed the basis of the award of any degree or diploma of any other University.

Ashweta Fondekar
Project Supervisor

Work Record

1 Acknowledgements

Contents

1	Acknowledgements	4
2	Introduction	10
3	Project Rationale	12
4	Software Requirements	14
4.1	Functionalities for Doctors	14
4.2	Functionalities for Patients	14
4.3	Functionalities for Master Nodes	14
4.4	Functionalities for Slave Nodes	14
5	Hardware Requirements	15
5.1	Raspberry Pi	15
5.2	Switch	15
5.3	Router	15
5.4	Master Node	15
5.5	Slave Node	16
5.6	Database	16
5.7	Web Server	16
6	Technology Stack	17
6.1	Python	17
6.1.1	Pandas	17
6.1.2	Sklearn	17
6.1.3	Numpy	17
6.1.4	Itertools	17
6.2	MYSQL	17
6.3	Apache Web Server	18
6.4	PHP	18
6.5	AJAX	18
6.6	GitHub	19
6.7	SSH	19
7	Methodology	20
8	Deploying to Cluster	21
8.1	Benefits of Clustering	21
9	Future Scope	23

10 Tables	24
11 Appendix I	31
11.1 K Nearest Neighbours	31
11.2 Logistic Regression	31
11.3 Decision Tree	31
11.4 Random Forest Classifier	31
11.5 Gradient Boosting	31
11.6 Multi Layered Perceptron	32
11.7 SVM	32
13 Appendix II	33
14 Conclusion	59

List of Figures

1	Diabetic v/s Healthy Subject Count.	33
2	Subject distribution across Body Mass Index range.	34
3	Subject distribution across age.	34
4	Subjects Insulin Distribution across ranges.	35
5	Subjects distribution via Pregnancies.	35
6	Subjects distribution via Diabetes Pedigree Function.	36
7	Plot of N0	36
8	Diabetic v/s Healthy Subjects Percentage.	37
9	Boxplot for all attributes with outliers.	38
10	Heatmap using Pearsons Correlation Coefficient.	39
11	Number of missing values in count and percentage.	39
12	Glucose v/s Age scatterplot.	40
13	Subjects glucose distribution.	40
14	Skin Thickness distribution of subjects.	41
15	Blood Pressure distribution of subjects.	41
16	Subjects Insulin distributon.	42
17	New feature N1.	42
18	New feature N3.	43
19	New feature N4.	43
20	New feature N6.	44
21	New feature N7.	44
22	N1 barplot for diabetic and healthy population.	45
23	N1 distribution in percentage.	45
24	N2 barplot for diabetic and healthy population.	46
25	N2 distribution in percentage.	46
26	Pregnancies v/s age scatterplot.	47
27	N3 barplot for diabetic and healthy population.	47
28	N3 distribution in percentage.	48
29	Glucose v/s Blood Pressure scatterplot.	48
30	N4 barplot for diabetic and healthy population.	49
31	N4 distribution by target.	49
32	N5 barplot for diabetic and healthy population.	50
33	N5 distribution by target.	50
34	Skin Thickness v/s BMI scatterplot.	51
35	N6 barplot for diabetic and healthy population.	51
36	N6 distribution by target.	52
37	Glucose v/s Body Mass Index scatterplot.	52

38	N7 barplot for diabetic and healthy population.	53
39	N7 distribution by target.	53
40	N9 barplot for diabetic and healthy population.	54
41	N9 distribution by target.	54
42	N10 barplot for diabetic and healthy population.	55
43	N10 distribution by target.	55
44	N11 barplot for diabetic and healthy population.	56
45	N11 distribution by target.	56
46	N15 barplot for diabetic and healthy population.	57
47	N15 distribution by target.	57
48	Extended heatmap with new features combined.	58

List of Tables

1	Pima Indians dataset header.	24
2	Train accuracy using various machine learning algorithms with various parameters.	25
3	Test accuracy using various machine learning algorithms with various parameters.	26
4	Null Value Check.	27
5	First 15 rows of the raw dataset.	28
6	Correlation Plot.	29
7	MLP 5 fold verification.	30
8	Insulin Median Comparison.	30
9	Glucose Median Comparison.	30
10	Skin Thickness Median Comparison.	30
11	Blood Pressure Median Comparison.	30
12	Body Mass Index Median Comparison.	30

Abstract

Diabetes Mellitus is a disease that prevents the body from properly expanding the energy stored from the food consumed. The purpose of this project was to select machine learning algorithms that are able to predict or classify a person as diabetic or healthy based on the legacy data. The algorithms compared were KNN Classifier, Logistic Regression, Decision Tree Classifier, Random Forest Classifier, Gradient Boosting Classifier, Support Vector Classifier and the Multi-Layered Perceptron. From all the above the Multi-Layered Perceptron gave an accuracy of prediction on the dataset as 79.70%. To improve the performance of the classifier, I have considered new features deduced from the currently existing feature set and re-trained the classifier on the new dataset that I generated, which is now able to classify the subject as diabetic or healthy with a new accuracy of 93.10%. A significant change which boosted the accuracy by 13.4%. After selection of the algorithm I further advanced the platform of cluster computing to deploy the algorithm onto and generate predictions in without any human interference (apart from entering the data itself) and also made the data available to the users via an easy to use web application which gives them access to the observations then stored in the database after being predicted by the algorithm deployed on the nodes.

2 Introduction

In recent days, there has been a sharp increase in the cases of diabetes mellitus. Diabetes mellitus is on the rise amongst many people and the rate of contracting this lifestyle disease could be reduced significantly if proper measures and precautions were to be instilled amongst people the number of people can be reduced.

Machine learning is a growing field in computer science. With the development and introduction of many algorithms the prediction and accuracy of the predictions itself has improved substantially. Machine learning and healthcare systems are also becoming increasingly popular in the healthcare sector.

The project encompasses the qualities of Remote Patient Monitoring (RPM) and Clinical Decision Support (CDS). RPM provides medical facilities that have the ability to transmit patient data to healthcare professionals who might very well be halfway around the world. RPM can monitor blood glucose levels and blood pressure. It is particularly helpful for patients with chronic conditions such as type 2 diabetes, hypertension, or cardiac disease. Data collected and transmitted via PRM can be used by a healthcare professional or a healthcare team to detect medical events such as stroke or heart attack that require immediate and aggressive medical intervention. Data collected may be used as part of a research project or health study. RPM is a life-saving system for patients in remote areas who cannot access face-to-face health

care. CDS analyzes data from clinical and administrative systems. The aim is to assist healthcare providers in making informed clinical decisions. Data available can provide information to medical professions who are preparing diagnoses or predicting medical conditions like drug interactions and reactions. CDS tools filter information to assist healthcare professionals in caring for individual clients.

Cluster computing or High-Performance computing frameworks is a form of computing in which bunch of computers (often called nodes) that are connected through a LAN (local area network) so that, they behave like a single machine. A computer cluster help to solve complex operations more efficiently with much faster processing speed, better data integrity than a single computer and they only used for mission-critical applications.

The objective of this project is to create a system that is able to use the machine learning algorithms and predict the outcome of the parameters entered into the algorithm and help the patient draw a conclusion whether or not he/she has the same traits exhibited by similar patients that have diabetes. Also the system should have a UI that is capable of displaying the data of the patients to the doctor and to the patients themselves for further interpretation.

3 Project Rationale

Diabetes is a chronic disease that occurs either when the pancreas does not produce enough insulin or when the body cannot effectively use the insulin it produces. Insulin is a hormone that regulates blood sugar. Hyperglycaemia, or raised blood sugar, is a common effect of uncontrolled diabetes and over time leads to serious damage to many of the body's systems, especially the nerves and blood vessels.

Taking a look at the number of people with diabetes, the count has risen from 108 million in 1980 to 422 million in 2014. The global prevalence of diabetes among adults over 18 years of age has risen from 4.7% in 1980 to 8.5% in 2014. Diabetes prevalence has been rising more rapidly in middle-income and low-income countries. Diabetes as a disease is a major cause of blindness, kidney failure, heart attacks, stroke and lower limb amputation amongst those who recognise the fatal disease late in life. In 2016, an estimated 1.6 million deaths were directly caused by diabetes. Another 2.2 million deaths were attributable to high blood glucose in 2012. Almost half of all deaths attributable to high blood glucose occur before the age of 70 years. WHO estimates that diabetes was the seventh leading cause of death in 2016. A healthy diet, regular physical activity, maintaining a normal body weight and avoiding tobacco use are ways to prevent or delay the onset of type 2 diabetes. Diabetes can be treated and its consequences avoided or delayed with diet, physical activity, medication and regular screening and treatment for complications. In 2014, 8.5% of adults aged 18 years and older had diabetes. In 2016, diabetes was the direct cause of 1.6 million deaths and in 2012 high blood glucose was the cause of another 2.2 million deaths.

Diabetes patients are segregated into types. Type 1 diabetes is characterized by deficient insulin production and requires daily administration of insulin. The cause of type 1 diabetes is not known and it is not preventable with current knowledge. A person with type 1 diabetes may have symptoms like excessive excretion of urine, thirst, constant hunger, weight loss, vision changes, and fatigue. These symptoms may occur suddenly.

Type 2 diabetes results from the body's ineffective use of insulin. Type 2 diabetes comprises the majority of people with diabetes around the world, and is largely the result of excess body weight and physical inactivity. Symptoms may be similar to those of type 1 diabetes, but are often less marked. As a result, the disease may be diagnosed several years after onset, once complications have already arisen.

Until recently, this type of diabetes was seen only in adults but it is now also occurring increasingly frequently in children.

Gestational diabetes is another type of diabetes where hyperglycaemia with blood glucose values above normal but below those diagnostic of diabetes, occurring during pregnancy. Women with gestational diabetes are at an increased risk of complications

during pregnancy and at delivery. They and their children are also at increased risk of type 2 diabetes in the future. Gestational diabetes is diagnosed through prenatal screening, rather than through reported symptoms.

Untreated diabetes can damage the heart, blood vessels, eyes, kidneys, and nerves. Adults with diabetes have a two- to three-fold increased risk of heart attacks and strokes. Combined with reduced blood flow, neuropathy (nerve damage) in the feet increases the chance of foot ulcers, infection and eventual need for limb amputation. Diabetic retinopathy is an important cause of blindness, and occurs as a result of long-term accumulated damage to the small blood vessels in the retina. 2.6% of global blindness can be attributed to diabetes. Diabetes is among the leading causes of kidney failure.

Early diagnosis can be accomplished through relatively inexpensive testing of blood sugar. Treatment of diabetes involves diet and physical activity along with lowering blood glucose and the levels of other known risk factors that damage blood vessels. Tobacco use cessation is also important to avoid complications.

4 Software Requirements

The main function of this project is to enable the doctors to advise their patients with the help of the prediction software. The system should be accessible to the patient as well as the doctor, therefore it should have a web interface for the two parties to interact with.

4.1 Functionalities for Doctors

As an owner of a doctors account a doctor should be able to:

- Add new observations for the machine learning algorithm to predict.
- Analyse the patients previous records.
- Have a dashboard for patient performance monitoring.

4.2 Functionalities for Patients

As the patient, one should be able to:

- Should be able to see the predicted risk of developing diabetes.
- Should be able to view historic data.

4.3 Functionalities for Master Nodes

As the master nodes:

- Control chunk size.
- Remotely update the slave nodes.
- Check the status of nodes.
- Control the SQL database.

4.4 Functionalities for Slave Nodes

- Should have the machine learning algorithm tuned as per specifications.
- Remote update should be possible.
- Remote access should be possible.

5 Hardware Requirements

5.1 Raspberry Pi

According to raspberrypi.org, the Raspberry Pi 3 Model B is the earliest model of the third-generation Raspberry Pi. It replaced the Raspberry Pi 2 Model B in February 2016. Some of the key features of this single board computer or SBC are:

- Quad Core 1.2GHz Broadcom BCM2837 64bit CPU.
- 1GB RAM.
- BCM43438 wireless LAN and Bluetooth Low Energy (BLE) on board.
- 100 Base Ethernet.
- 40-pin extended GPIO.
- 4 USB 2 ports.
- Full size HDMI.
- Micro SD port for loading the operating system and storing data.

In this project I will be using 3 Raspberry Pi's to implement a cluster and deploy the machine learning algorithm on.

5.2 Switch

A network switch was used in the project to make up for the lack of ethernet ports available on the router.

5.3 Router

A router is required to assign internet protocol addresses to the nodes using dynamic host control protocol (DHCP). The router also is responsible for displaying the nodes connected to the network thereby displaying host names and making it more easier to capture the addresses of each node.

5.4 Master Node

The master node is assigned the task of managing the the slave nodes connected to the network. The master node and the slave nodes should be connected to the same database to run and execute queries. The master node shall also be assigned with the task of killing a particular process or shutting down a particular node if required.

5.5 Slave Node

The slave nodes are to be configured with the selected algorithm and are the most vital part of the structure. The slave nodes will be responsible for receiving instructions from the master node and will be responsible for learning from the dataset and then can generate predicted outcomes for the patient records received from the master.

5.6 Database

The database will store all the records for the patients and doctors. The database is a SQL database that will not be subjected to a change in the database schema structure for consistency.

5.7 Web Server

The web server provides an interface for the doctor and the patients to read the predictions and legacy data of the patient from the website.

6 Technology Stack

6.1 Python

6.1.1 Pandas

Pandas, is a library that is required for loading the comma separated value file into python. Pandas is a package for data manipulation and analysis. In particular, it offers data structures and operations for manipulating numerical tables and time series.

6.1.2 Sklearn

Scikit-learn is a free software machine learning library for the Python programming language. It features various classification, regression and clustering algorithms including support vector machines, random forests, gradient boosting, k-means and DBSCAN, and is designed to inter-operate with the Python numerical and scientific libraries NumPy and SciPy.

6.1.3 Numpy

NumPy is a library for the Python programming language, adding support for large, multi-dimensional arrays and matrices, along with a large collection of high-level mathematical functions to operate on these arrays. NumPy is open-source software and has many contributors.

6.1.4 Itertools

The itertools module standardizes a core set of fast, memory efficient tools that are useful by themselves or in combination. Together, they form an “iterator algebra” making it possible to construct specialized tools succinctly and efficiently in pure Python.

6.2 MYSQL

MySQL is an open-source relational database management system. MySQL is free and open-source software under the terms of the GNU General Public License, and is also available under a variety of proprietary licenses. MySQL was owned and sponsored by the Swedish company MySQL AB, which was bought by Sun Microsystems (now Oracle Corporation). MySQL is a component of the LAMP web application software stack (and others), which is an acronym for Linux, Apache, MySQL, Perl/PHP/Python. MySQL is used by many database-driven web applications, including Drupal, Joomla, phpBB, and WordPress. MySQL is also used by many popular websites, including Facebook, Flickr, MediaWiki, Twitter and YouTube.

6.3 Apache Web Server

The Apache HTTP Server, colloquially called Apache, is free and open-source cross-platform web server software, released under the terms of Apache License 2.0. Apache is developed and maintained by an open community of developers under the auspices of the Apache Software Foundation. The vast majority of Apache HTTP Server instances run on a Linux distribution, but current versions also run on Microsoft Windows and a wide variety of Unix-like systems. Past versions also ran on OpenVMS, NetWare, OS/2 and other operating systems, including ports to mainframes.

6.4 PHP

PHP is a general-purpose programming language originally designed for web development. PHP originally stood for Personal Home Page, but it now stands for the recursive initialism PHP: Hypertext Preprocessor. PHP code may be executed with a command line interface (CLI), embedded into HTML code, or used in combination with various web template systems, web content management systems, and web frameworks. PHP code is usually processed by a PHP interpreter implemented as a module in a web server or as a Common Gateway Interface (CGI) executable. The web server outputs the results of the interpreted and executed PHP code, which may be any type of data, such as generated HTML code or binary image data. PHP can be used for many programming tasks outside of the web context, such as standalone graphical applications and robotic drone control.

6.5 AJAX

Ajax is a set of web development techniques using many web technologies on the client side to create asynchronous web applications. With Ajax, web applications can send and retrieve data from a server asynchronously (in the background) without interfering with the display and behavior of the existing page. By decoupling the data interchange layer from the presentation layer, Ajax allows web pages and, by extension, web applications, to change content dynamically without the need to reload the entire page.[3] In practice, modern implementations commonly utilize JSON instead of XML.

Ajax is not a single technology, but rather a group of technologies. HTML and CSS can be used in combination to mark up and style information. The webpage can then be modified by JavaScript to dynamically display—and allow the user to interact with—the new information. The built-in XMLHttpRequest object, or since 2017 the new "fetch()" function within JavaScript, is commonly used to execute Ajax on web pages allowing websites to load content onto the screen without refreshing the page. Ajax is not a new technology, or different language, just existing technologies used in

new ways.

6.6 GitHub

GitHub is a global company that provides hosting for software development version control using Git. It offers all of the distributed version control and source code management (SCM) functionality of Git as well as adding its own features. It provides access control and several collaboration features such as bug tracking, feature requests, task management, and wikis for every project.

6.7 SSH

Secure Shell is a cryptography network protocol for operating network services securely over an unsecured network. Typical applications include remote command-line, login, and remote command execution, but any network service can be secured with SSH. The SSH is used to communicate between the master and slave nodes.

7 Methodology

The dataset used for this project is from the UCI Machine learning repository and can be found at their respective website. Table 2 is an outcome of the accuracy of the algorithms subjected to train and 3 is the accuracy obtained from the predicted outcomes for given sets of patients.

Firstly, I began the experiment by examining the amount of people that were recorded by the dataset. We can see that the dataset contains a total of 768 observations. We can see that 268 people are healthy, while the rest of the 500 subjects are diabetic. If converted to percentage, 65.1% of the subjects are healthy and 34.9% of subjects are diabetic.

To proceed with the exploratory data analysis, we can examine the missing values in the dataset. We can see that the dataset has multiple missing values across the columns of the dataset. Plotting of boxplots reveals that the dataset has multiple outliers. To replace these missing values with values to make the predictions more accurate we need to understand the relation of the data to the target variable. To get an understanding we can use the table 6 for clarity on the correlation of values.

To fill the missing values in the dataset we have to consider each of the attributes. For example in the insulin parameter, I have observed that the median value for diabetic people is 169.50 muU/ml, while for the section of healthy participants the median insulin value is recorded at 102.50 muU/ml. Likewise, I carried out the observations on the other attributes of the dataset to fill the missing values via the median of healthy subjects for the healthy subjects whose insulin record is missing and diabetic insulin median for the subjects that are diabetic and do not have values in the columns.

In order to generate new values for the dataset to render more accurate outcomes we look at scatterplots for our attributes. An example would be the Glucose vs Age scatterplot. According to the scatter plot of glucose to age we can see that higher glucose value with increase in age, depicts the glucose level of a diabetic person.

8 Deploying to Cluster

A cluster is a group of two or more servers connected to each other in such a way that they behave like a single server. Each machine in the cluster is called a node. Because each machine in the cluster runs the same services as other machines in the cluster, any machine can stand in for any other machine in the cluster. This becomes important when one machine goes down or must be taken out of service for a time. The remaining machines in the cluster can seamlessly take over the work of the downed machine, providing users with uninterrupted access to services and data.

8.1 Benefits of Clustering

- **Increased resource availability** If one Intelligence Server in a cluster fails, the other Intelligence Servers in the cluster can pick up the workload. This prevents the loss of valuable time and information if a server fails. Strategic resource usage: You can distribute projects across nodes in whatever configuration you prefer. This reduces overhead because not all machines need to be running all projects, and allows you to use your resources flexibly.
- **Increased performance** Multiple machines provide greater processing power.
- **Greater scalability** As your user base grows and report complexity increases, your resources can grow.
- **Simplified management** Clustering simplifies the management of large or rapidly growing systems.

Fail-over support ensures that a business intelligence system remains available for use if an application or hardware failure occurs. Clustering provides failover support in two ways:

- **Load redistribution** When a node fails, the work for which it is responsible is directed to another node or set of nodes. Request recovery: When a node fails, the system attempts to reconnect MicroStrategy Web users with queued or processing requests to another node. Users must log in again to be authenticated on the new node. The user is prompted to resubmit job requests.
- **Load Balancing** Load balancing is a strategy aimed at achieving even distribution of user sessions across Intelligence Servers, so that no single machine is overwhelmed. This strategy is especially valuable when it is difficult to predict the number of requests a server will receive. MicroStrategy achieves four-tier load balancing by incorporating load balancers into the MicroStrategy Web and Web products. Load is calculated as the number of user sessions connected to a

node. The load balancers collect information on the number of user sessions each node is carrying. Using this information at the time a user logs in to a project, MicroStrategy Web connects them to the Intelligence Server node that is carrying the lightest session load. All requests by that user are routed to the node to which they are connected until the user disconnects from the MicroStrategy Web product.

- **Project Distribution and Project Fail over** When you set up several server machines in a cluster, you can distribute projects across those clustered machines or nodes in any configuration, in both Windows and Linux environments. All servers in a cluster do not need to be running all projects. Each node in the cluster can host a different set of projects, which means only a subset of projects need to be loaded on a specific Intelligence Server machine. This feature provides you with flexibility in using your resources, and it provides better scalability and performance because of less overhead on each Intelligence Server machine. Distributing projects across nodes also provides project fail-over support. For example, one server is hosting project A and another server is hosting projects B and C. If the first server fails, the other server can host all three projects to ensure project availability. Project creation, duplication, and deletion in a three-tier, or server, connection are automatically broadcast to all nodes during run-time to ensure synchronization across the cluster.
- **Work Fencing** User fences and workload fences allow you to reserve nodes of a cluster for either users or a project subscriptions.

9 Future Scope

10 Tables

Column Name	Description
Pregnancies	Number of times pregnant.
Glucose	Plasma glucose concentration a 2 hours in an oral glucose tolerance test.
Blood Pressure	Diastolic blood pressure (mm Hg)
Skin Thickness	Triceps skin fold thickness (mm)
Insulin	2-Hour serum insulin ($\mu\text{U}/\text{ml}$)
Body Mass Index	Body mass index ($\text{weight in kg}/(\text{height in m})^2$)
Diabetes Pedigree Function	Diabetes pedigree function
Age	Age (years)
Outcome	Class variable (0 or 1) 268 of 768 are 1, the others are 0

Table 1: Pima Indians dataset header.

Algorithm	Additional Parameters	Train Set Accuracy
K Nearest Neighbour	-	0.790
Logistic Regression	$C = 1$	0.781
Logistic Regression	$C = 0.01$	0.700
Logistic Regression	$C = 100$	0.785
Decision Tree	-	1.000
Decision Tree	Max Depth = 3	0.773
Random Forest	Estimators = 100	1.000
Random Forest	Estimators = 100; Max Depth = 3	0.800
Gradient Boosting	-	0.917
Gradient Boosting	Max Depth = 1	0.804
Gradient Boosting	Learning Rate = 0.01	0.802
Support Vector Machine	-	1.000
Support Vector Machine	Train and Test set scaled using MinMaxScaler	0.770
Support Vector Machine	$C = 1000$	0.790
MLP Classifier	Random State = 42	0.730
MLP Classifier	Random State = 0	0.823
MLP Classifier	Max Iterations = 1000	0.908
MLP Classifier	Max Iterations = 1000; Alpha = 1; Random State = 0	0.806

Table 2: Train accuracy using various machine learning algorithms with various parameters.

Algorithm	Additional Parameters	Test Set Accuracy
K Nearest Neighbour	-	0.780
Logistic Regression	$C = 1$	0.771
Logistic Regression	$C = 0.01$	0.703
Logistic Regression	$C = 100$	0.766
Decision Tree	-	0.714
Decision Tree	Max Depth = 3	0.740
Random Forest	Estimators = 100	0.786
Random Forest	Estimators = 100; Max Depth = 3	0.755
Gradient Boosting	-	0.792
Gradient Boosting	Max Depth = 1	0.781
Gradient Boosting	Learning Rate = 0.01	0.776
Support Vector Machine	-	0.650
Support Vector Machine	Train and Test set scaled using MinMaxScaler	0.770
Support Vector Machine	$C = 1000$	0.797
MLP Classifier	Random State = 42	0.720
MLP Classifier	Random State = 0	0.802
MLP Classifier	Max Iterations = 1000	0.792
MLP Classifier	Max Iterations = 1000; Alpha = 1; Random State = 0	0.797

Table 3: Test accuracy using various machine learning algorithms with various parameters.

Column	Number of Non-Zero Values	Percentage
Pregnancies	0	0
Glucose	763	0.65
Blood Pressure	733	4.56
Skin Thickness	541	29.56
Insulin	394	48.7
Body Mass Index	757	1.43
Diabetes Pedigree Function	0	0
Age	0	0
Outcome	0	0

Table 4: Null Value Check.

Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
6	148	72	35	0	33.6	0.627	50	1
1	85	66	29	0	26.6	0.351	31	0
8	183	64	0	0	23.3	0.672	32	1
1	89	66	23	94	28.1	0.167	21	0
0	137	40	35	168	43.1	2.288	33	1
5	116	74	0	0	25.6	0.201	30	0
3	78	50	32	88	31.0	0.248	26	1
10	115	0	0	0	35.3	0.134	29	0
2	197	70	45	543	30.5	0.158	53	1
8	125	96	0	0	0.0	0.232	54	1
4	110	92	0	0	37.6	0.191	30	0
10	168	74	0	0	38.0	0.537	34	1
10	139	80	0	0	27.1	1.441	57	0
1	189	60	23	846	30.1	0.398	59	1
5	166	72	19	175	25.8	0.587	51	1

Table 5: First 15 rows of the raw dataset.

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DPF	Age	Outcome
Pregnancies	1.000000	0.129459	0.141282	-0.081672	-0.073535	0.017683	-0.03352	0.544341	0.221898
Glucose	0.129459	1.000000	0.152590	0.057328	0.331357	0.221071	0.137337	0.263514	0.466581
BloodPressure	0.141282	0.152590	1.000000	0.207371	0.088933	0.281805	0.041265	0.239528	0.065068
SkinThickness	-0.081672	0.057328	0.207371	1.000000	0.436783	0.392573	0.183928	-0.113970	0.074752
Insulin	-0.073535	0.331357	0.088933	0.436783	1.000000	0.197859	0.185071	-0.042163	0.130548
BMI	0.017683	0.221071	0.281805	0.392573	0.197859	1.000000	0.140647	0.036242	0.292695
DPF	-0.033523	0.137337	0.041265	0.183928	0.185071	0.140647	1.000000	0.033561	0.173844
Age	0.544341	0.263514	0.239528	-0.113970	-0.042163	0.036242	0.033561	1.000000	0.238356
Outcome	0.221898	0.466581	0.065068	0.074752	0.130548	0.292695	0.173844	0.238356	1.000000

Table 6: Correlation Plot.

Fold	Precision	F1 Score	Accuracy	Recall	ROC-AUC Curve
1	0.811	0.804	0.864	0.796	0.908
2	0.787	0.733	0.825	0.685	0.896
3	0.841	0.755	0.844	0.685	0.902
4	0.820	0.796	0.863	0.774	0.920
5	0.783	0.832	0.876	0.887	0.933
mean	0.809	0.784	0.854	0.765	0.912
std	0.021	0.035	0.018	0.076	0.013

Table 7: MLP 5 fold verification.

Insulin	
Subject Type	Level
Healthy	102.5
Diabetic	169.5

Table 8: Insulin Median Comparison.

Glucose	
Subject Type	Level
Healthy	107.0
Diabetic	140.0

Table 9: Glucose Median Comparison.

Skin Thickness	
Subject Type	Level
Healthy	27.0
Diabetic	32.0

Table 10: Skin Thickness Median Comparison.

Blood Pressure	
Subject Type	Level
Healthy	70.0
Diabetic	74.0

Table 11: Blood Pressure Median Comparison.

Body Mass Index	
Subject Type	Level
Healthy	30.1
Diabetic	34.3

Table 12: Body Mass Index Median Comparison.

11 Appendix I

11.1 K Nearest Neighbours

```
1 class sklearn.neighbors.KNeighborsClassifier(n_neighbors=5, weights='uniform', algorithm='auto', leaf_size=30, p=2, metric='minkowski', metric_params=None, n_jobs=None, **kwargs)
```

11.2 Logistic Regression

```
1 class sklearn.linear_model.LogisticRegression(penalty='l2', dual=False, tol=0.0001, C=1.0, fit_intercept=True, intercept_scaling=1, class_weight=None, random_state=None, solver='lbfgs', max_iter=100, multi_class='auto', verbose=0, warm_start=False, n_jobs=None, l1_ratio=None)
```

11.3 Decision Tree

```
1 class sklearn.tree.DecisionTreeClassifier(criterion='gini', splitter='best', max_depth=None, min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features=None, random_state=None, max_leaf_nodes=None, min_impurity_decrease=0.0, min_impurity_split=None, class_weight=None, presort='deprecated', ccp_alpha=0.0)
```

11.4 Random Forest Classifier

```
1 class sklearn.ensemble.RandomForestClassifier(n_estimators=100, criterion='gini', max_depth=None, min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0, min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=None, random_state=None, verbose=0, warm_start=False, class_weight=None, ccp_alpha=0.0, max_samples=None)
```

11.5 Gradient Boosting

```
1 class sklearn.ensemble.GradientBoostingClassifier(loss='deviance', learning_rate=0.1, n_estimators=100, subsample=1.0, criterion='friedman_mse', min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_depth=3, min_impurity_decrease=0.0, min_impurity_split=None, init=None, random_state=None, max_features=None, verbose=0, max_leaf_nodes=None, warm_start=False, presort='deprecated', validation_fraction=0.1, n_iter_no_change=None, tol=0.0001, ccp_alpha=0.0)
```


11.6 Multi Layered Perceptron

```
1 class sklearn.neural_network.MLPClassifier(hidden_layer_sizes=(100, ),
      activation='relu', solver='adam', alpha=0.0001, batch_size='auto',
      learning_rate='constant', learning_rate_init=0.001, power_t
      =0.5, max_iter=200, shuffle=True, random_state=None, tol=0.0001,
      verbose=False, warm_start=False, momentum=0.9, nesterovs_momentum=
      True, early_stopping=False, validation_fraction=0.1, beta_1=0.9,
      beta_2=0.999, epsilon=1e-08, n_iter_no_change=10, max_fun=15000)
```

11.7 SVM

```
1 class sklearn.svm.SVC(C=1.0, kernel='rbf', degree=3, gamma='scale',
      coef0=0.0, shrinking=True, probability=False, tol=0.001,
      cache_size=200, class_weight=None, verbose=False, max_iter=-1,
      decision_function_shape='ovr', break_ties=False, random_state=None
      )
```

13 Appendix II

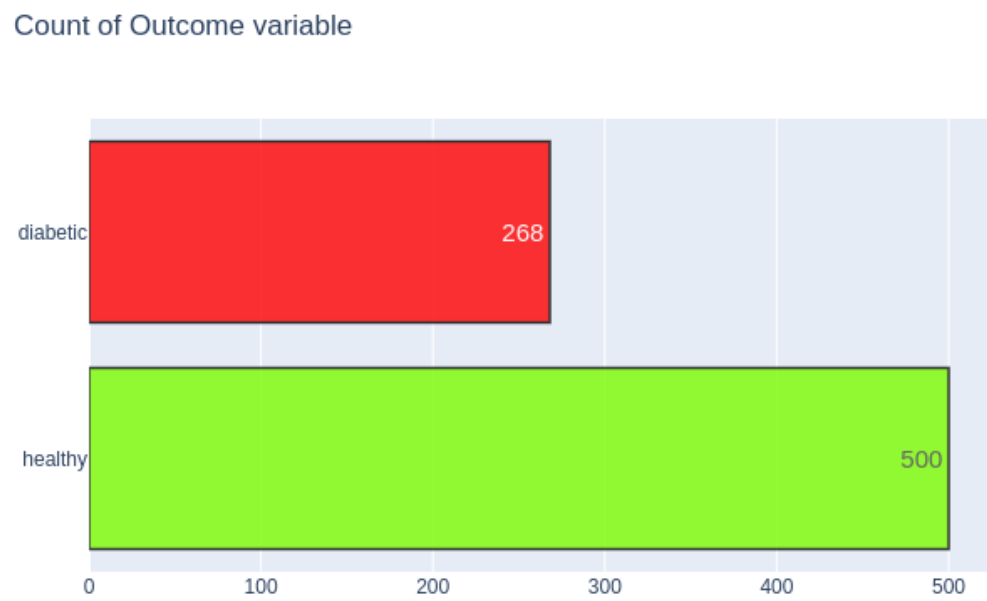


Figure 1: Diabetic v/s Healthy Subject Count.

BMI

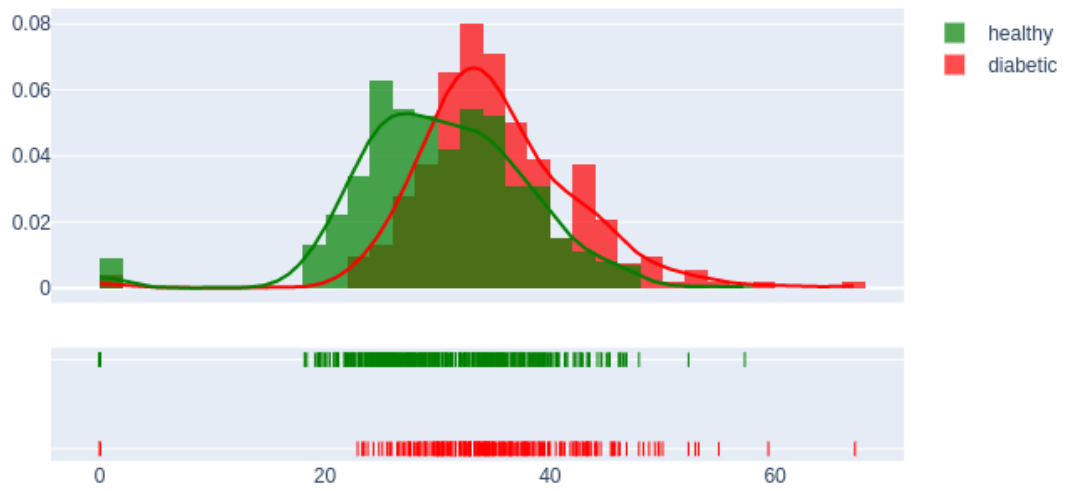


Figure 2: Subject distribution across Body Mass Index range.

Age

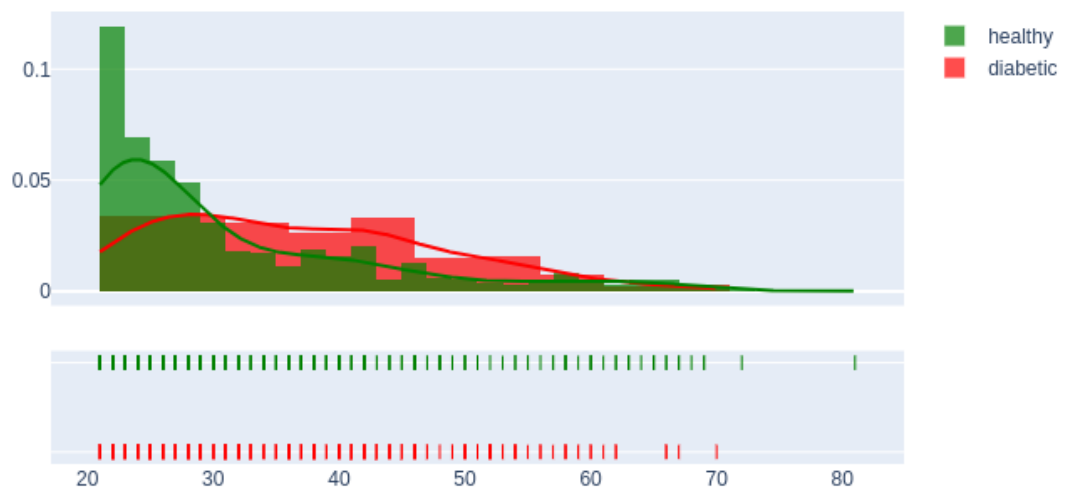


Figure 3: Subject distribution across age.

Insulin

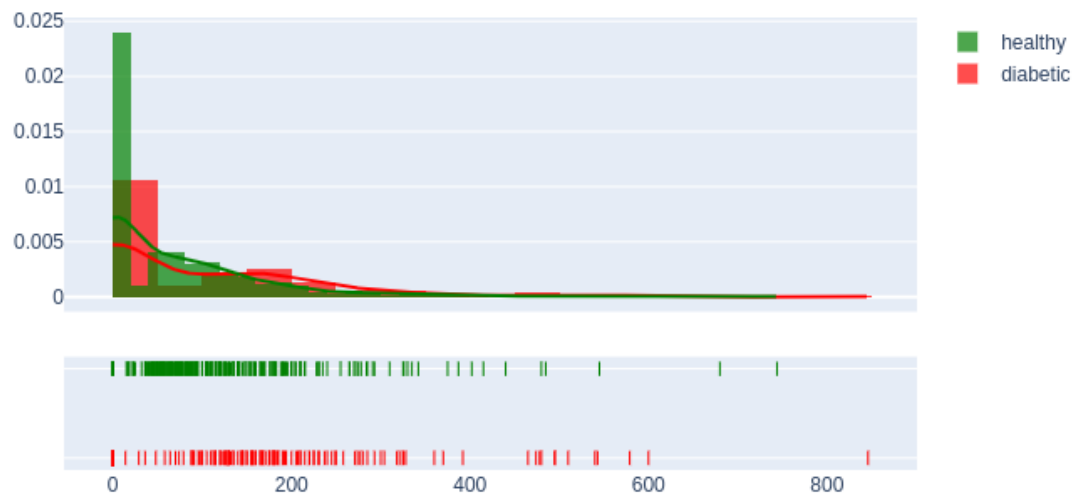


Figure 4: Subjects Insulin Distribution across ranges.

Pregnancies

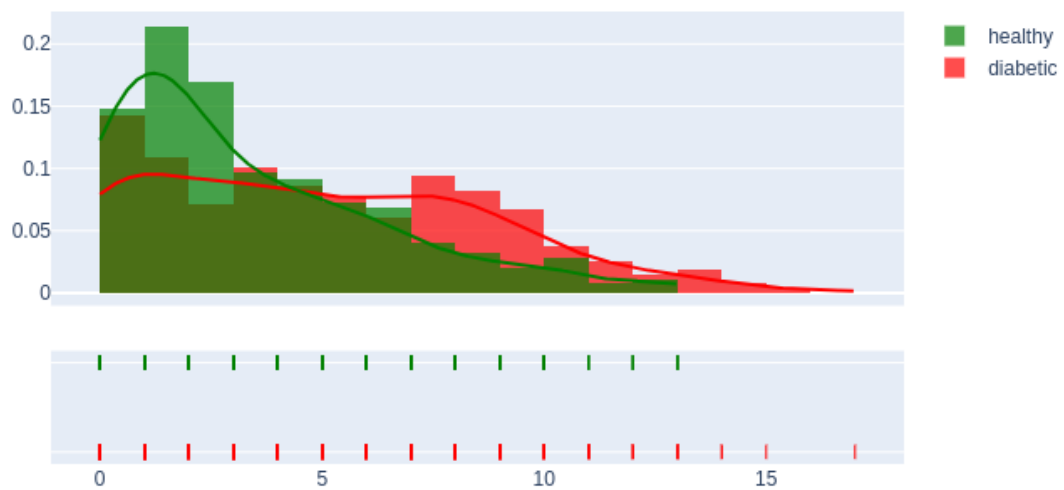


Figure 5: Subjects distribution via Pregnancies.

DiabetesPedigreeFunction

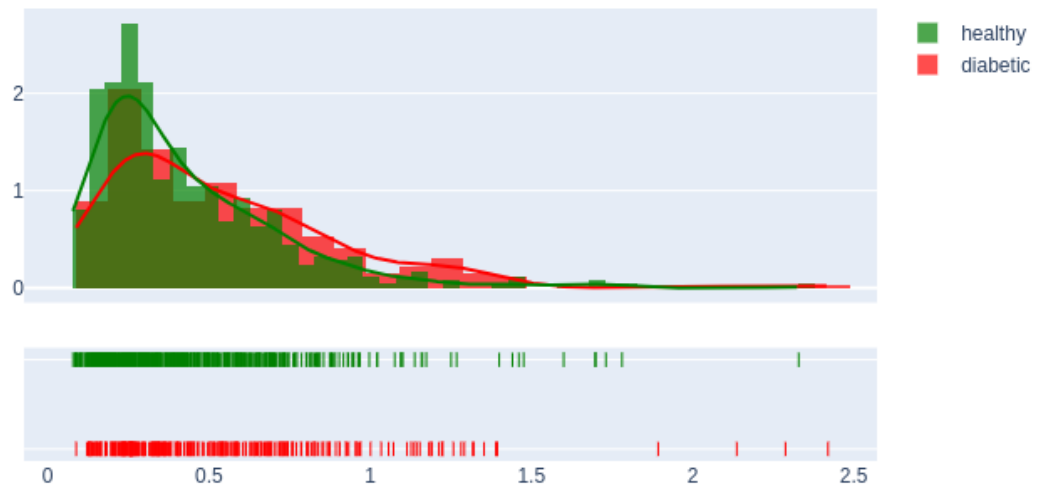


Figure 6: Subjects distribution via Diabetes Pedigree Function.

N0

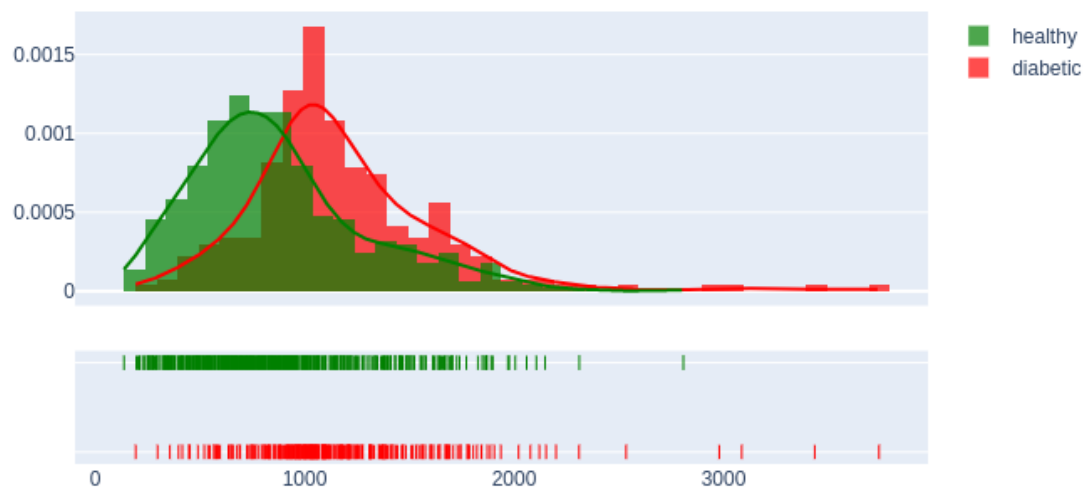


Figure 7: Plot of N0

Distribution of Outcome variable

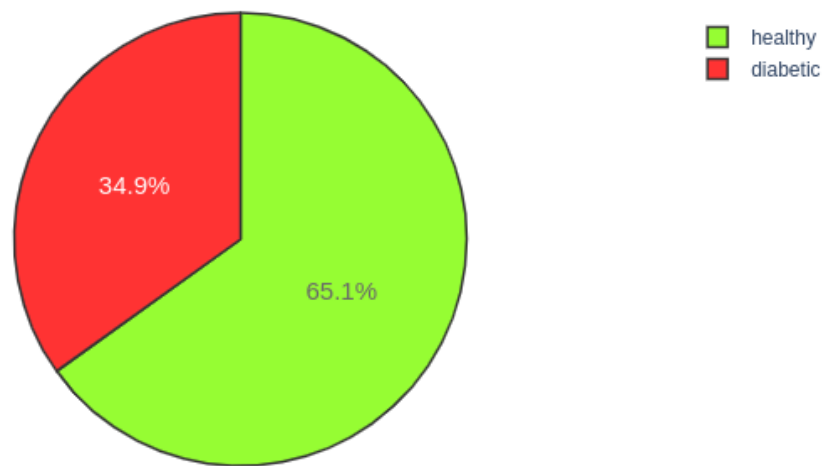


Figure 8: Diabetic v/s Healthy Subjects Percentage.

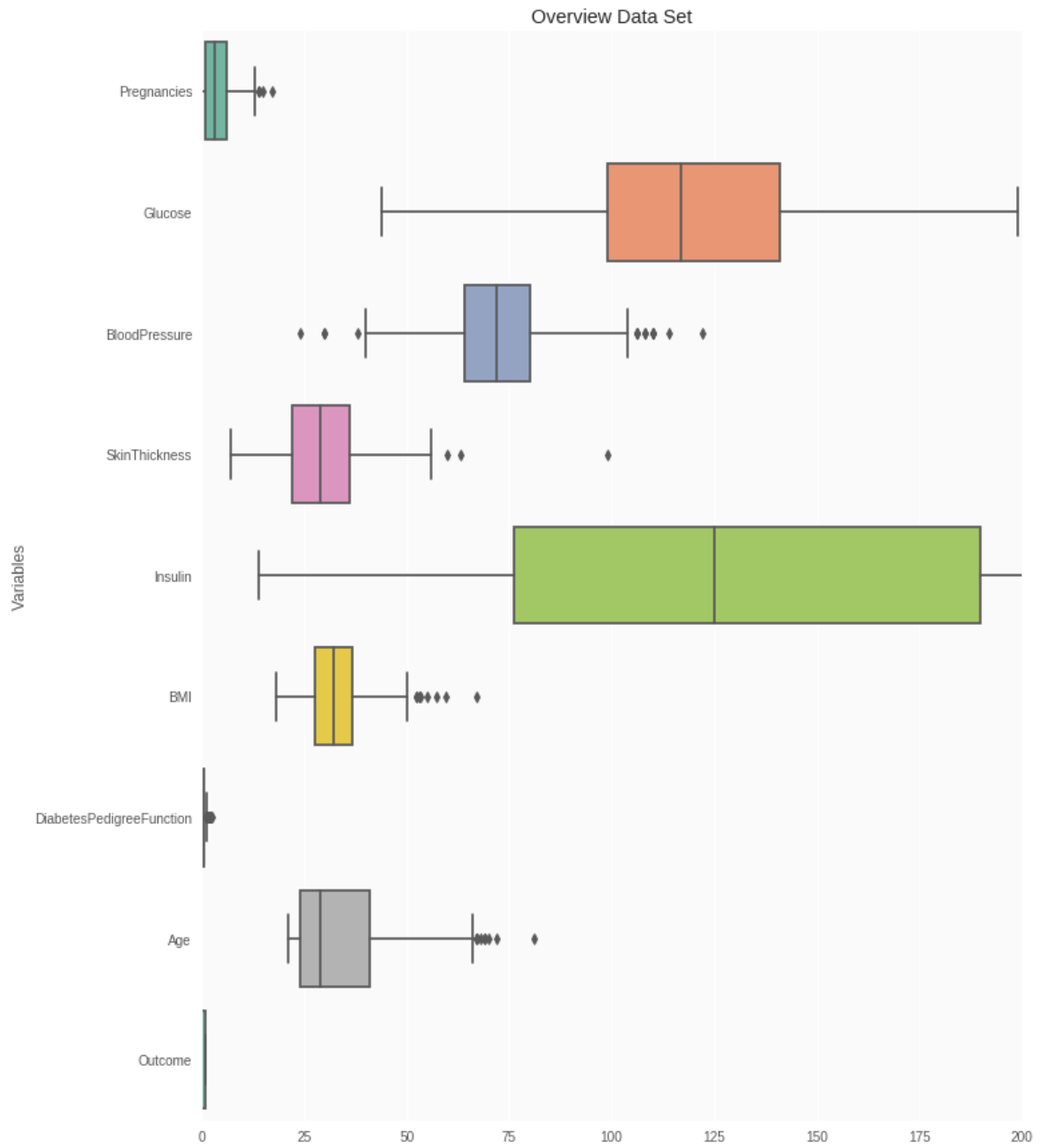


Figure 9: Boxplot for all attributes with outliers.

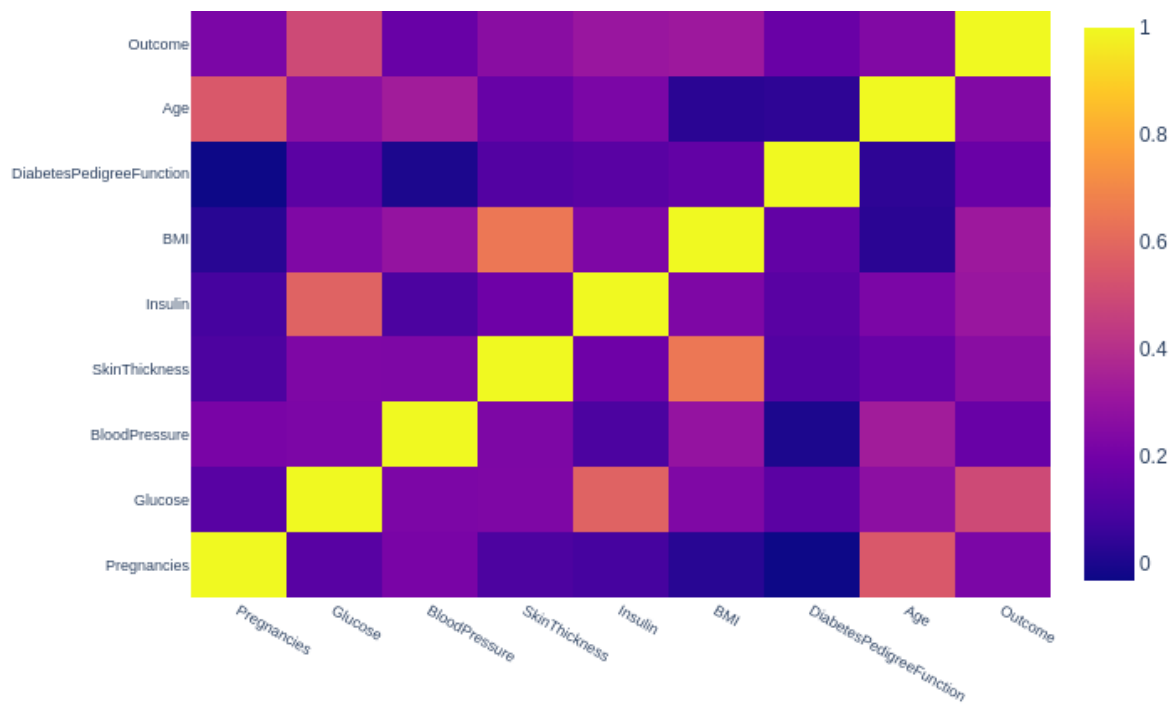


Figure 10: Heatmap using Pearsons Correlation Coefficient.

Missing Values (count & %)

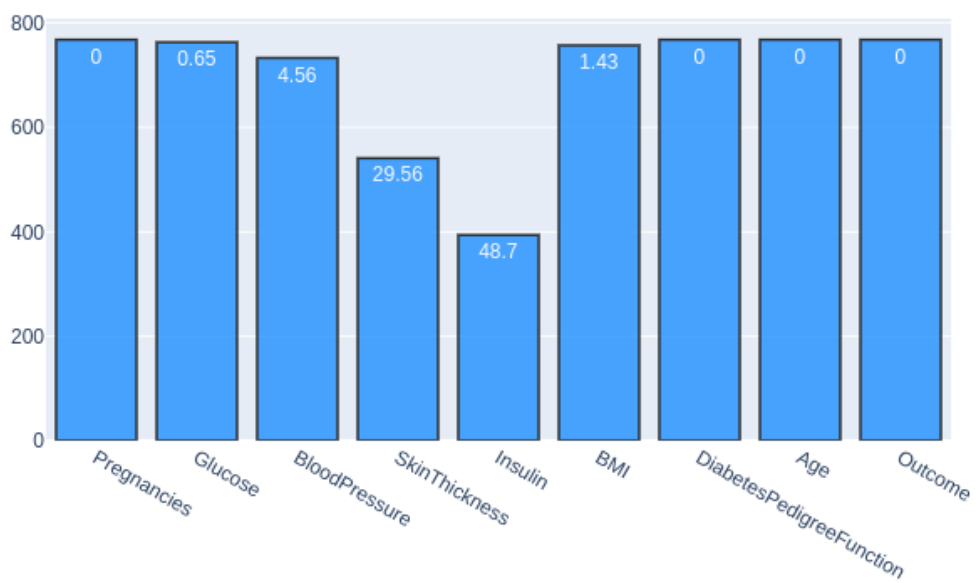


Figure 11: Number of missing values in count and percentage.

Glucose vs Age

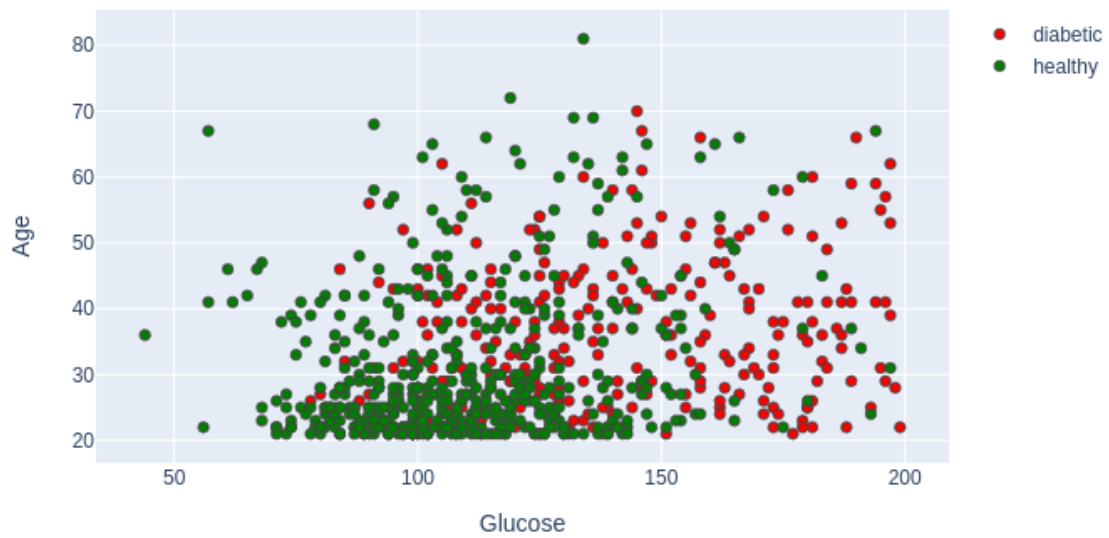


Figure 12: Glucose v/s Age scatterplot.

Glucose

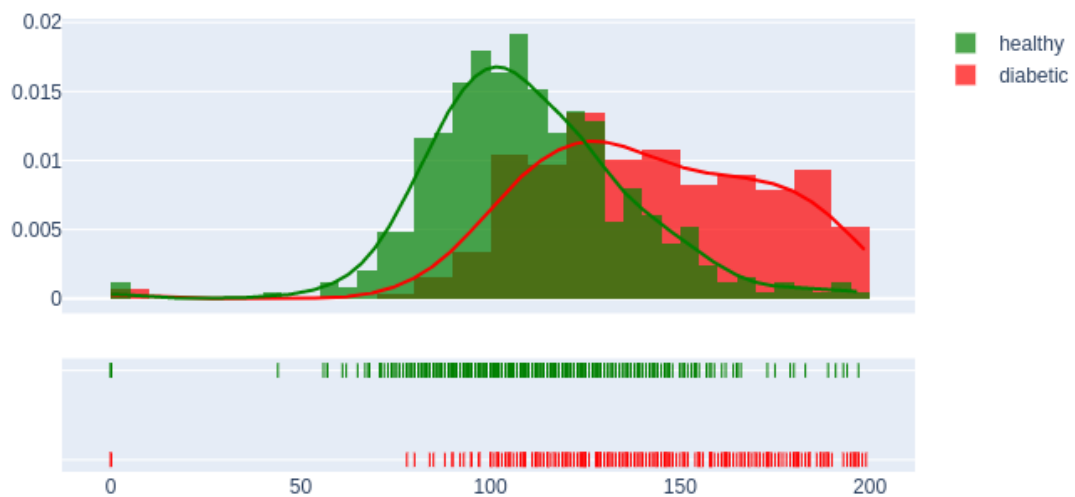


Figure 13: Subjects glucose distribution.

SkinThickness

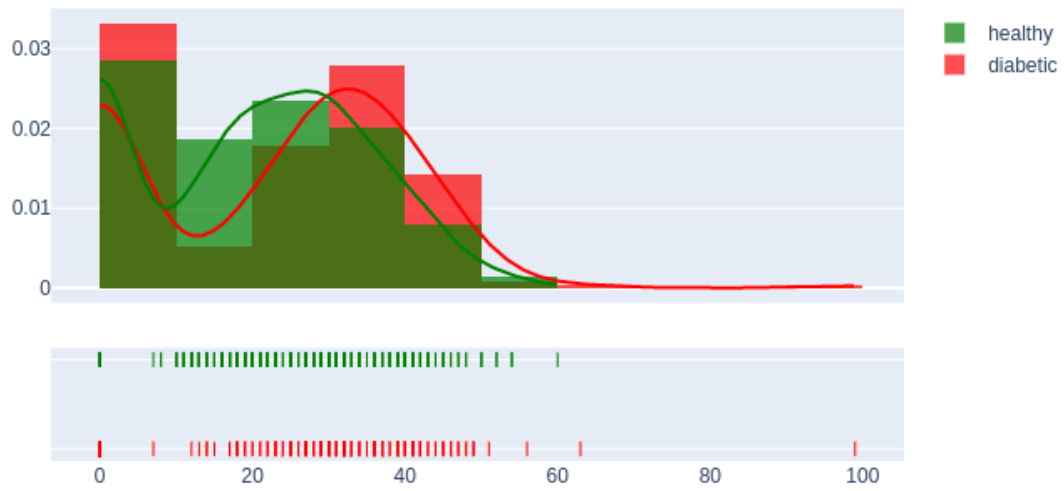


Figure 14: Skin Thickness distribution of subjects.

BloodPressure

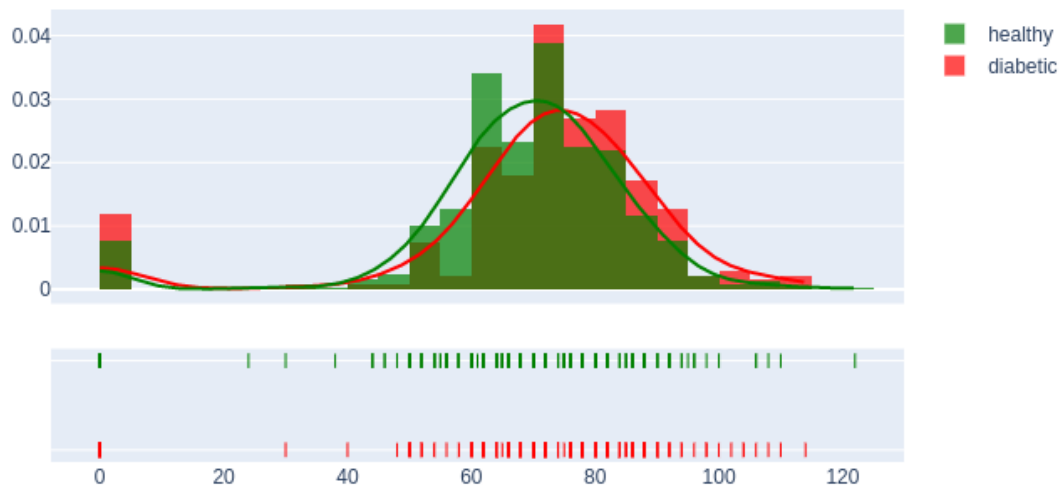


Figure 15: Blood Pressure distribution of subjects.

Insulin

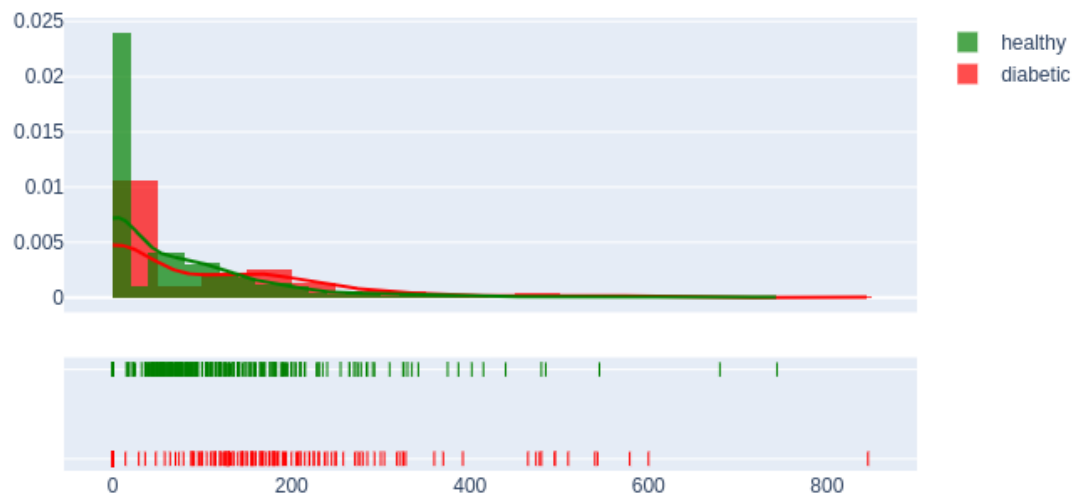


Figure 16: Subjects Insulin distributon.

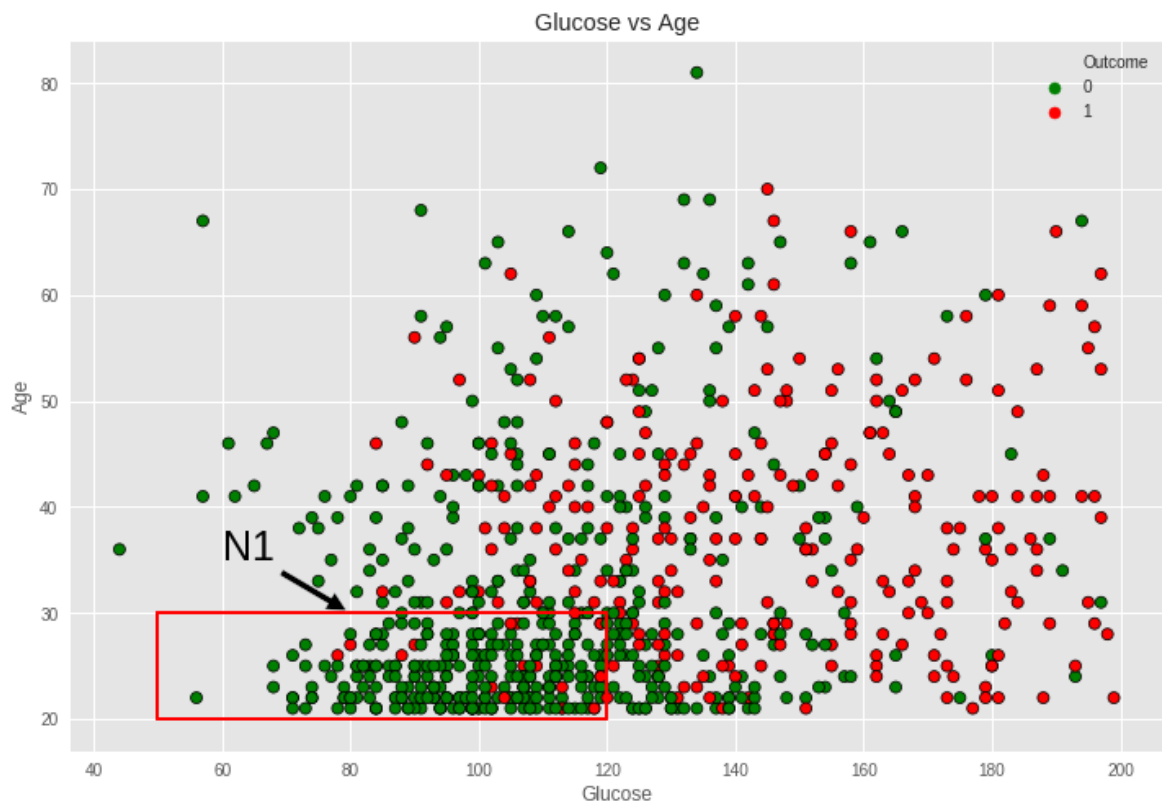


Figure 17: New feature N1.

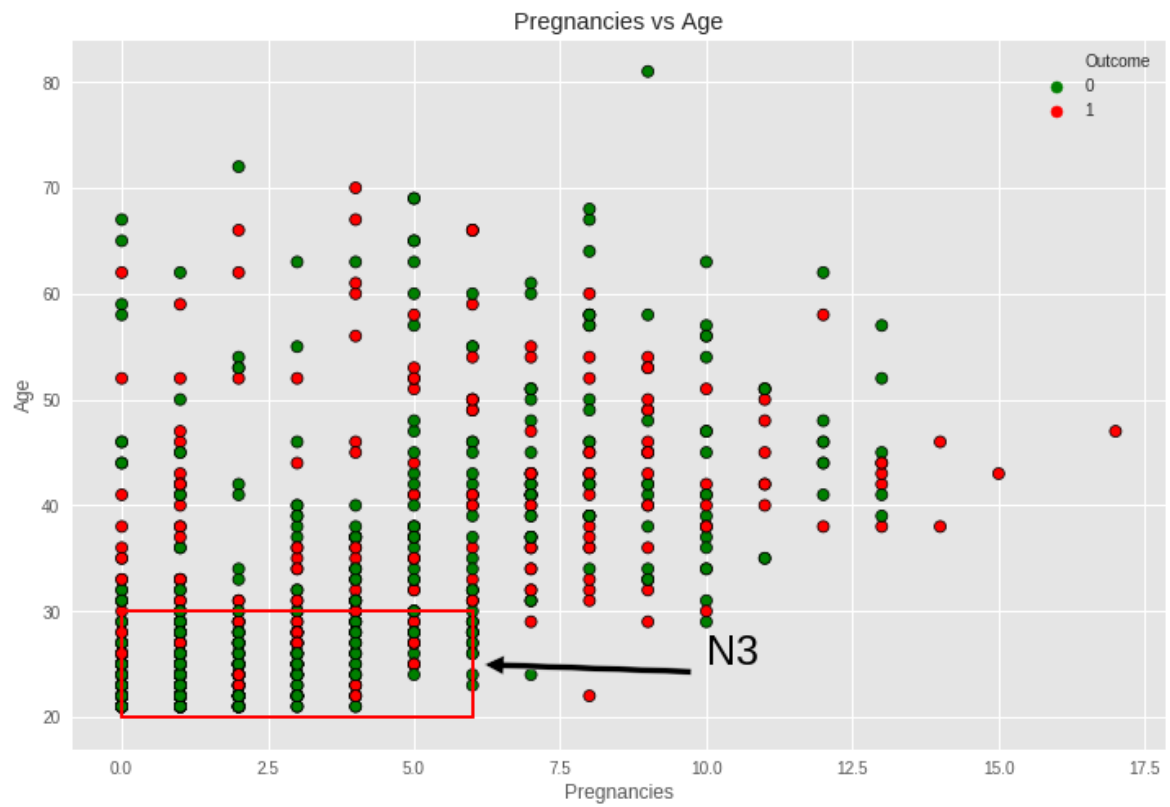


Figure 18: New feature N3.

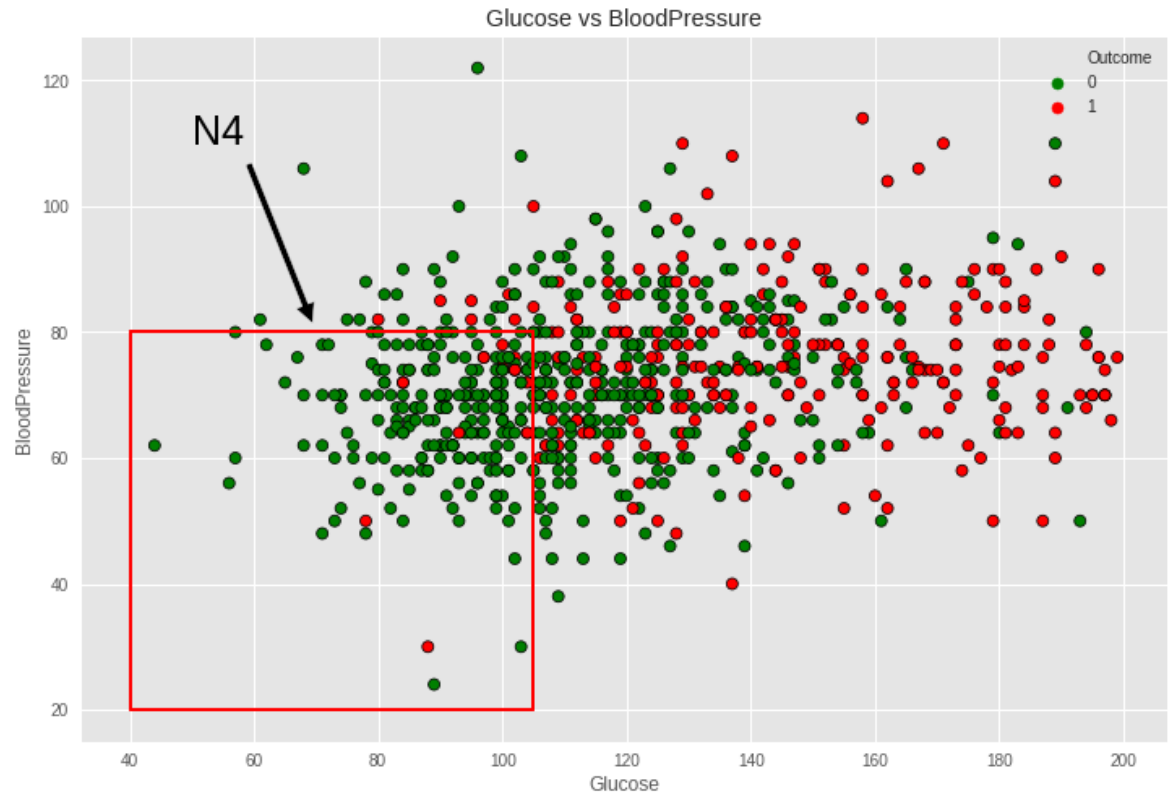


Figure 19: New feature N4.

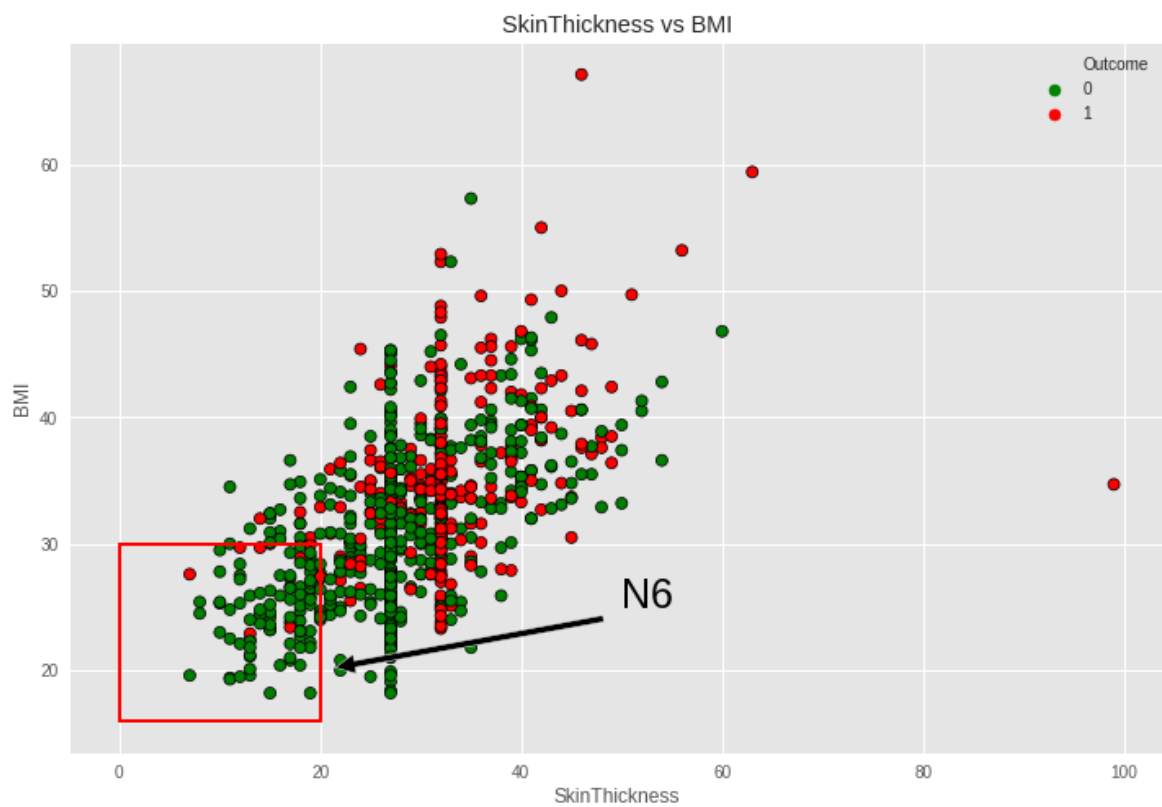


Figure 20: New feature N6.

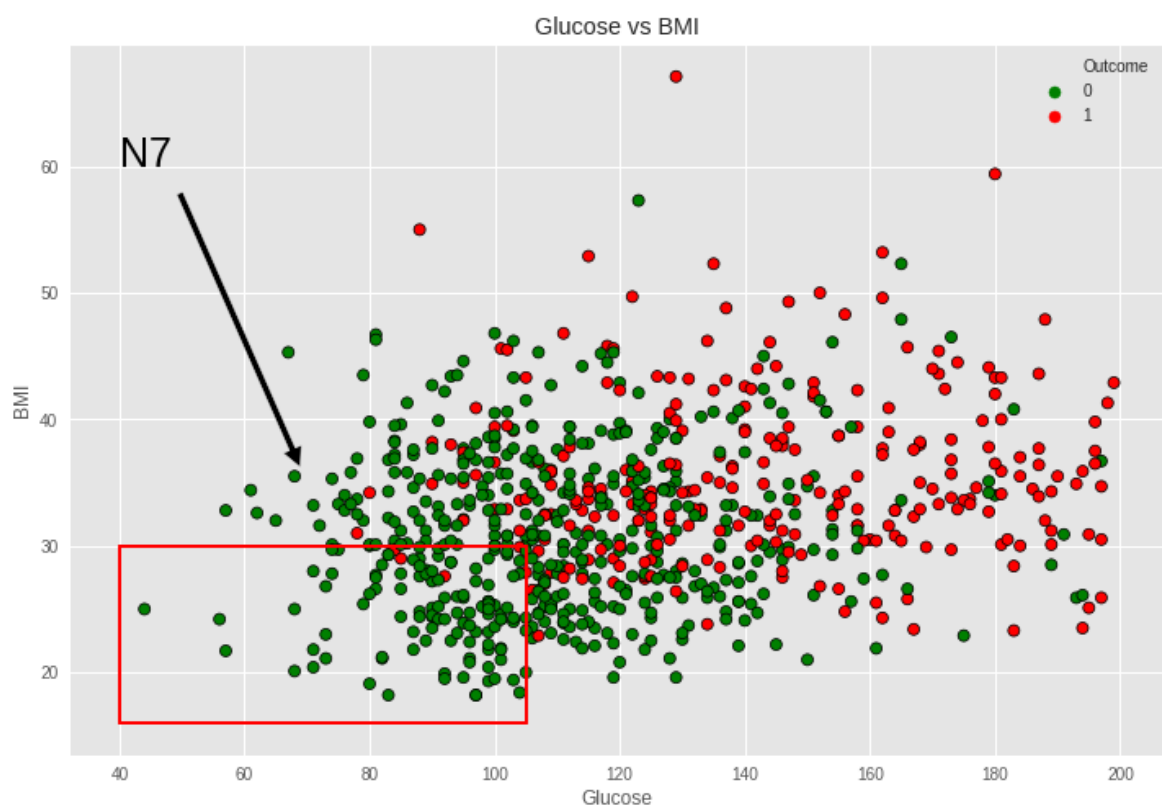


Figure 21: New feature N7.

N1 :Glucose <= 120 and Age <= 30

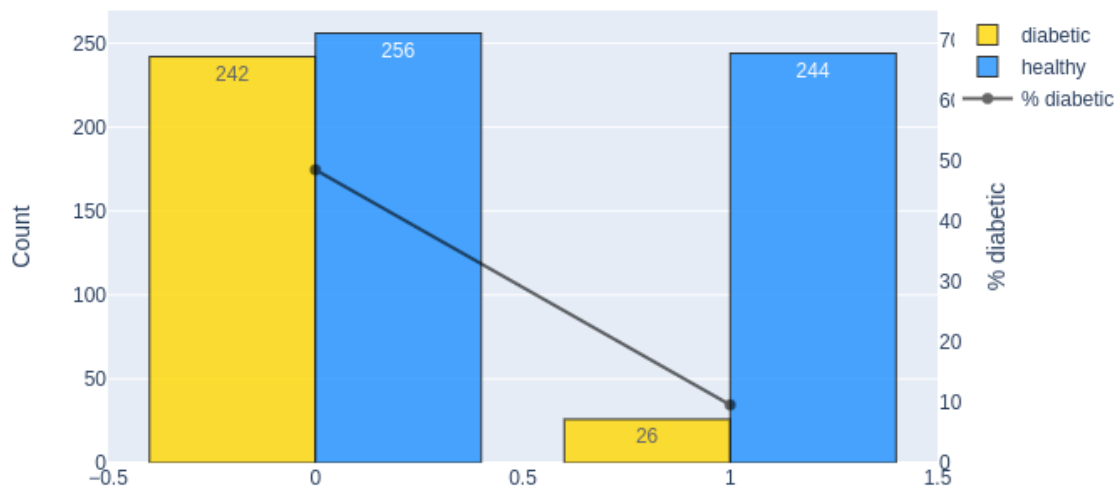


Figure 22: N1 barplot for diabetic and healthy population.

N1 distribution by target
(Glucose <= 120 and Age <= 30)

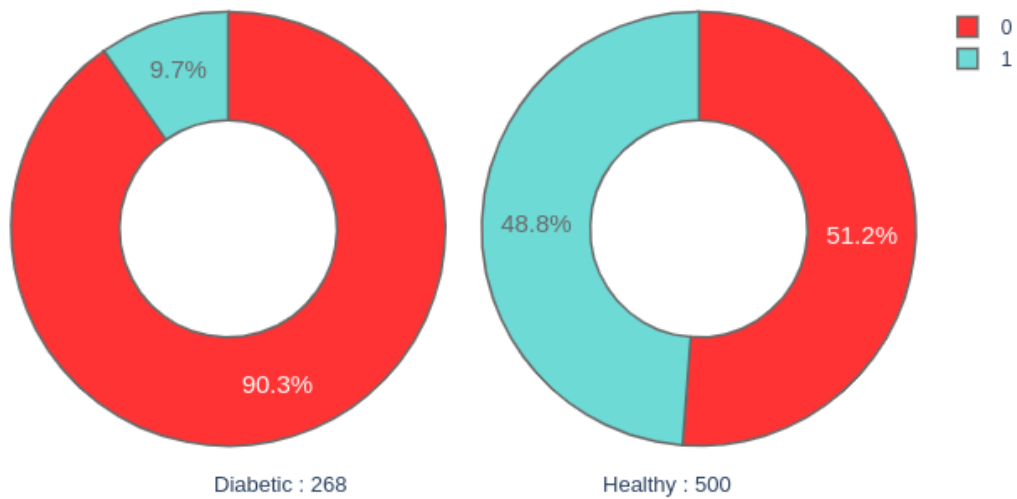


Figure 23: N1 distribution in percentage.

N2 : BMI <= 30

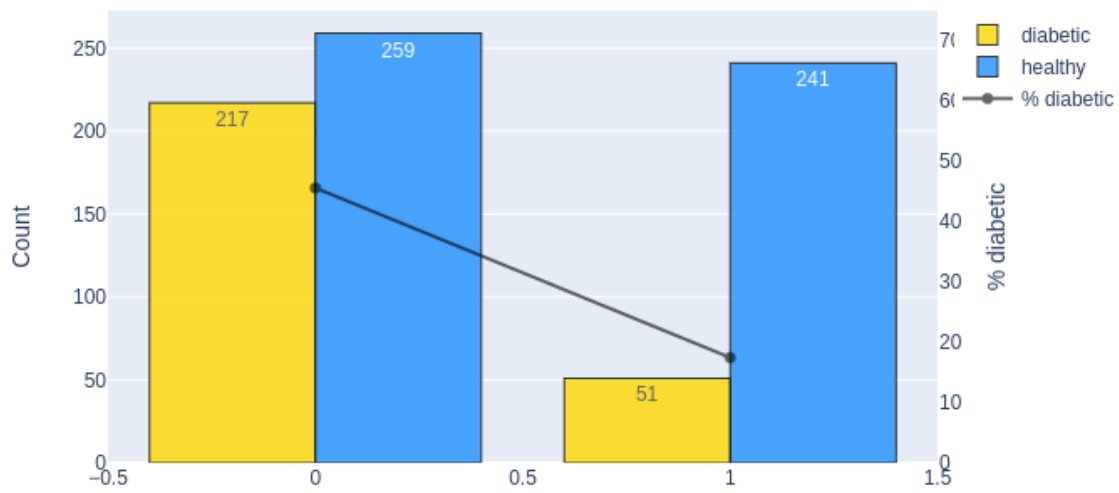


Figure 24: N2 barplot for diabetic and healthy population.

N2 distribution by target
BMI <= 30

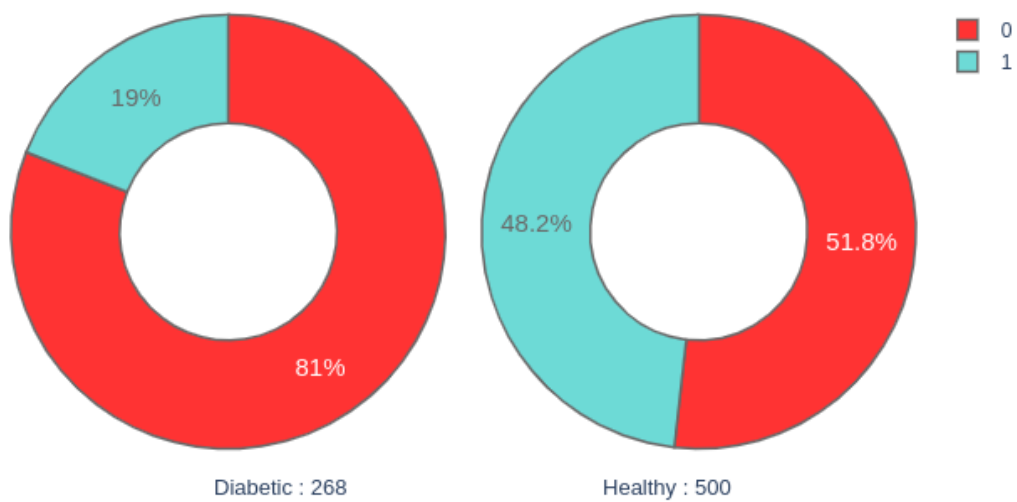


Figure 25: N2 distribution in percentage.

Pregnancies vs Age

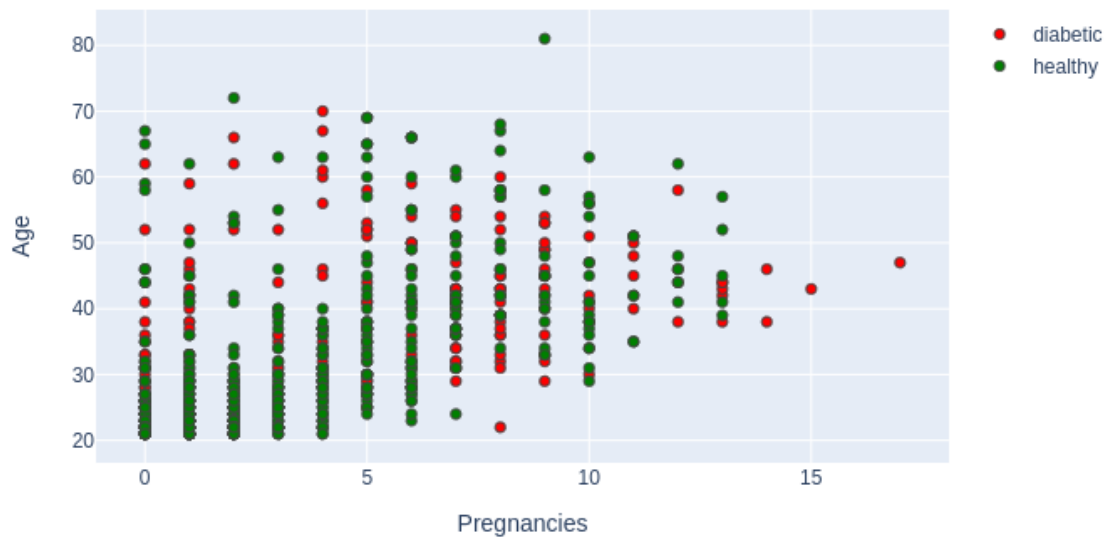


Figure 26: Pregnancies v/s age scatterplot.

N3 : Age ≤ 30 and Pregnancies ≤ 6

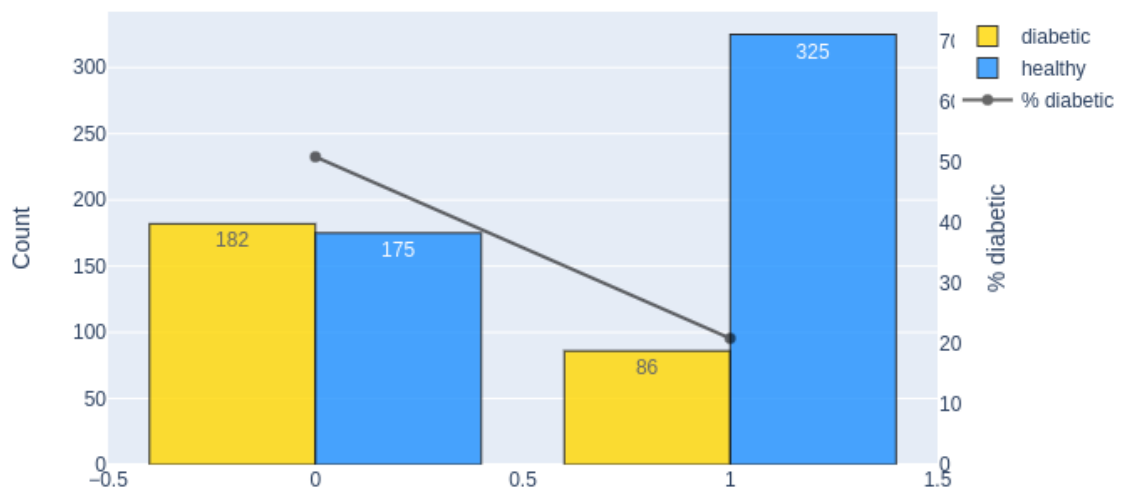


Figure 27: N3 barplot for diabetic and healthy population.

N3 distribution by target
Age ≤ 30 and Pregnancies ≤ 6

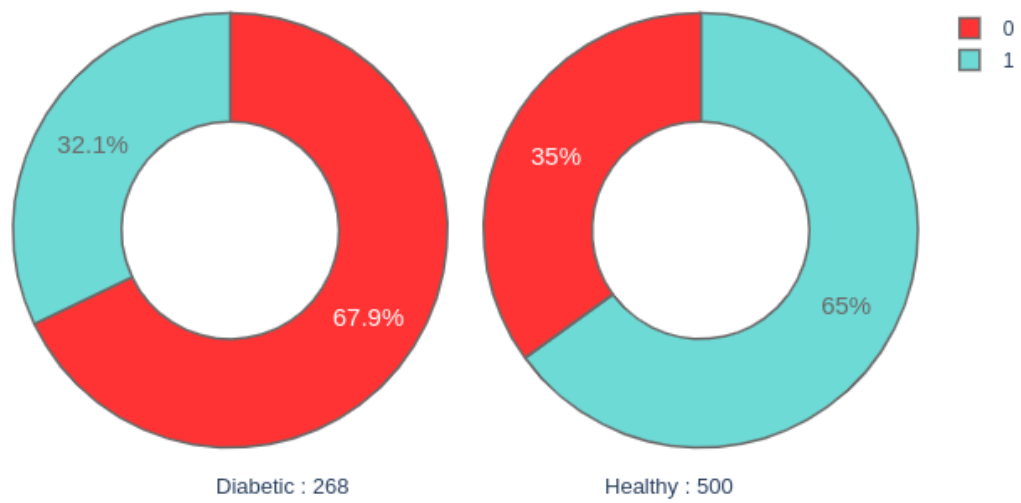


Figure 28: N3 distribution in percentage.

Glucose vs BloodPressure

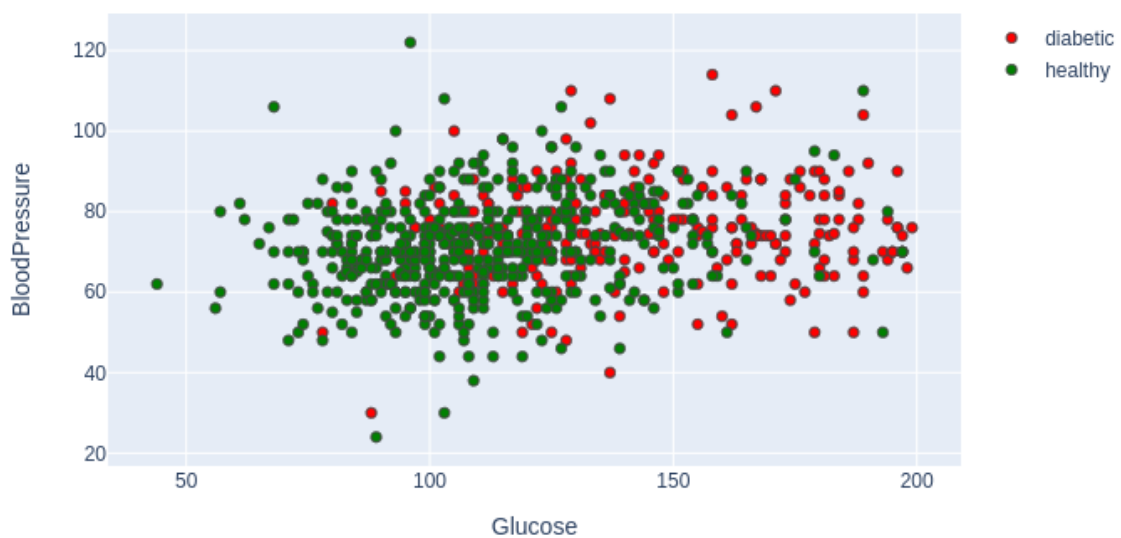


Figure 29: Glucose v/s Blood Pressure scatterplot.

N4 : Glucose \leq 105 and BloodPressure \leq 80

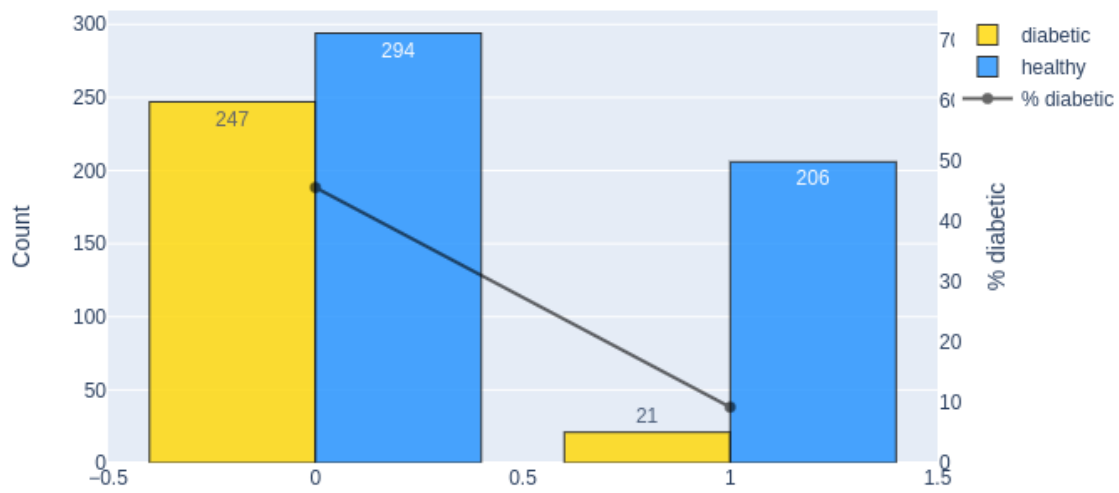


Figure 30: N4 barplot for diabetic and healthy population.

N4 distribution by target
Glucose \leq 105 and BloodPressure \leq 80

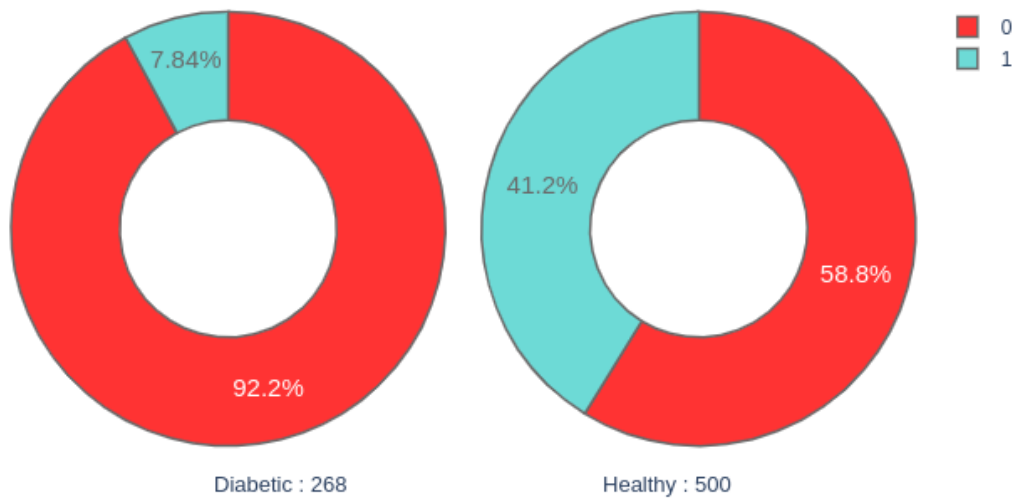


Figure 31: N4 distribution by target.

N5 :SkinThickness <= 20

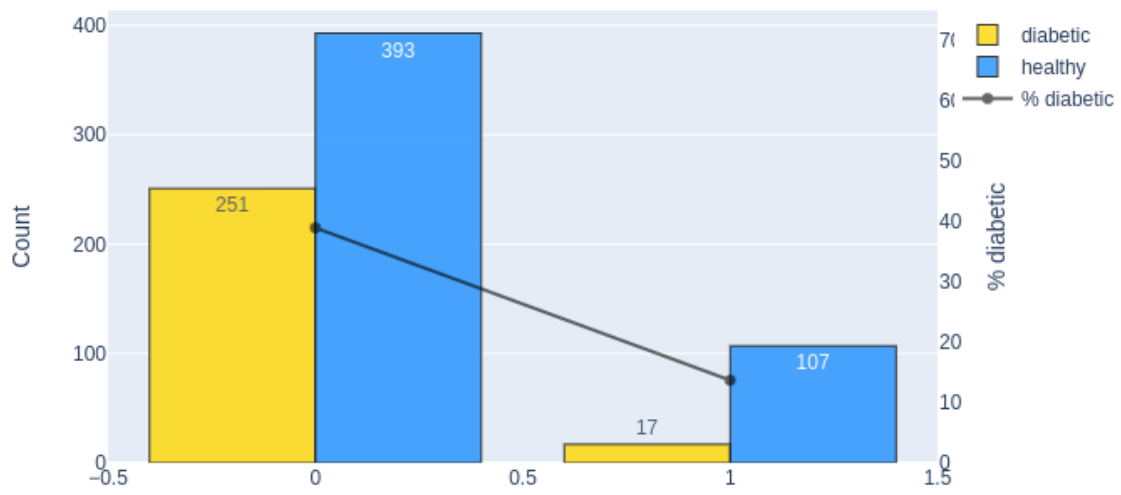


Figure 32: N5 barplot for diabetic and healthy population.

N5 distribution by target
SkinThickness <= 20

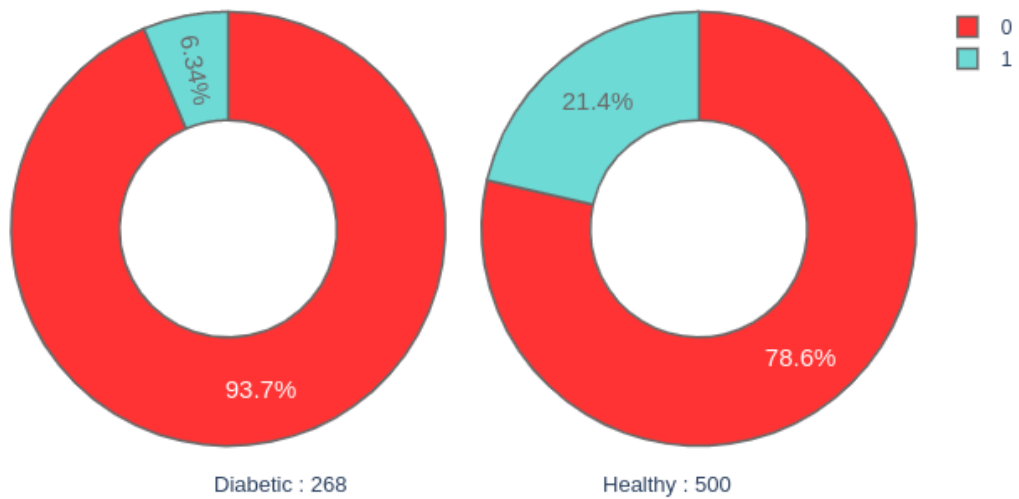


Figure 33: N5 distribution by target.

SkinThickness vs BMI

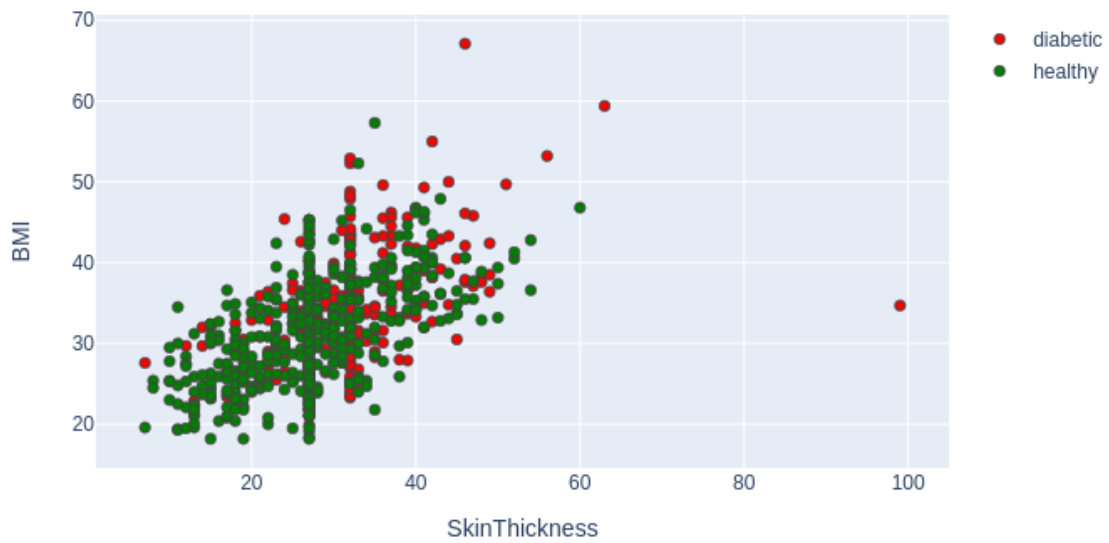


Figure 34: Skin Thickness v/s BMI scatterplot.

N6 : BMI < 30 and SkinThickness <= 20

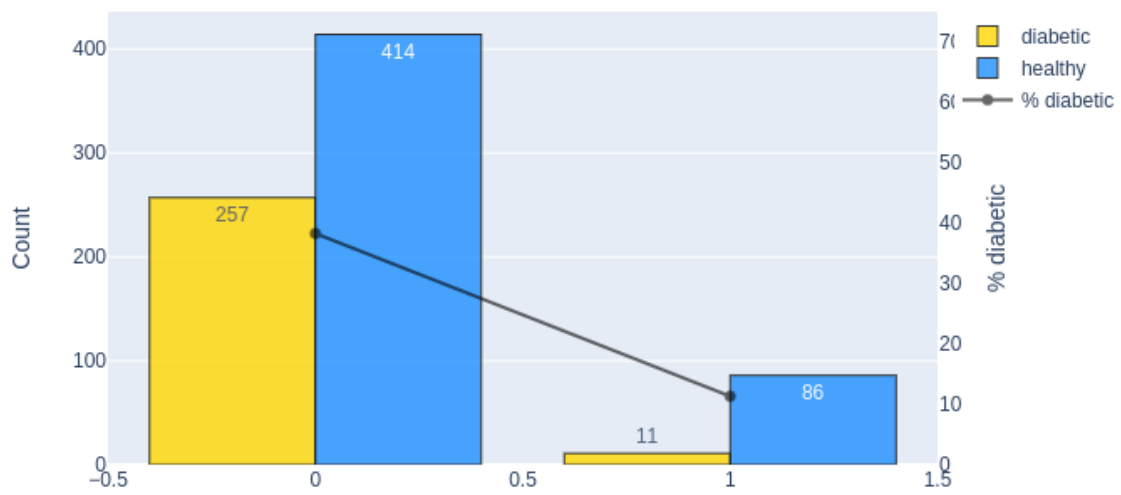


Figure 35: N6 barplot for diabetic and healthy population.

N6 distribution by target
BMI < 30 and SkinThickness <= 20

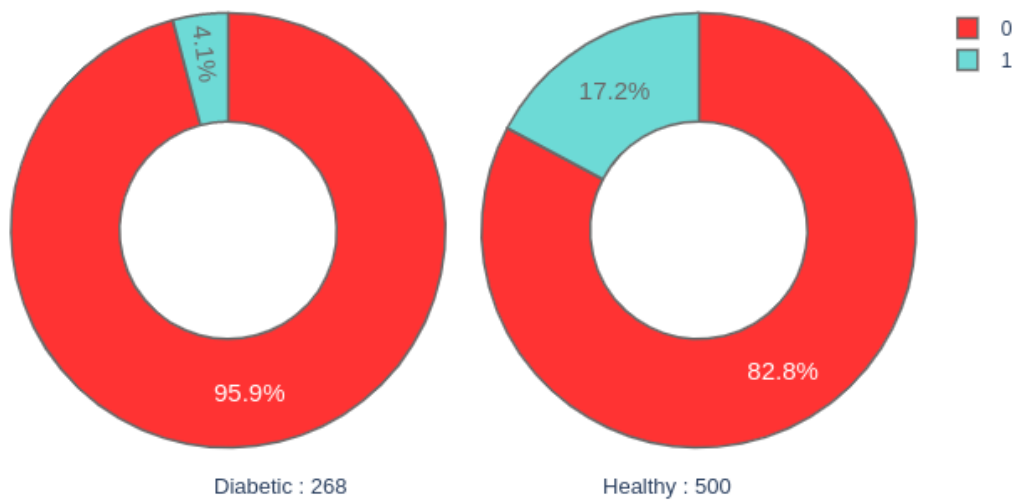


Figure 36: N6 distribution by target.

Glucose vs BMI

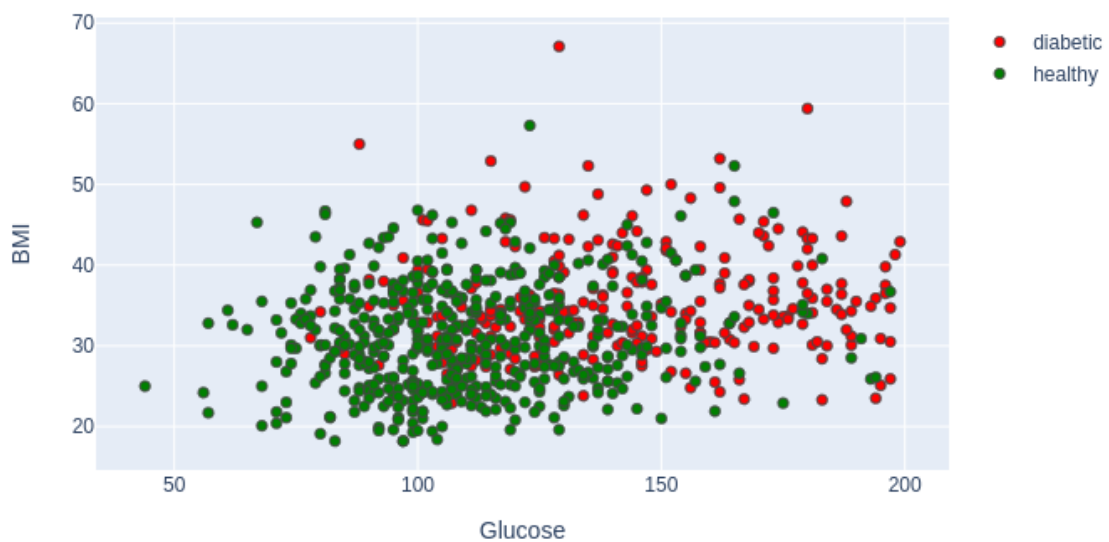


Figure 37: Glucose v/s Body Mass Index scatterplot.

N7 : Glucose \leq 105 and BMI \leq 30

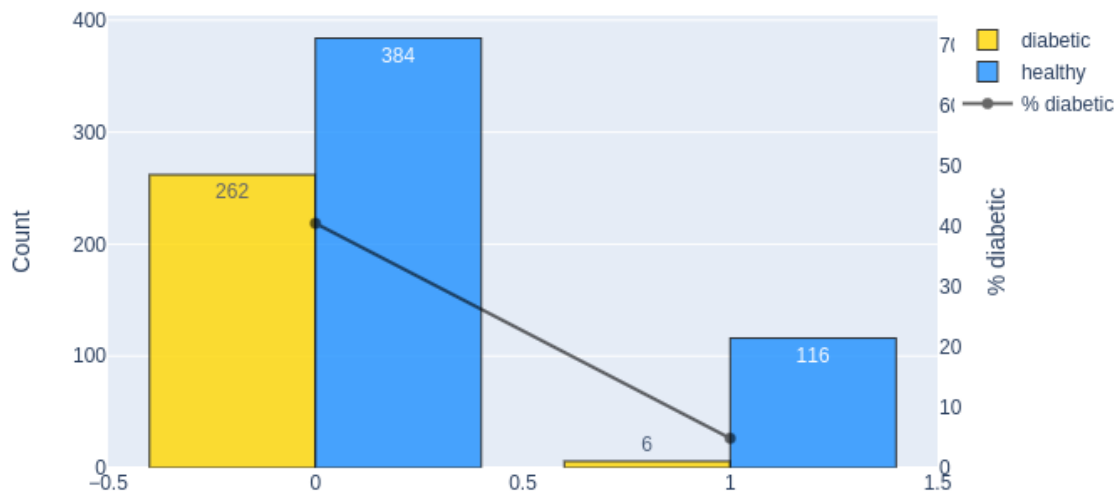


Figure 38: N7 barplot for diabetic and healthy population.

N7 distribution by target
Glucose \leq 105 and BMI \leq 30

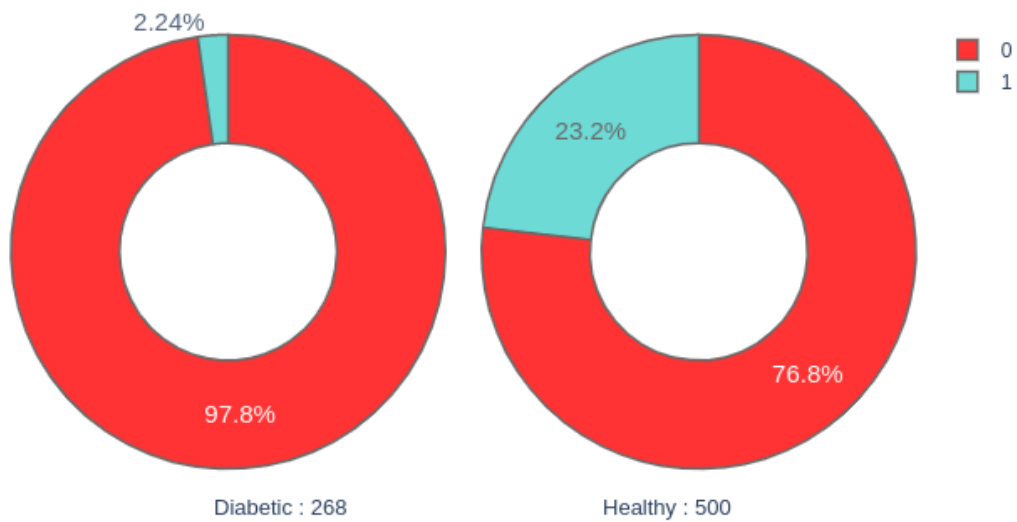


Figure 39: N7 distribution by target.

N9 : Insulin < 200

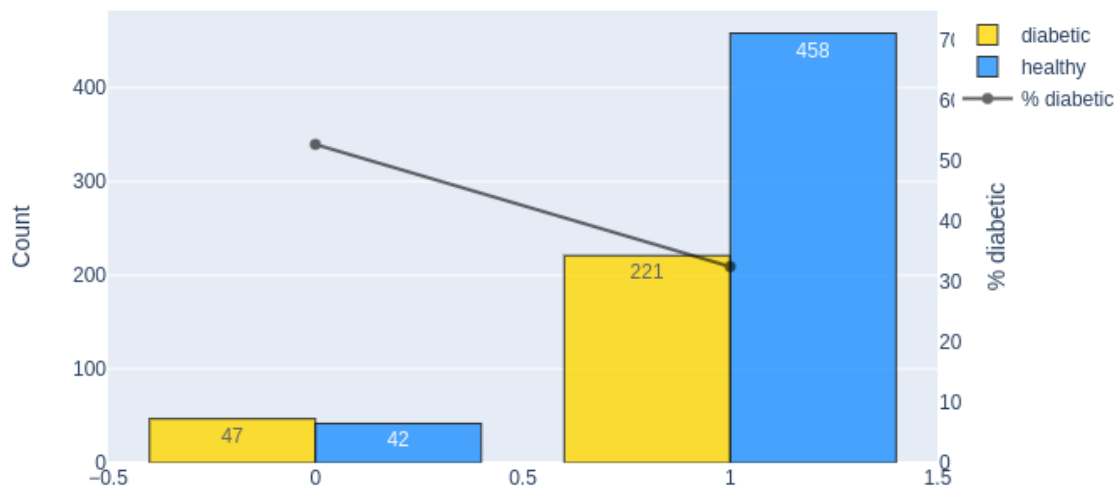


Figure 40: N9 barplot for diabetic and healthy population.

N9 distribution by target
Insulin < 200

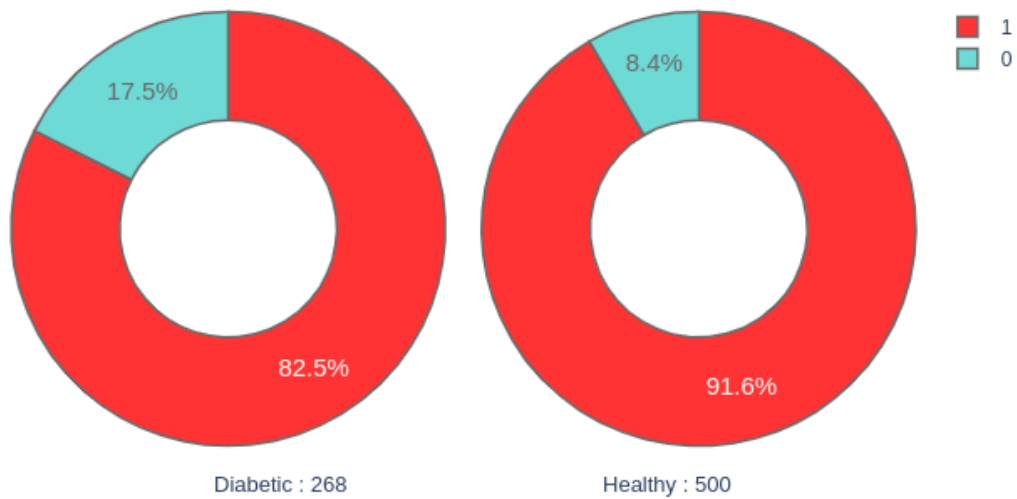


Figure 41: N9 distribution by target.

N10 : BloodPressure < 80

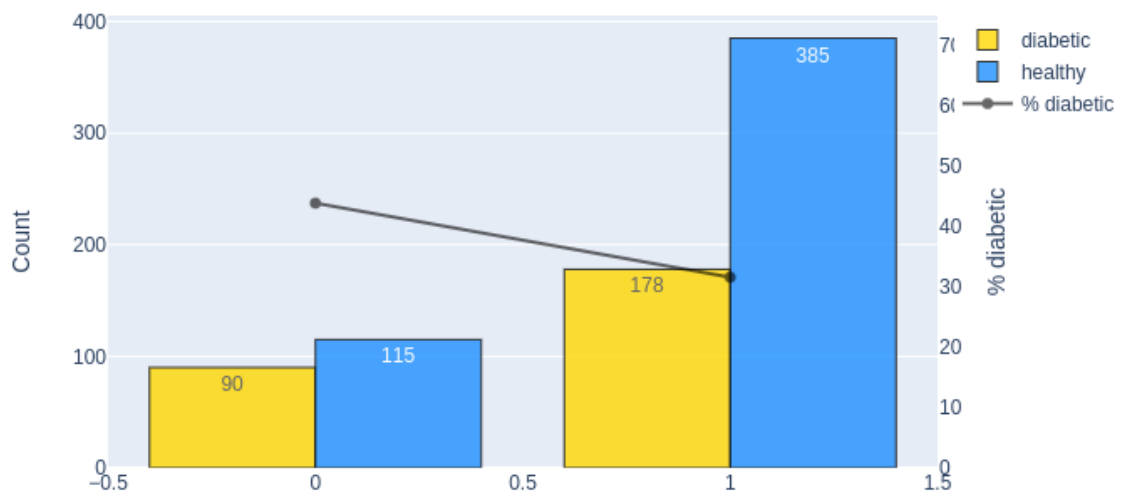


Figure 42: N10 barplot for diabetic and healthy population.

N10 distribution by target
BloodPressure < 80

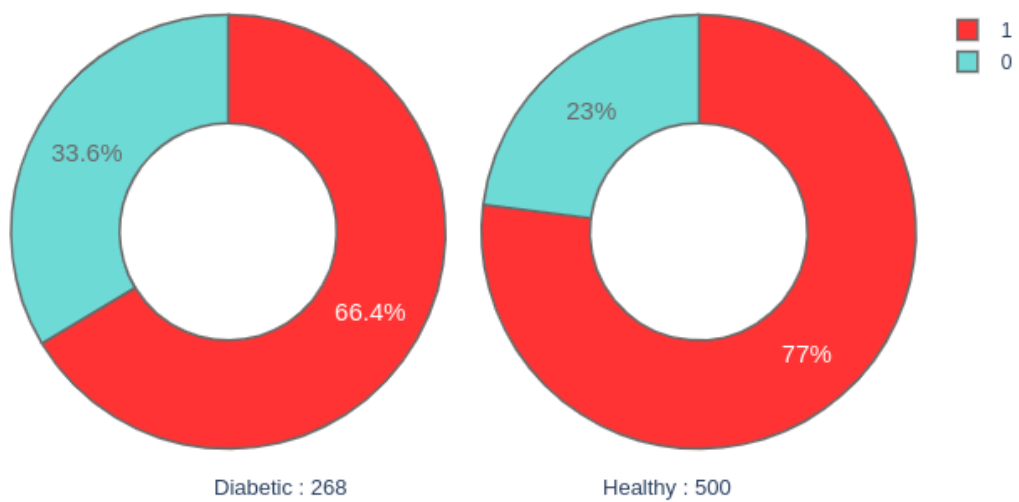


Figure 43: N10 distribution by target.

N11 : Pregnancies > 0 and < 4

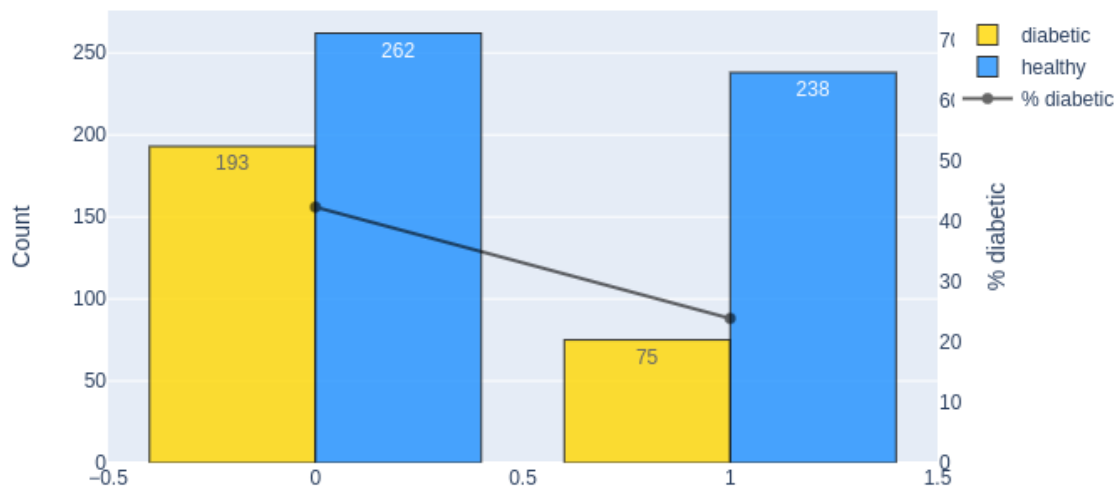


Figure 44: N11 barplot for diabetic and healthy population.

N11 distribution by target
Pregnancies > 0 and < 4

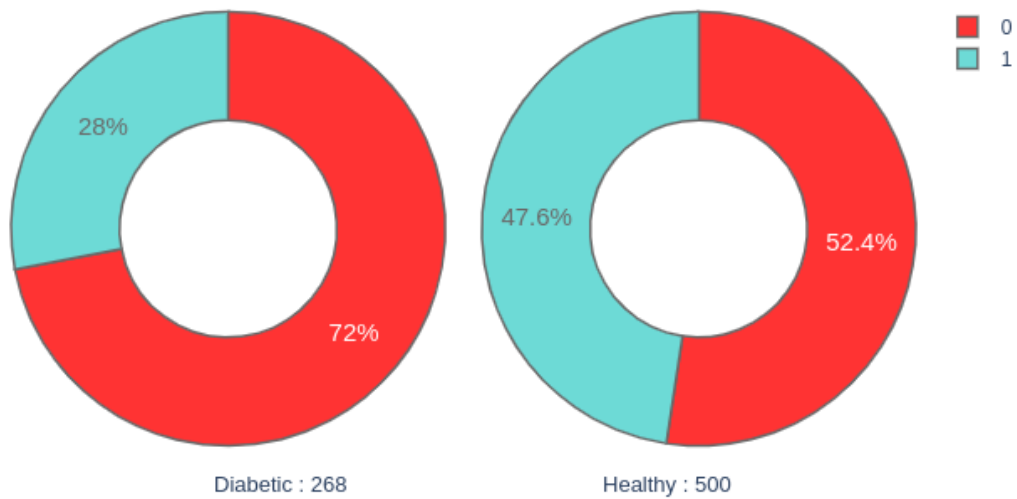


Figure 45: N11 distribution by target.

N15 : NO < 1034

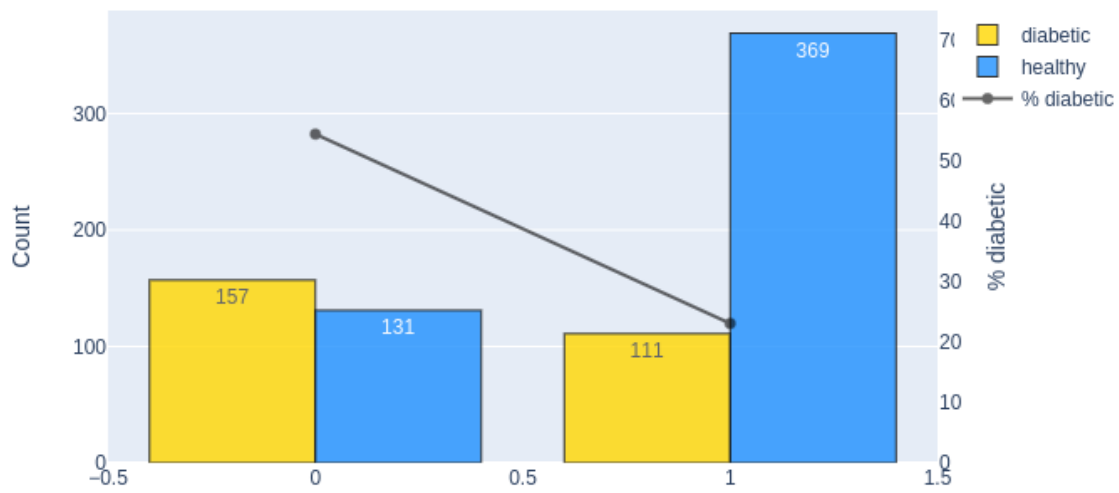


Figure 46: N15 barplot for diabetic and healthy population.

N15 distribution by target
NO < 1034

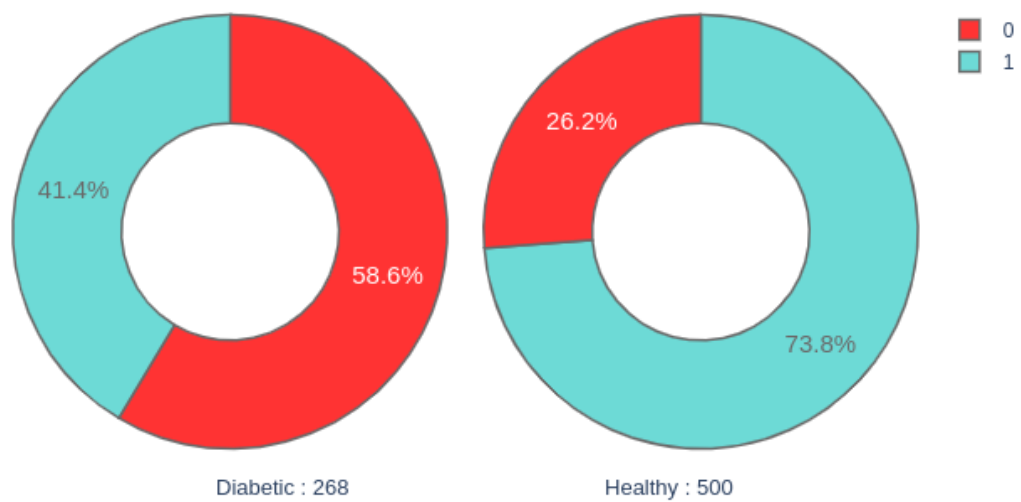


Figure 47: N15 distribution by target.

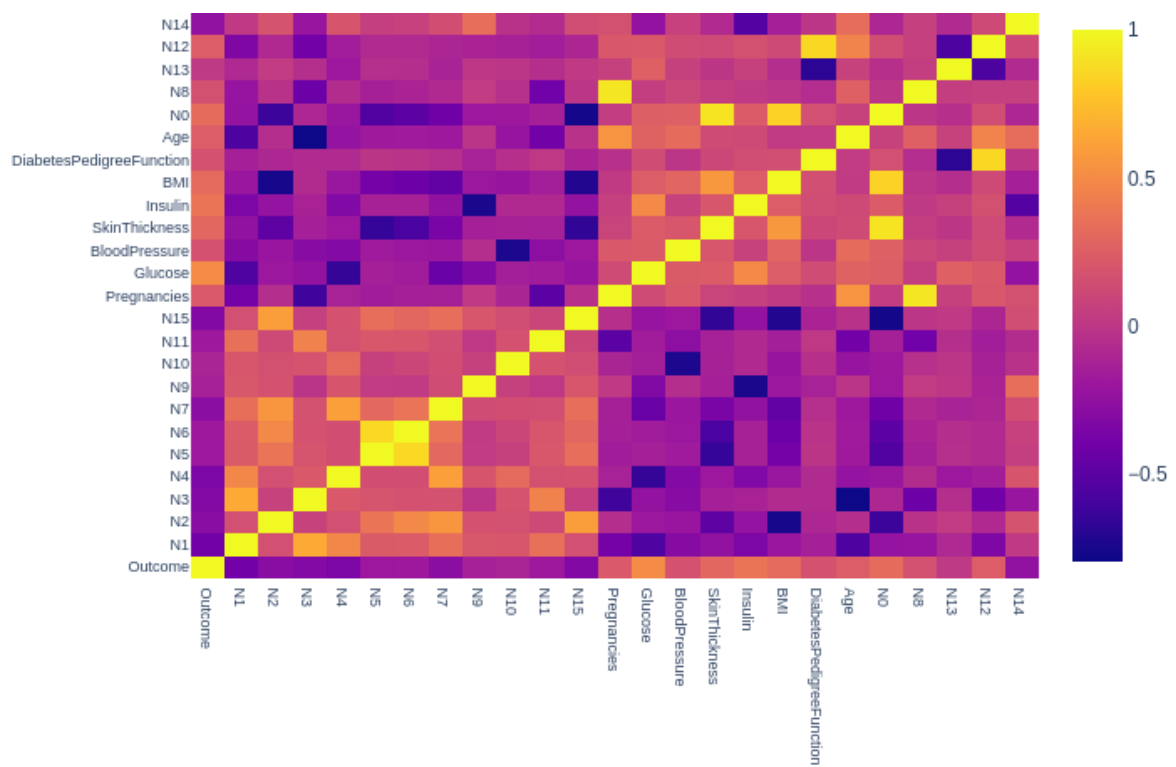


Figure 48: Extended heatmap with new features combined.

14 Conclusion

References

- [1] S. Watson, “Diabetes: Symptoms, causes, treatment, prevention, and more,” *Healthline*, Mar 2019.
- [2] Ministry for Primary Industries, “Kina sea urchin regions in NZ.” <http://fs.fish.govt.nz/Page.aspx?pk=7&sc=SUR>, 2013. Online; accessed 29 January 2014.
- [3] S. P. Chatrati, G. Hossain, A. Goyal, A. Bhan, S. Bhattacharya, D. Gaurav, and S. M. Tiwari, “Smart home health monitoring system for predicting type 2 diabetes and hypertension,” *Journal of King Saud University - Computer and Information Sciences*, Jan 2020.
- [4] D. Deva, “Qcon rio - machine learning for everyone,” *LinkedIn SlideShare*, Aug 2015.
- [5]
- [6] P. Kaur and R. Kaur, “Comparative analysis of classification techniques for diagnosis of diabetes,” *SpringerLink*, Jan 1970.
- [7] K. Alberti and P. Zimmet, “Definition, diagnosis and classification of diabetes mellitus and its complications. part 1: diagnosis and classification of diabetes mellitus. provisional report of a who consultation,” Jul 2004.
- [8] L. G. Lisa, L. Gladman, and Lisa, “Types of healthcare information systems - scott-clark medical,” *Scott*, Nov 2019.
- [9] K. Maladkar, “Why is random search better than grid search for machine learning,” *Analytics India Magazine*, Sep 2019.
- [10] A. Mandal, “What is diabetes?,” *News*, Feb 2019.
- [11] P. Mandot, “What is lightgbm, how to implement it? how to fine tune the parameters?,” *Medium*, Dec 2018.
- [12] S. Norena, “Python model tuning methods using cross validation and grid search,” *Medium*, Jun 2018.
- [13] “sklearn.preprocessing.labelencoder,” *scikit*.
- [14]