

18CSC301J-INFORMATION VISUALIZATION

SEMESTER VI

YEAR: 21-22

Diabetes predictions and its contributors

Report Submitted by

S.AJAY [RA1911003010102]

TARUUN.P [RA1911003010103]

SAIPAVAN TIPPANU [RA1911003010121]

S.KARTHIKEYAN [RA1911003010120]

Faculty in-charge

Debajyoti Modal

Kavitha



SRM
INSTITUTE OF SCIENCE & TECHNOLOGY
Deemed to be University u/s 3 of UGC Act, 1956

**DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING
COLLEGE OF ENGINEERING AND TECHNOLOGY**

SRM INSTITUTE OF SCIENCE AND TECHNOLOGY

S.R.M. Nagar, Kattankulathur – 603203, Kancheepuram District

Agenda:

To take a dataset of people with diabetes and see the factors that contribute to diabetes and other ailments.

Abstract:

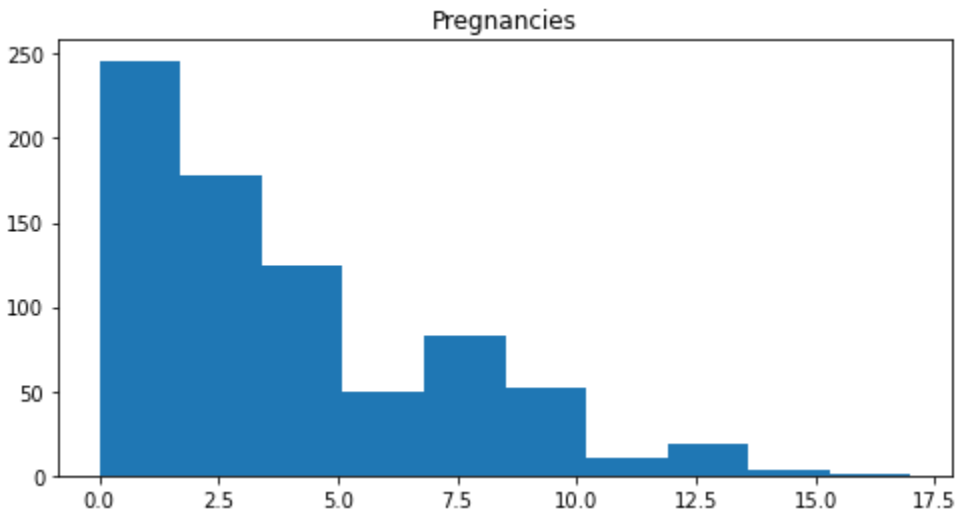
Diabetes is a prolonged infection or set of metabolic ailments where somebody suffers from a prolonged level of Blood Glucose(BG) in the body, due to deficiency of, or due to the non-reaction of the cells and do not react appropriately to insulin. Nowadays, this disease is leading to long-term impediments and stern health problems. Healthcare industry contains huge volumes of highly sensitive information that has to be dealt with appropriately. Diabetes Mellitus(DM) is considered to be a deadly sickness on the planet. Clinical experts need a reliable framework of expectations for the analysis of diabetes.

Introduction:

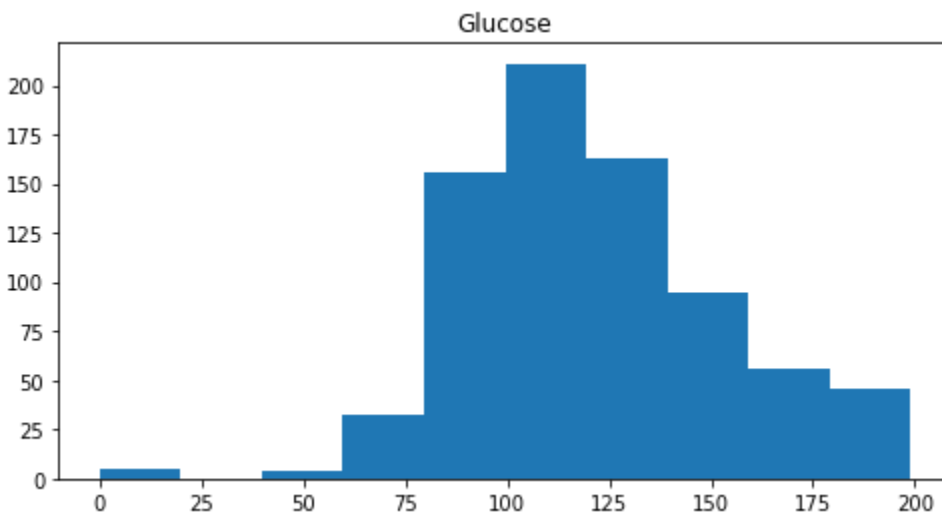
In India, diabetes is a major issue. Between 1971 and 2000, the incidence of diabetes rose ten times, from 1.2% to 12.1%. 61.3 million people 20–79 years of age in India are estimated living with diabetes (Expectations of 2011). It is expected that by 2030 this number will rise to 101.2 million. There are a few AI strategies that are utilized to perform prescient analysis of big data in a variety of domains. Prognostic health examination is a daunting work, but it can ultimately aid experts to make timely and well-versed judgments about the health and handling of patients. The core objective of the study is to assist physicians and experts in the primary diagnosis of diabetes by means of ML methods and visualizations.

OBSERVATION:

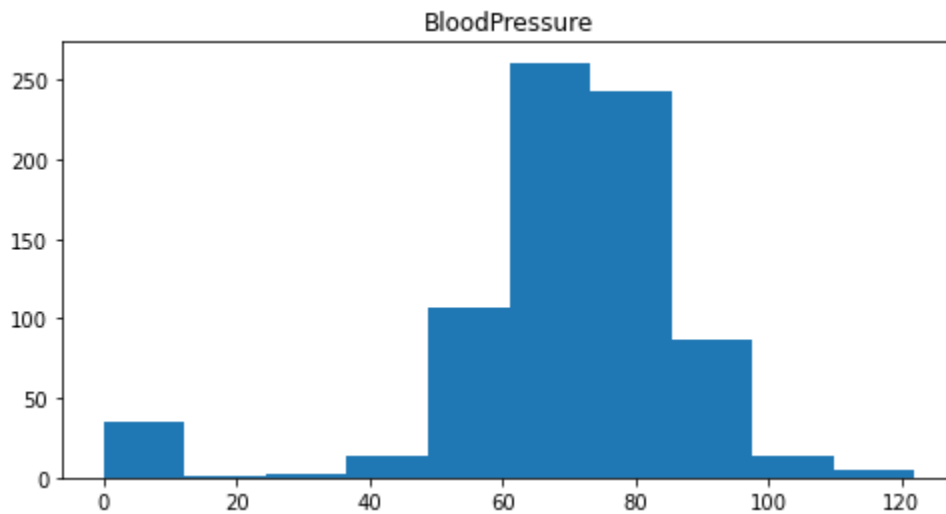
- **UNIVARIATE ANALYSIS**



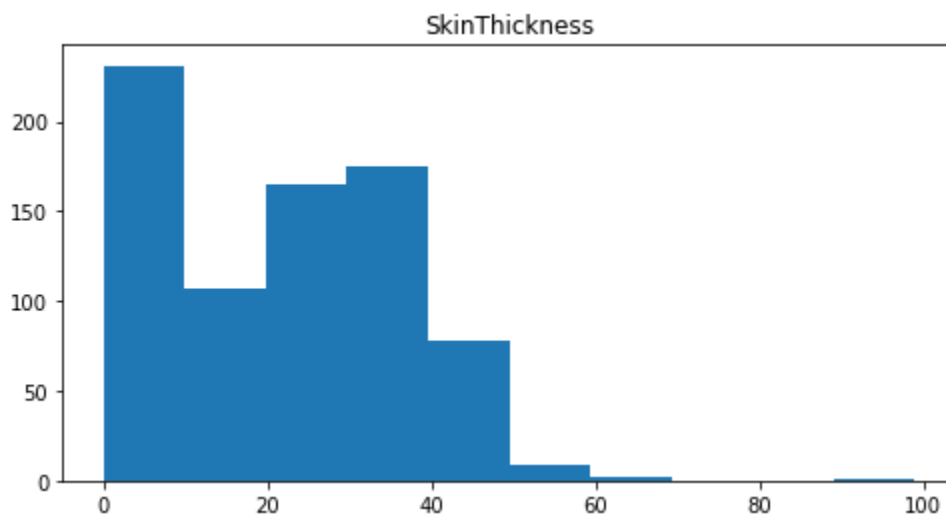
This depicts the number of pregnancies that are there amongst the dataset. Here, we can clearly see the number of pregnancies rates. i.e How many people have a certain number of pregnancy periods.



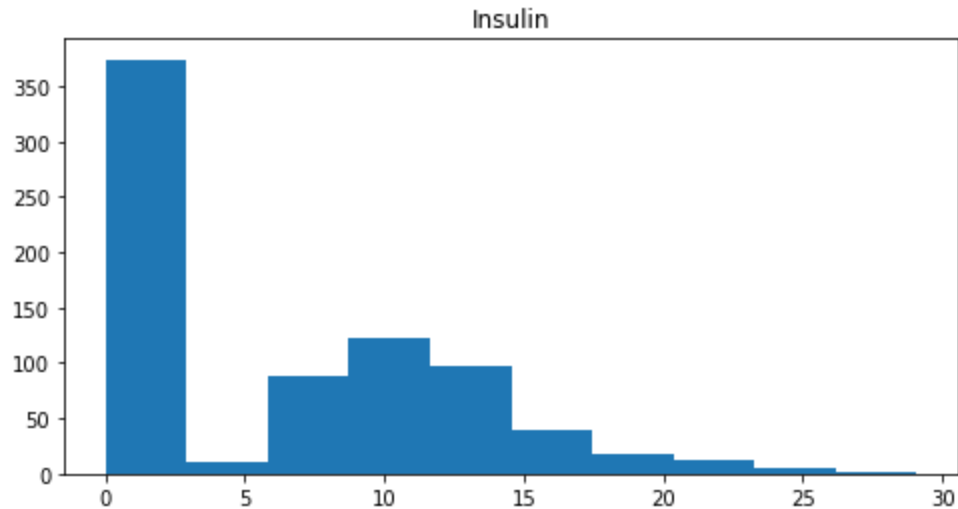
This graph depicts the glucose levels of the people pertaining to the dataset. There are many people with high glucose levels in the middle sector of the graph.



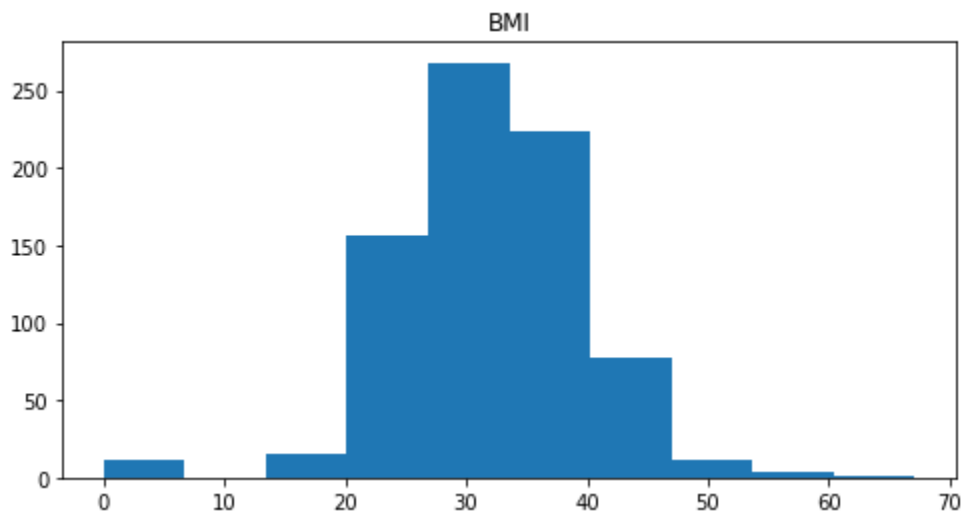
This graph depicts the blood pressure levels of the people pertaining to the dataset. Here, we can see that many people have BP within the range of 60-85 and a smaller number of people have less than 40.



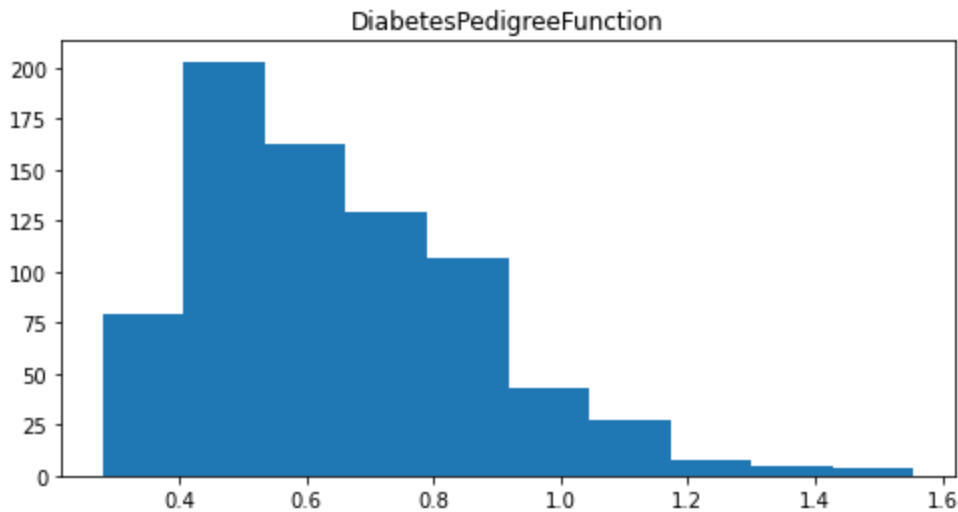
This graph depicts the skin thickness of the people pertaining to the dataset. Here we can see that many people have skin thickness ranging from 0-40 and very minimal people have it in the range of 50+.



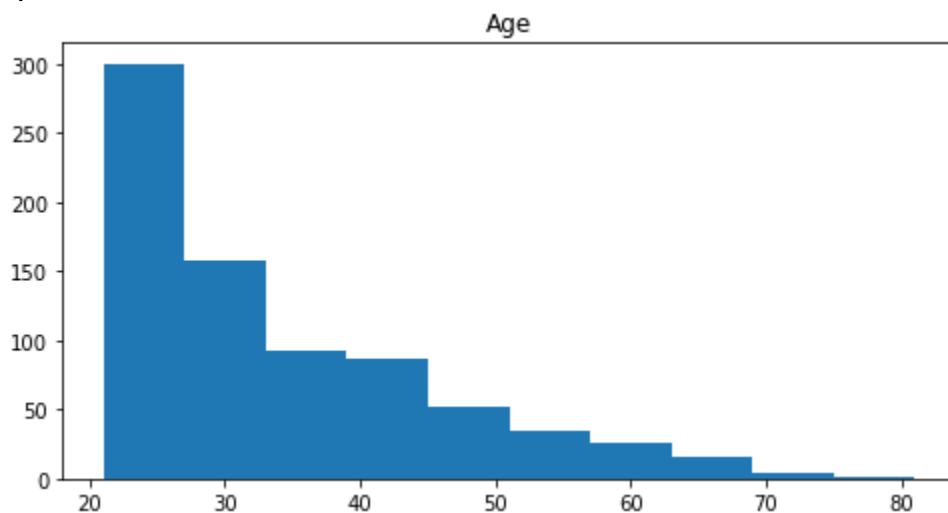
This graph depicts the insulin levels of the people pertaining to the dataset. Here we can see that people have high insulin ranging from 0-3 and starts to gradually decrease from its second peak at 10.



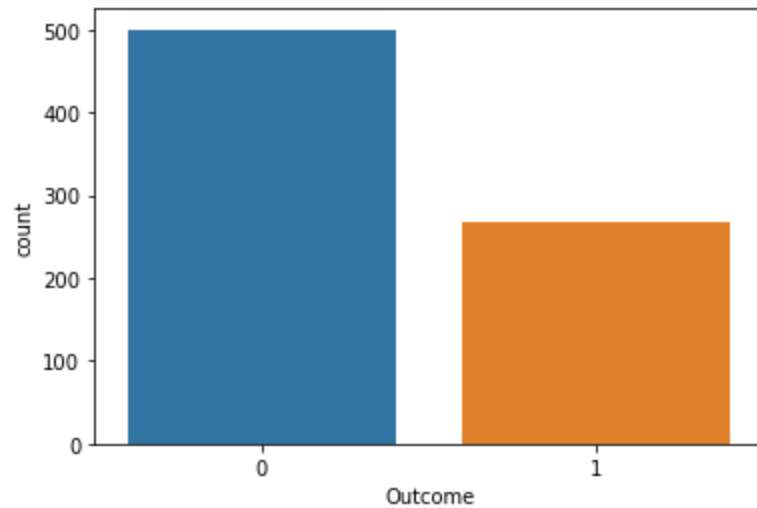
This graph depicts the BMI of the people pertaining to the dataset. Here we can see that bmi peaks out at 30 and is high in the range of 20-45



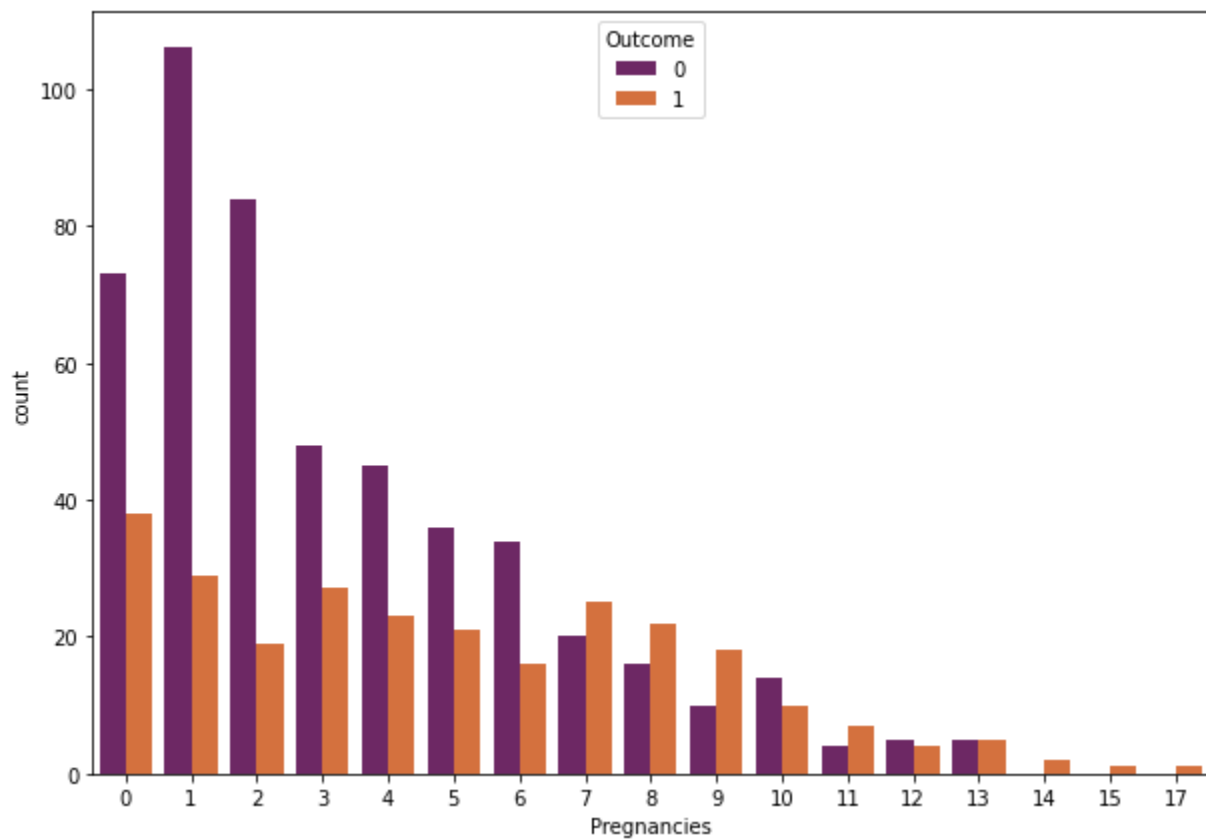
This graph depicts the Diabetes Pedigree function of the people pertaining to the dataset. Here we can see the graph peaking at 0.5 and gradually decreases upon further examination.



This graph depicts the Age of the people pertaining to the dataset. Here we can see the graph peaking at 0.5 and gradually decreases upon further examination.

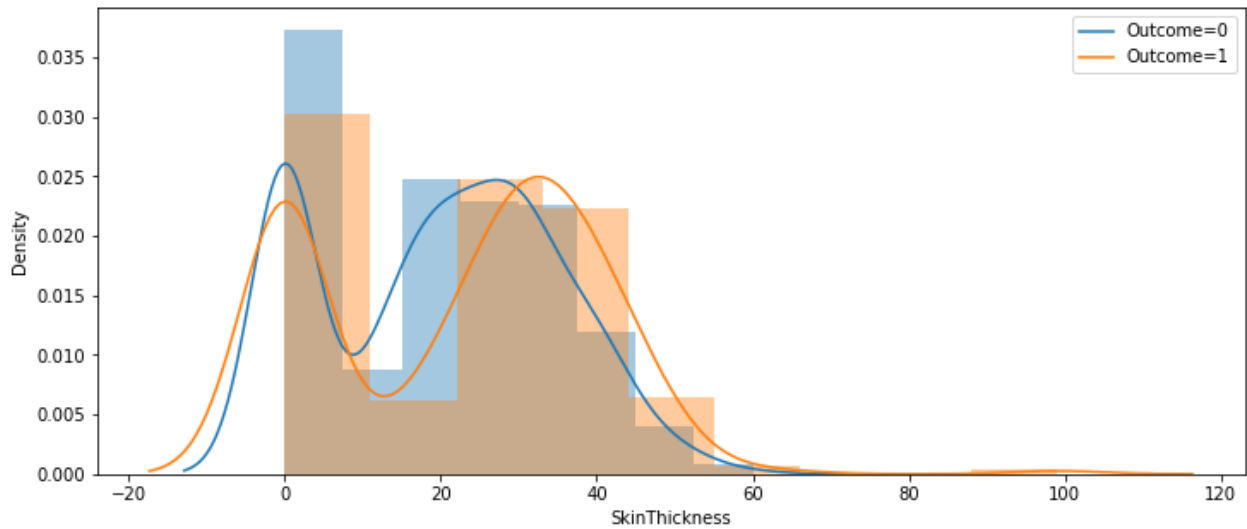


In this graph we can observe the outcome of the univariate analysis based on the previous graph's data as input. This graph can be deciphered by taking the orange as has diabetes and blue as does not have diabetes

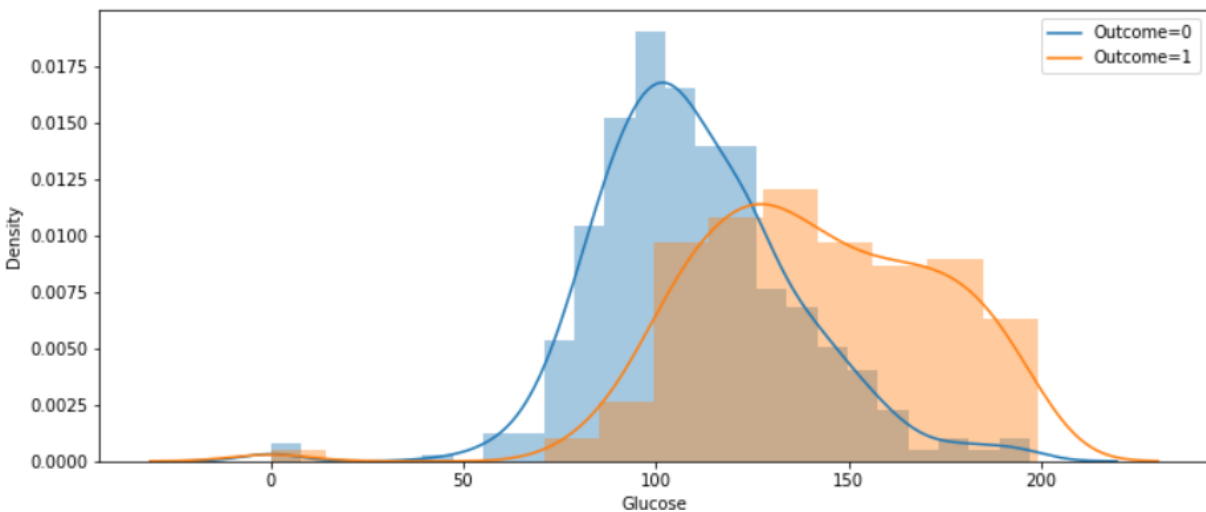


Via this viz, we can see that pregnancies might not always contribute to diabetes as the outcome of most of the pregnant women are not diabetic. This shows us that this might not only be the main contributing factor if the patient is diabetic.

- **BIVARIATE ANALYSIS**

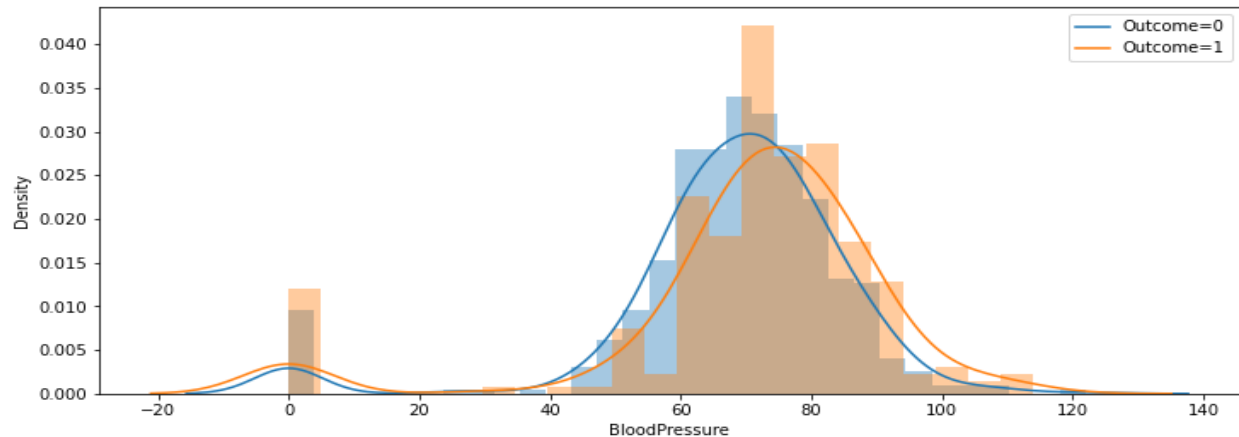


In this Graph we are able to grasp the correlation between thickness of skin and having diabetes. The graph also shows that, the outcome of diabetes is nearly the same for all the people of different skin thickness. Only in the vicinity of 0.025 density of 0 skin thickness and around 10-20, there is a clear correlation where we can see that the outcome is 0. And from around 30-60, the outcome is 1 where, we can conclude that this might be a contributing factor for diabetes.

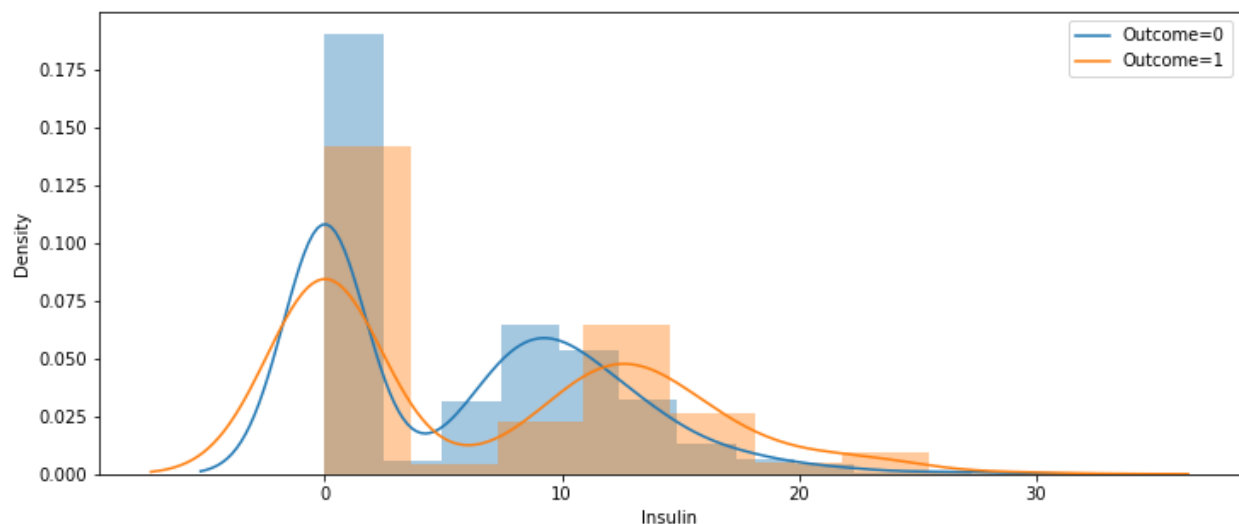


In this Graph we are able to grasp the correlation between Glucose and having Diabetes. The graph depicts the glucose level is directly proportional to the chances of diabetes. We can clearly see the graphs outcome around 100-250,

the outcome 1. where, we can conclude that this might be a contributing factor for diabetes.

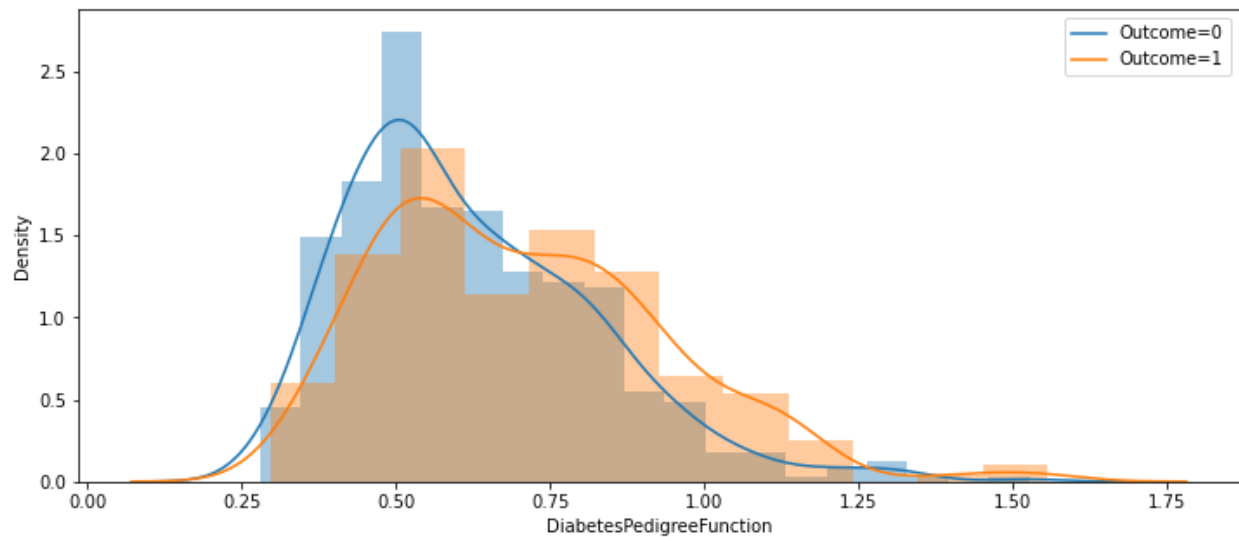


In this Graph we are able to see the correlation between blood pressure and having diabetes. The graph hasn't shown much variance till 40 but from 40-120 there is a slight variance. The dataset shows that people who have a bp higher than 80 have a higher chance of diabetes.

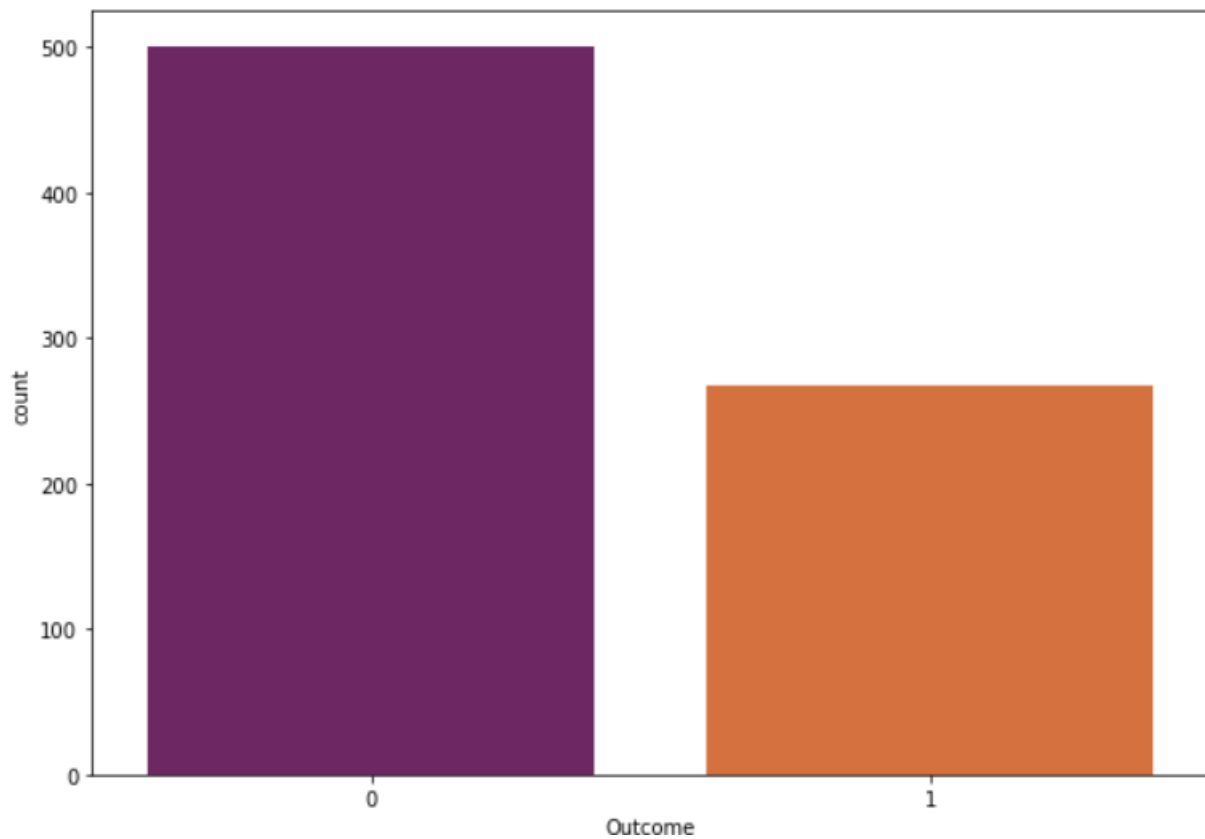


In this Graph we are able to grasp the correlation between Insulin and having Diabetes. The variance in insulin levels play a crucial role in telling a person if he has diabetes or not. As we can see from the graph the region around 5-12 has the outcome of 0 but from 15 and higher the outcome is 1. So higher or lower the

insulin level, the more likely a person has diabetes. ie. having less or excess amount of insulin might lead to diabetes.

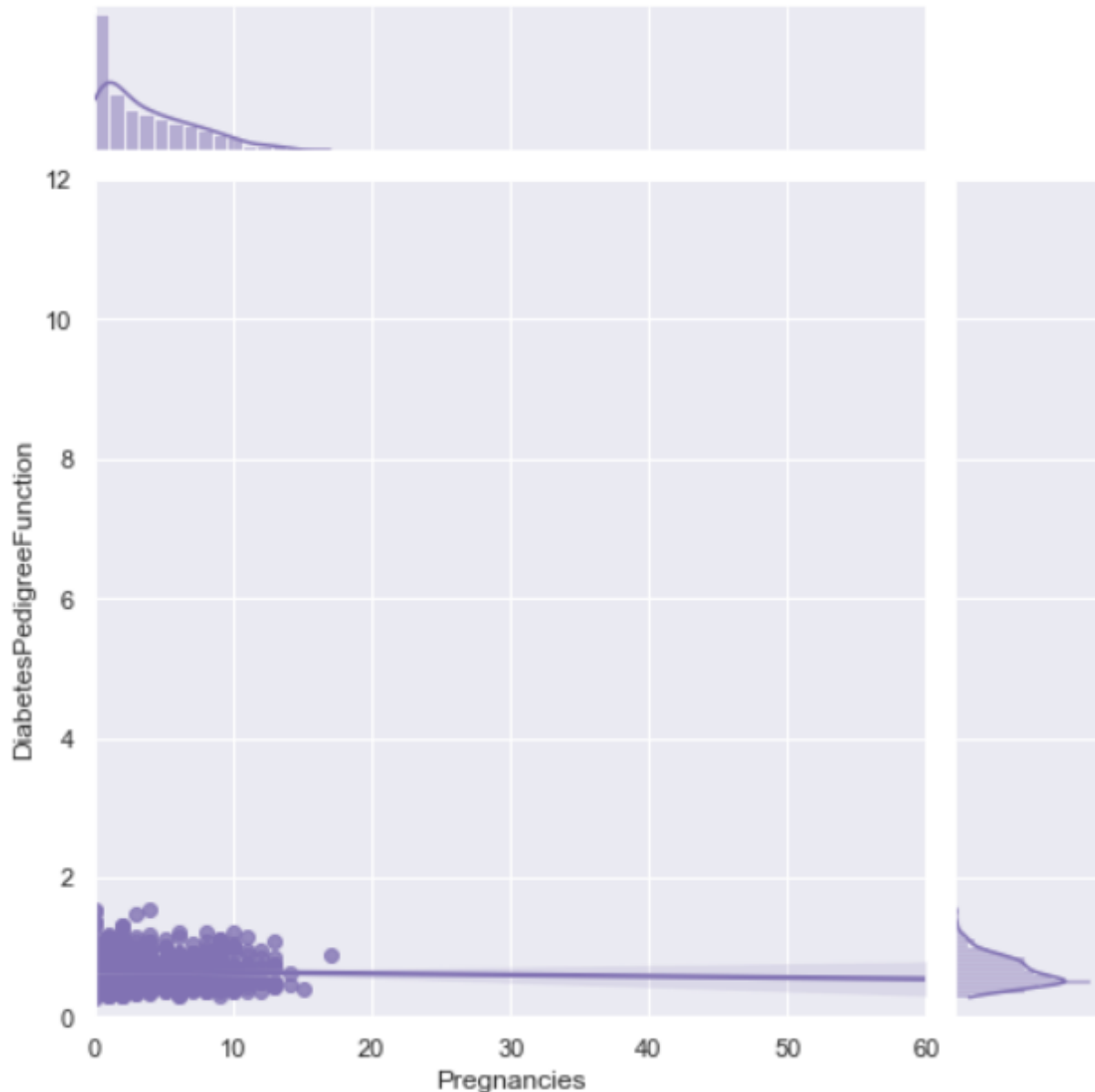


In this Graph we are able to grasp the correlation between Diabetes and having diabetes. From the given graph we can see from the graph around the region 0.25-0.75 the outcome is 0 but as the numbers get higher ie. from 0.75 - 1,25 the outcome is 1. We can conclude that with the increase in diabetes pedigree function, the outcome of diabetes increases.

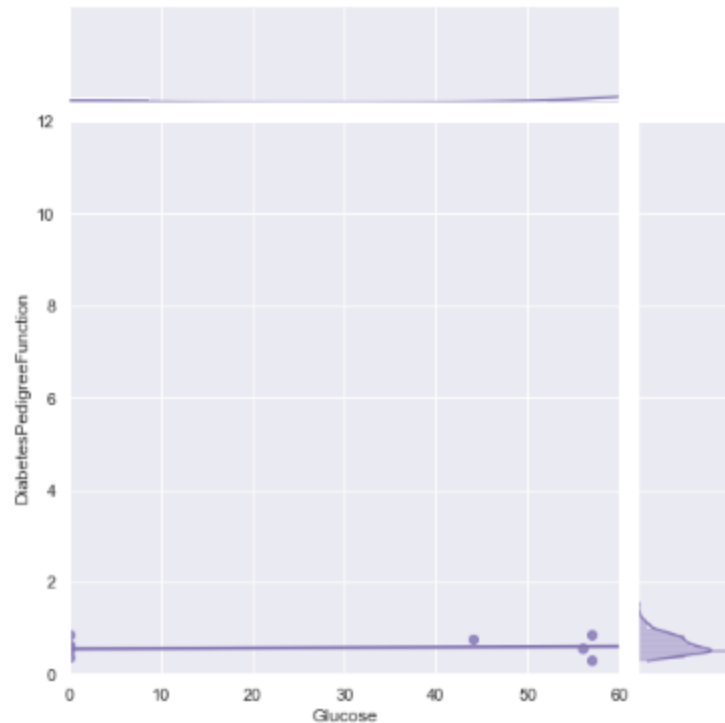


In this graph we can observe the outcome of the Bivariate analysis based on the previous graph's data as input. This graph can be deciphered by taking the orange as has diabetes and purple as does not have diabetes.

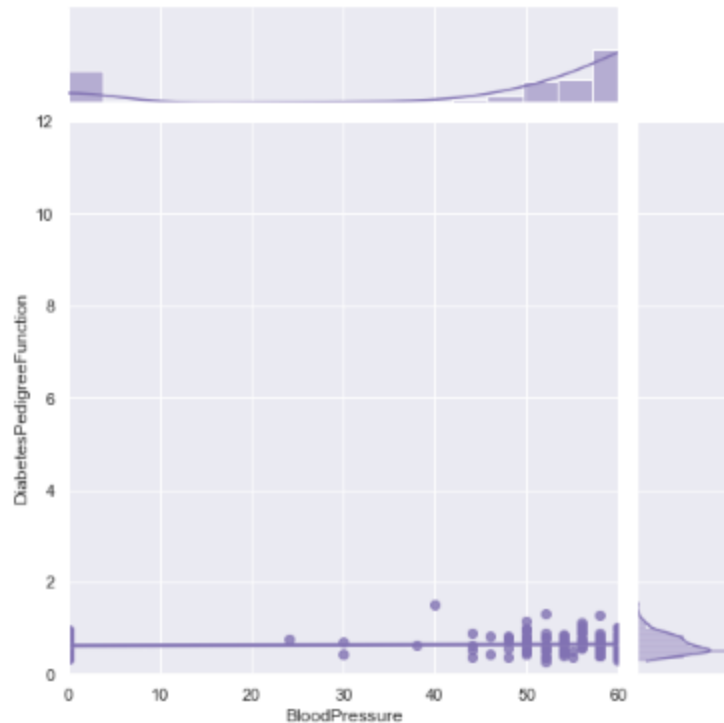
- **Linear Regression With Marginal Distribution**



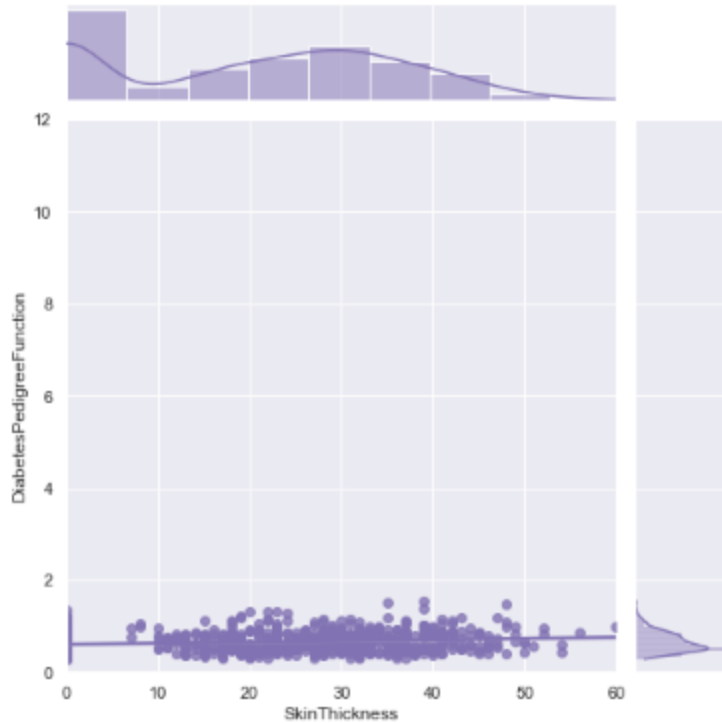
In this graph we are able to perceive the correlation between Diabetes Pedigree Function and Pregnancies. The number of pregnancies doesn't matter in this case as the diabetes pedigree function stays constant throughout. On the side we have the marginal side plots, wherein we can see the marginal distribution pertaining to diabetes pedigree function and pregnancies. It increases at the start and slopes down as the density decreases.



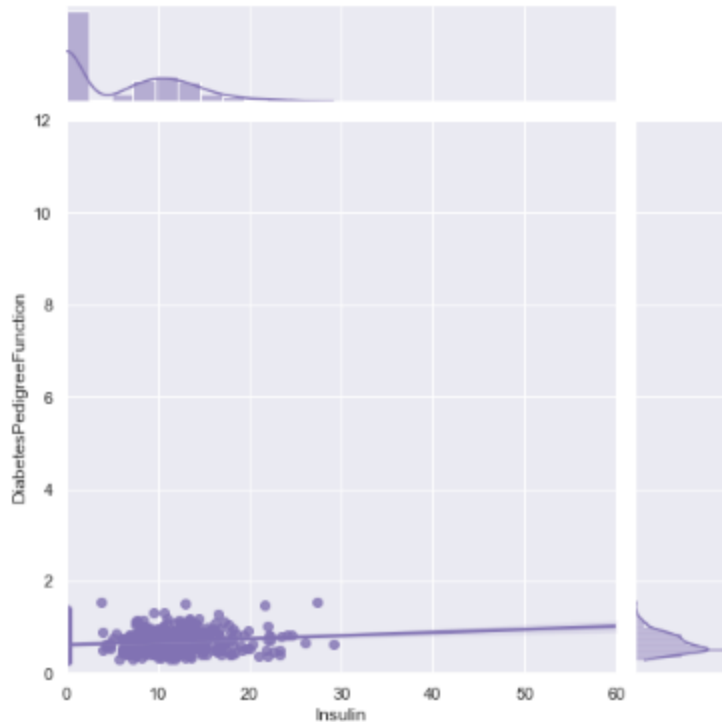
In this graph we are able to perceive the correlation between Diabetes Pedigree Function and Glucose. The factor of glucose doesn't contribute much in this case as the diabetes pedigree function stays constant throughout. On the side we have the marginal side plots, wherein we can see the marginal distribution pertaining to diabetes pedigree function and glucose. In this only at the end of the graph does the density and the graph have some minimal density and value because of which we can say that this is not that big of a factor of diabetes.



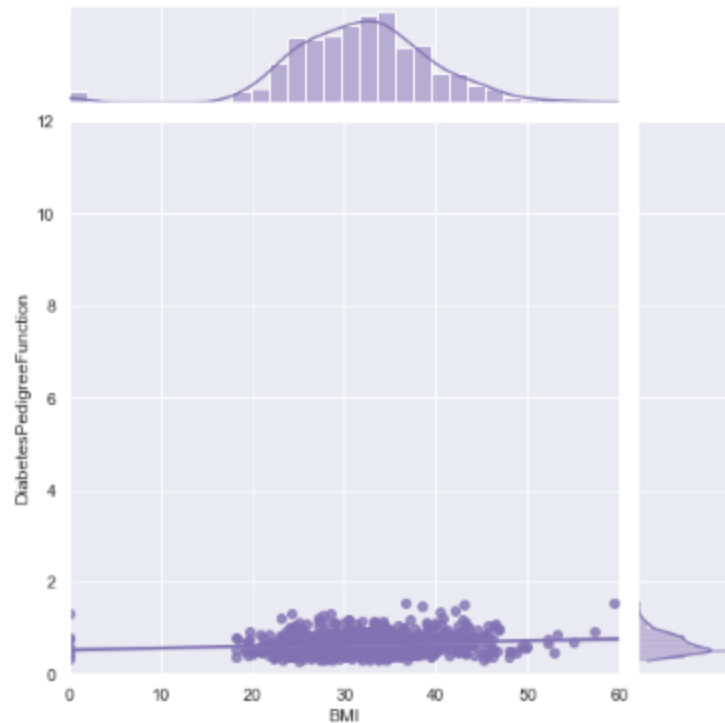
In this graph we are able to perceive the correlation between Diabetes Pedigree Function and Blood Pressure. The factor of Blood pressure doesn't contribute much as in this case as the diabetes pedigree function stays constant throughout at the beginning. On the side we have the marginal side plots, wherein we can see the marginal distribution pertaining to diabetes pedigree function and Blood Pressure. It starts to increase in the middle and spike up as the density increases. So, we can conclude that only above a certain margin, diabetes pedigree function might matter.



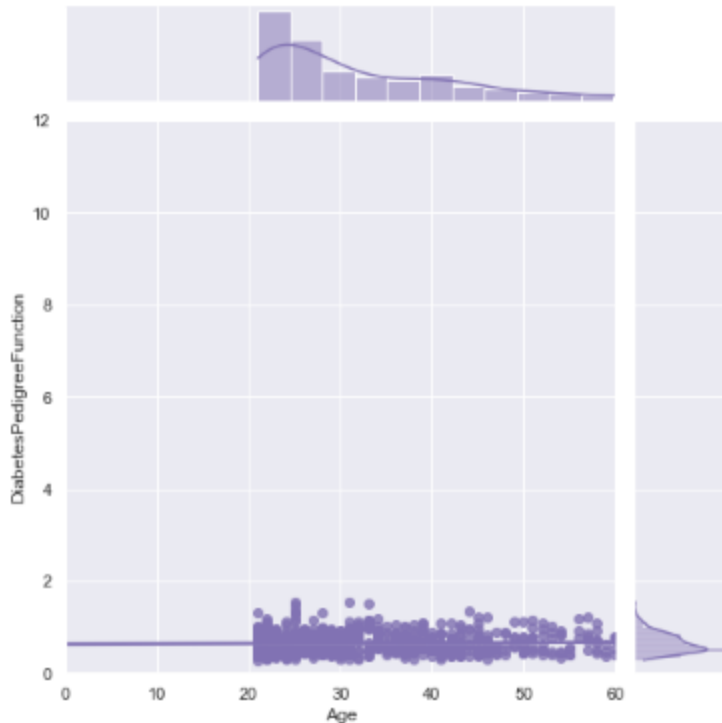
In this graph we are able to perceive the correlation between Diabetes Pedigree Function and Skin Thickness. The factor of Skin thickness contributes a lot as in this case as the diabetes pedigree function stays in the middle throughout and has a lot of density. On the side we have the marginal side plots, wherein we can see the marginal distribution pertaining to diabetes pedigree function and Skin Thickness. It increases at the start and slopes down a little then shoots up again as the density increases and gradually slopes down as the density decreases.



In this graph we are able to perceive the correlation between Diabetes Pedigree Function and Insulin. The factor of Insulin contributes initially as in this case as the diabetes pedigree function varies throughout at start. On the side we have the marginal side plots, wherein we can see the marginal distribution pertaining to diabetes pedigree function and Insulin. It increases at the start and slopes down a little then shoots up again as the density increases and gradually slopes down as the density decreases and flatlines. We can conclude this by saying that, a lot of people take in insulin when they have diabetes, wherein Diabetes pedigree function might be a factor.

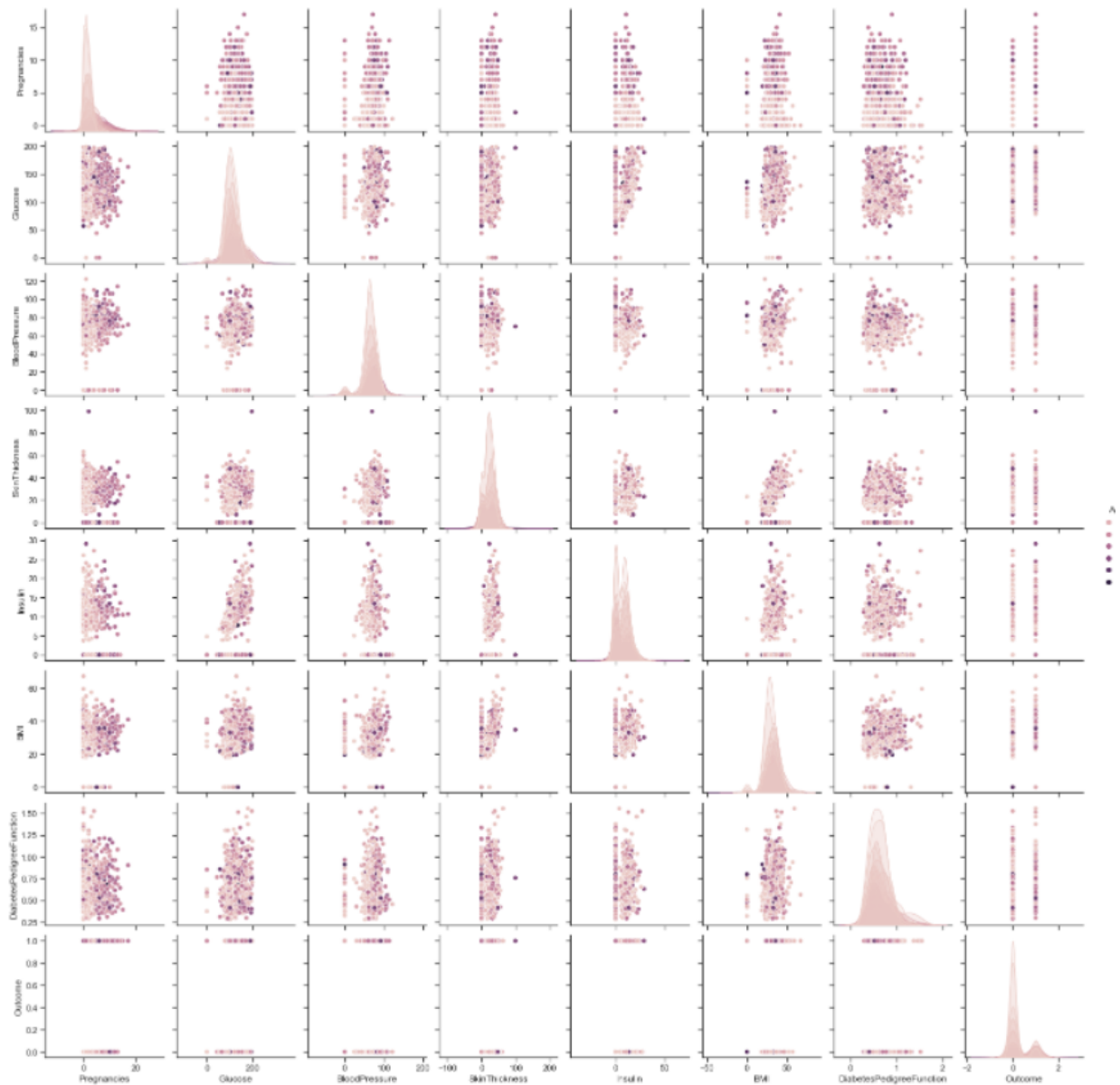


In this graph we are able to perceive the correlation between Diabetes Pedigree Function and BMI. The factor of BMI contributes a lot in this case as the diabetes pedigree function stays high throughout the middle sector. On the side we have the marginal side plots, wherein we can see the marginal distribution pertaining to diabetes pedigree function and BMI. In this we see a small rise in the beginning that flatlines immediately till we reach the middle where the graph peaks when it is having the maximum density and gradually slopes down as the density decreases. So, the BMI of diabetic people does vary between 20-50 heavily where the diabetes pedigree function might have a role in it as most of the people have at most 1 as the function value.



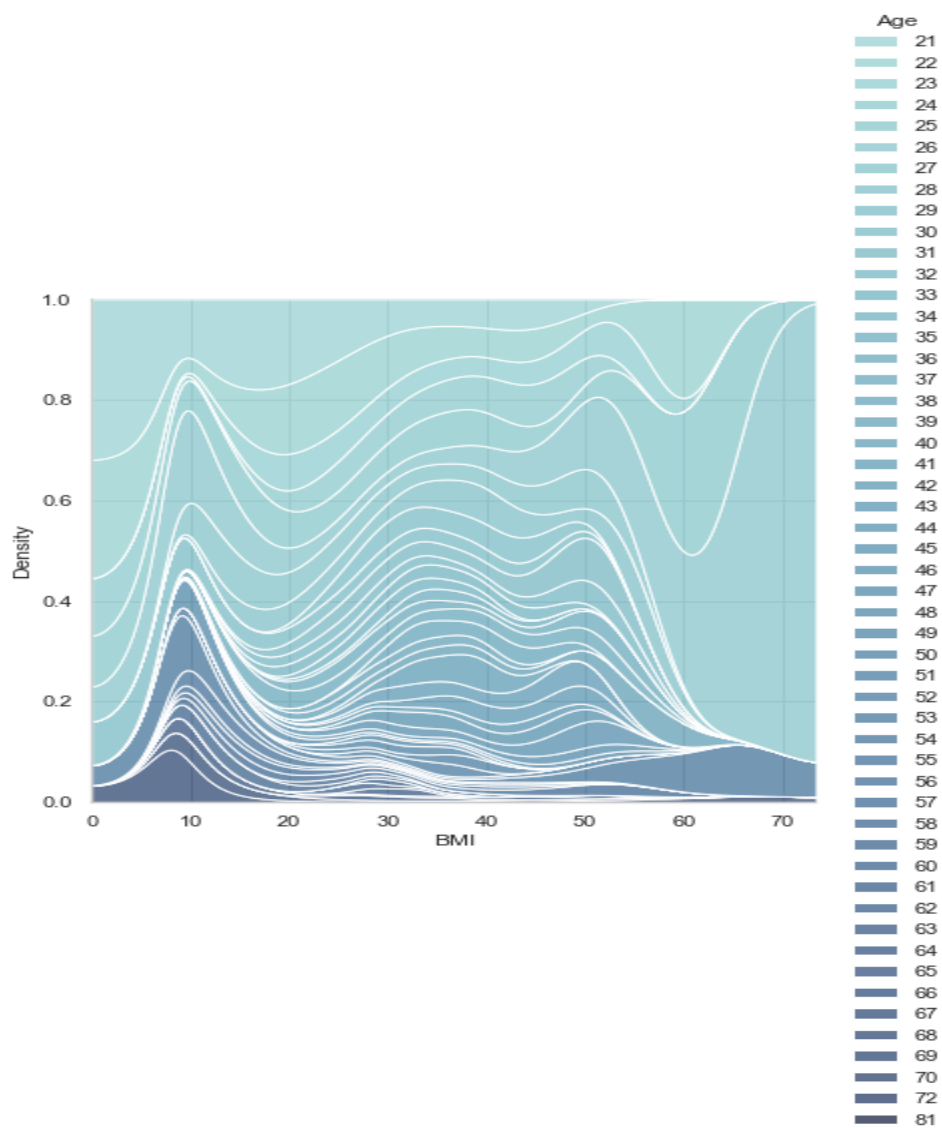
In this graph we are able to perceive the correlation between Diabetes Pedigree Function and Age. The factor of AGE contributes much in this case as the diabetes pedigree function stays increases and then decreases throughout. On the side we have the marginal side plots, wherein we can see the marginal distribution pertaining to diabetes pedigree function and AGE. in this we are able to see some data on the graph from the middle around the age of 20 as that is where the density begins to increase and maxes out within no time and gradually slopes down as the density decreases. The above marginal graph gives a great idea on how the density increases then decreases.

● Scatterplot Matrix

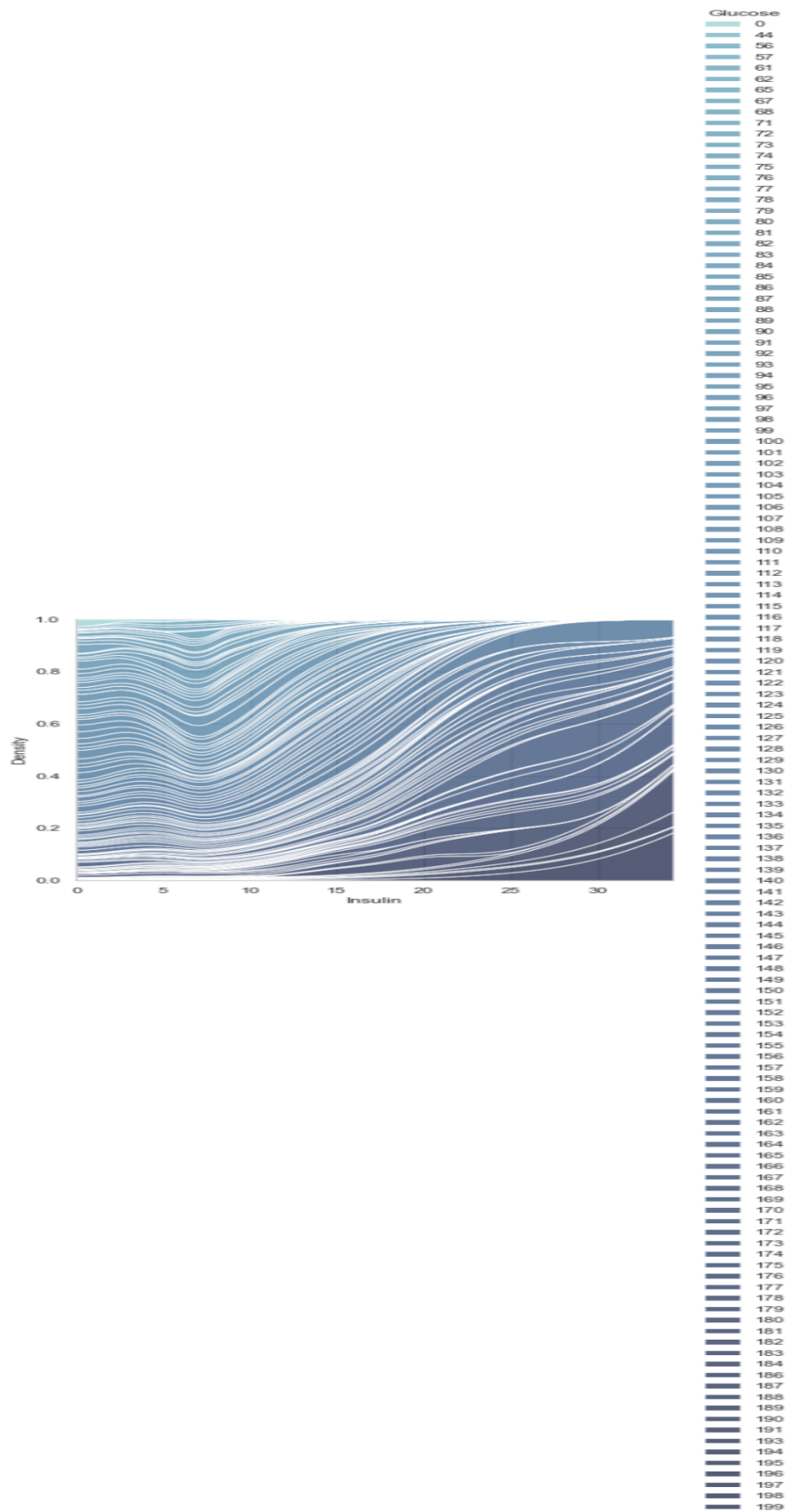


In this we are comparing all the attributes that are provided in the dataset and plotting a scatter plot matrix. Here, we have a scatterplot matrix wherein we use age as the base comparison feature. This plot gives a clear idea on how age contributes to each and every feature or body problems. Age is labeled in the form as hue where in when the age increases, the color darkness increases.

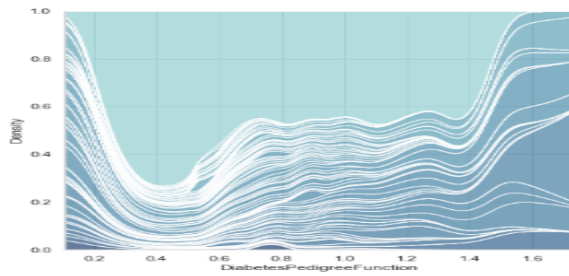
- **Conditional kernel density estimate**



In this graph we are able to visualize the correlation between age and bmi as a factor of having diabetes. The bottom quarter of the relation is dominated by the old people, since their BMI is usually quite low as they tend to lose weight. ie. in the area around 10 BMI. The middle part consists of the age group of around 40-55. The upper portion of the viz. Consists of lower aged people since they tend to be of different weight and height proportions.



In this graph we are able to visualize the correlation between the insulin level and glucose as a factor of having diabetes. Consider the levels from 5-15 the glucose levels are in the middle level but as the number goes up around 20 to 30 the graph gets darker. We can conclude from the graph that with the increasing levels of insulin, the glucose also goes higher.



Insulin

0.0

3.7416573667739413

3.872986346207417

4.0

4.242640687119285

4.69041575982343

4.795831523312719

5.0

5.365164807134504

5.659854249492381

6.0

6.082762530298219

6.164414022928975

6.324555320336759

6.4031242374328485

6.48017468940786

6.557438524302

6.6532495807108

6.708203832480389

6.782329863125268

6.928203230278509

7.0

7.0710678118654755

7.14142842894285

7.211102550927978

7.280108889280516

7.3484692283480340

7.416198467095663

7.483314773547883

7.5498344327075

7.615773105863909

7.681145747868508

7.745966682414834

7.810245675906654

7.837253833153772

8.0

8.062257748218855

8.12403840463596

8.18535277187245

8.248211201230321

8.309800265340756

8.426148773176359

8.48528137423857

8.54400374531753

8.602325267042627

8.660224037844387

8.717797867081348

8.774864387382123

8.831760866327848

8.888194417315589

9.0

9.055385138137417

9.1104335791443

9.16515136981168

9.219544457292887

9.273818495485704

9.327379853088816

9.38083151964686

9.433881132058603

9.486832980505138

9.539382014169456

9.591663046625436

9.64359714832659

9.694794344809963

9.745958971132712

9.796743710662

10.0

10.24695076559598

10.295630140987

10.392304845413264

10.48808481701515

10.58300524258363

10.677078252031311

10.723805294763608

10.770329814288807

10.908712114635714

10.954451150103322

11.048381017187261

11.180339887498949

11.224972180321824

11.269427669084644

11.313708498984761

11.357816591605547

11.40175425099138

11.448125293076057

11.6188003862225

11.832158566198232

11.916375287812984

12.0

12.041594578792296

12.083045973594572

12.165525060596439

12.24744871391589

12.328628005837892

12.449899597988733

12.489955986796797

12.5698508976535

12.609520212918492

12.649110640673518

12.84523257866513

12.884088726725126

12.922847883320086

12.96148139681572

13.038404810405298

13.07696830622021

13.22875655322953

13.2684591614216

13.341664064126334

13.416407864898739

13.480737863232042

13.527749258468883

13.5646598962750536

13.601470508735444

13.711309200802088

13.784048752090222

13.820274861085254

13.856406460551018

13.892433989448804

13.92838827718412

14.0

14.142135623730951

14.2828568570857

14.317821083276303

14.38749456993816

14.491376746189438

14.66287829861518

14.832398974191326

15.0

15.0996688705415

15.16575088103101

15.198684153570664

15.329709716755891

15.394804318340852

15.49193384829668

15.652475842498529

15.7797338380595

15.811388300841896

15.968719422671311

16.06237840420901

16.278820596099706

16.431676725154983

16.49207763315433

16.492422502470642

16.55294535724685

16.583123861777

16.64331689709324

16.67332000533065

16.73320653068151

16.852299546352718

16.881943016134134

17.00872210923198

17.11724278662369

17.32050807568875

17.435595774162696

17.60681686165901

17.8325480127006

17.916472867168917

18.02775637319946

18.05547008526779

18.110770276274835

18.16882212484895

18.303005217723125

18.49324200880693

18.573665961010276

19.235384061671343

19.3849161731037084

19.672315572906

19.78989867322333

20.048937820763422

20.37154878746336

20.57617895340305

21.06380852847824

21.77154105707724

21.883211108075447

21.908802300206645

22.02271554554524

22.24850546128689

22.58317958127243

23.2378000772448

23.302366395462088

23.345235059857504

24.06241883103193

24.49489742783178

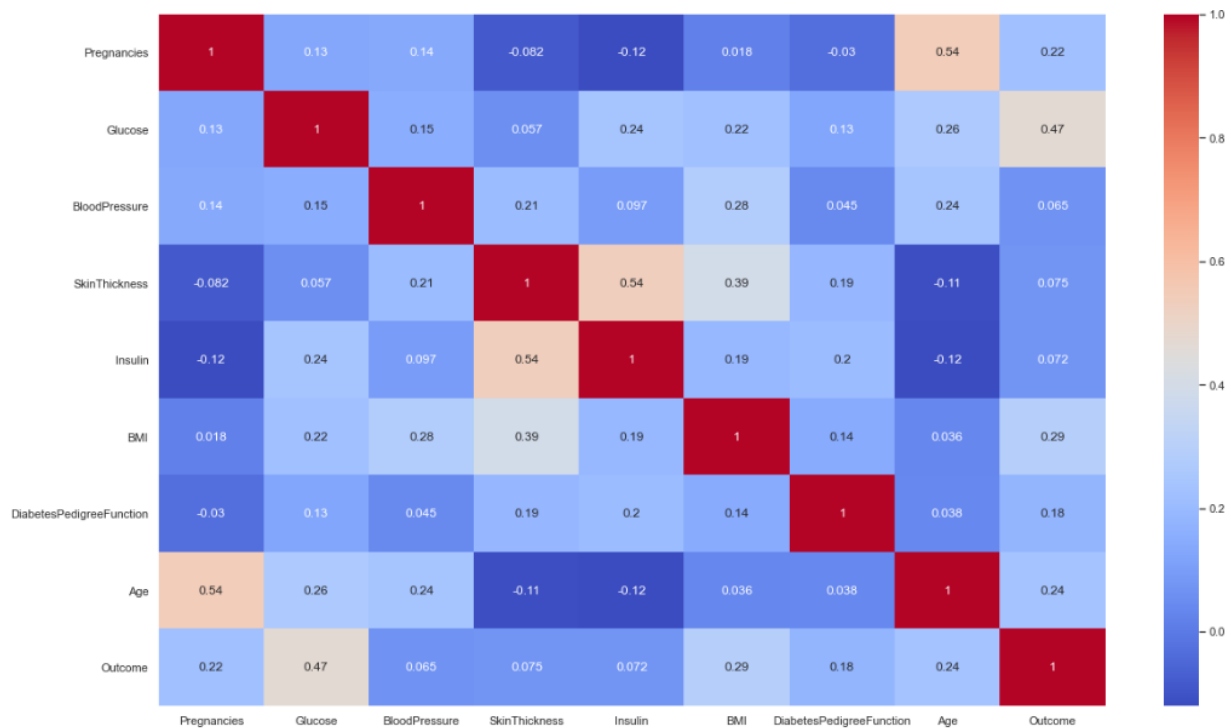
26.076809620810597

27.27836339397171

29.086079144497972

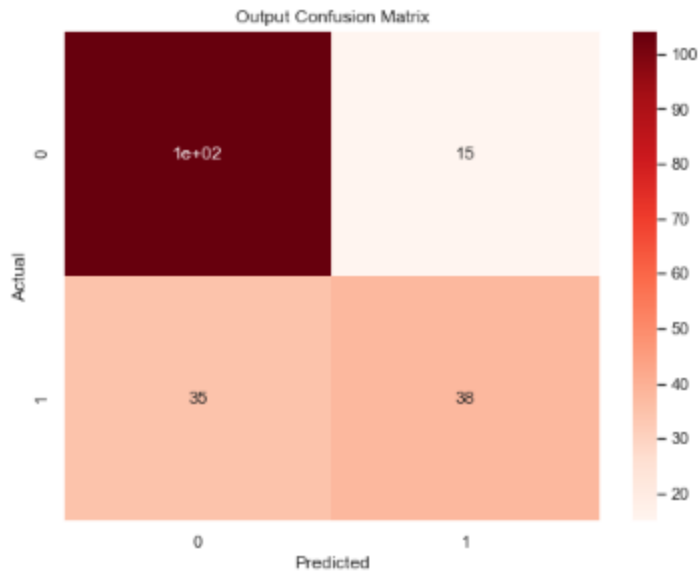
In this graph we are able to visualize the correlation between the insulin level and diabetes pedigree function as a factor of having diabetes. Diabetes pedigree function (**a function which scores likelihood of diabetes based on family history**). We can see from the graph that the low levels around 0 - 0.4 and 0.5 to higher levels can be found easily. All the other middle area looks pretty normal

● Correlation Matrix



This matrix gives us a clear picture of how each and every feature is correlated to the other. The warmer the label is, the more the contributing factor it is in that feature. This is also a diagonal matrix with value 1 where it has its own feature repeated. The other values correspond to the relationship between the feature and the other feature it is contributing to.

● Implementing Confusion Matrix



This matrix is derived after producing the supervised training methods. This confusion matrix gives a clear picture of how each and every factor has contributed. Eg. Accuracy, F1 score.

Conclusion:

With the help of algorithms like Univariate, Bivariate, Scatterplot Matrix, Linear Regression, Conditional Kernel density estimation, correlation matrix and confusion matrix we were able to visualize the relationship between the people can possess that could be a trait of having and possibly predicting a person having or getting Diabetes

For reference, visit the link provided below:

Github Link: <https://github.com/legendslayer01/Diabetes-IV-viz>

Refer the link for a more clear view of viz. Pictures and see the code.