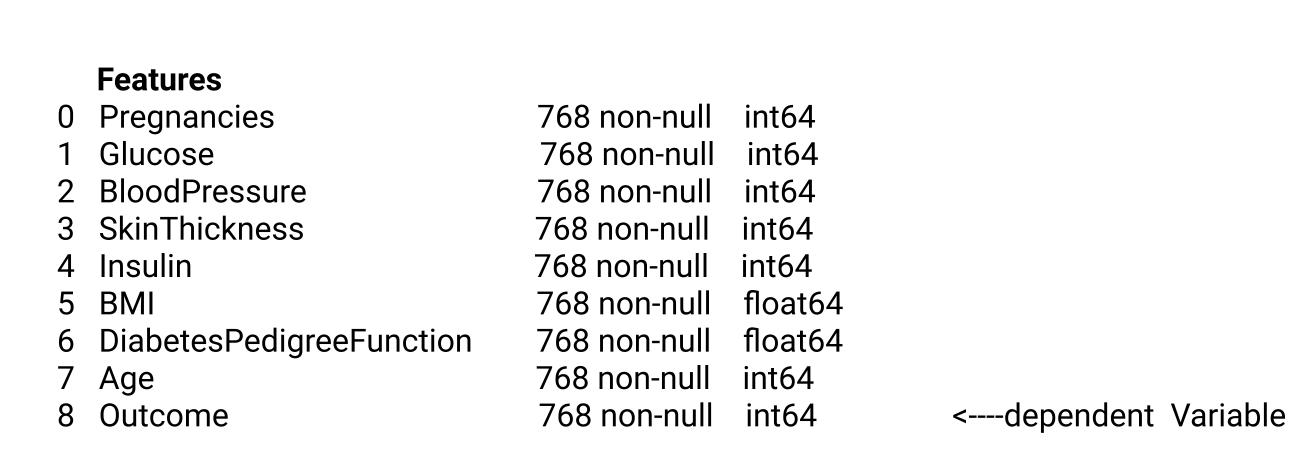
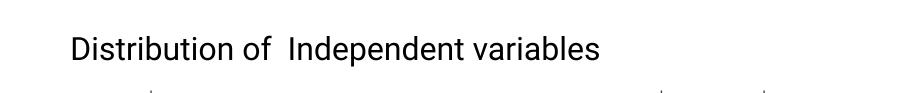
EDA and Comparision of Seed and Synthetic Data (Pima-indians-diabetes)

~Hitesh Kumar (hiteshkumar111@outlook.com)

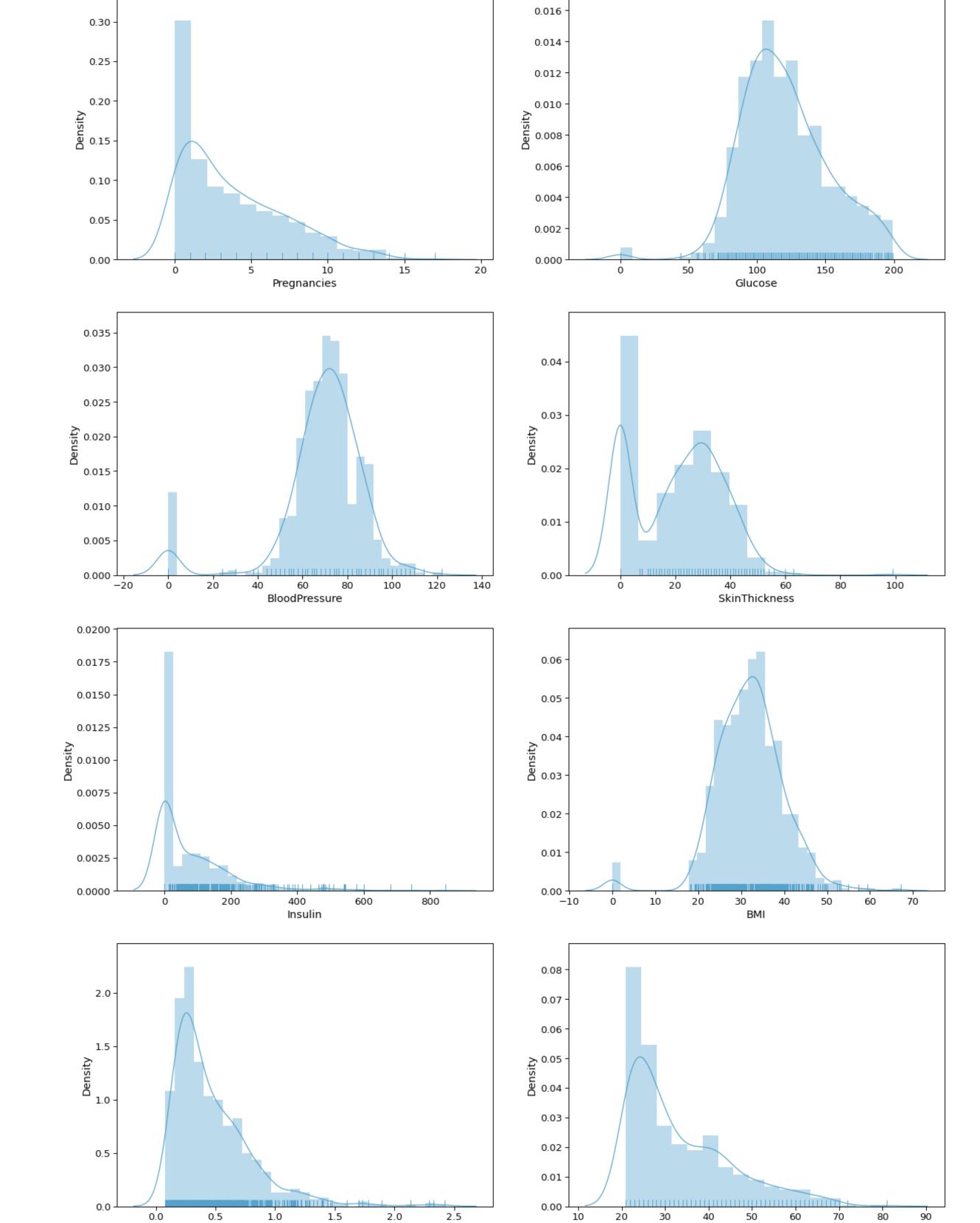
Seed Data

Shape of Seed data: 768 x 9





Distribution of all variables in Seed data



zero values which are in pregnancies, Glucose, Skin Thickness, Blood presssure, Insuline, BMI variable. # In case of Pregnancy variable zeros are explainable. # In other other case they are not, these are data discrepancy, So we will treat them as Missing values.

From here we can see that there are many

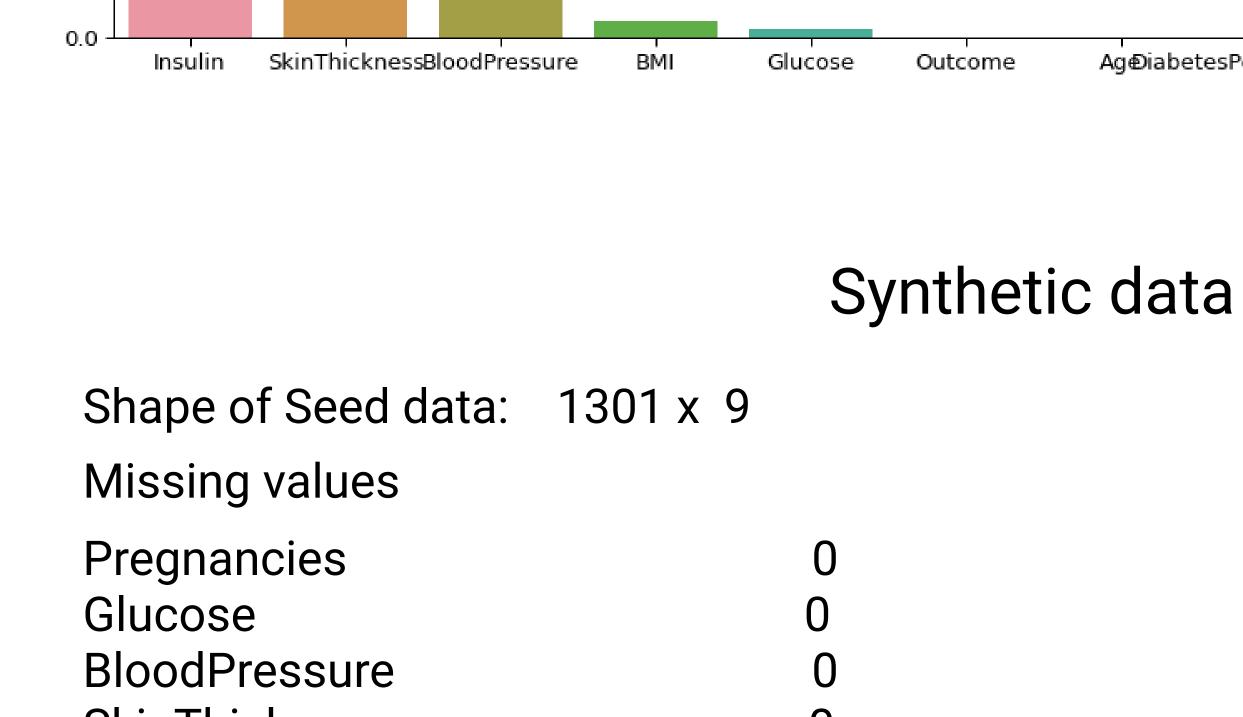
Percent of missing values

0.5

0.1

Missing Values

DiabetesPedigreeFunction



#This Synthetic Missing is generated using Gretel Api

Missing values Percentages

Insulin

BMI

Ag@iabetesPedigreeFun**Ptegm**ancies

Glucose

SkinThickness

BloodPressure

Total Percent

374

48.7%

29.6%

4.6%

1.4%

0.65%

Pregnancies	C
Glucose	0
BloodPressure	C
SkinThickness	0
Insulin	1
BMI	0
DiabetesPedigreeFunction	5
Age	0
Outcome	0
Sec	ed Dat
OC.	

Number of training lines duplicated in synthetic lines

DiabetesPedigreeFunction variable ta vs Synthetic Data

0

187

50

There are total of 6 missing

1 in Insulin variable and 5 in

value in synthetic data.

0.2114

0.2060

--- class_0

--- class_1

100 120

--- class_0

--- class_1

0.025

⊋ 0.020

ළ 0.015

0.010

0.005

0.000

0.06

≥ 0.04

ē 0.03

0.01

0.00

20

60

80

50

Synthetic

40

seed

synthetic

-respectively

Mean Absolute Error between training and	synthetic correlation matric	es 0.0853
Average field Jenson-Shannon distributiona	al distance	0.1827
Note: distance scores range from 0 (no difference) to	1 (maximally different)	
Field Name	Unique Values	Distance Score

Outcome 2 0.1479 BloodPressure 46 0.1418 Age 52 0.1149 Pregnancies 17 0.1148 BMI 247 0.1138 Distribution of Traget variable	<u>Glucose</u>	135	0.1616
Age 52 0.1149 Pregnancies 17 0.1148 BMI 247 0.1138	<u>Outcome</u>		0.1479
Pregnancies 17 0.1148 BMI 247 0.1138	BloodPressure	46	0.1418
BMI 247 0.1138	<u>Age</u>	52	0.1149
	Pregnancies	17	0.1148
	BMI	247	0.1138
	Distribution of Traget variable		

500

300

200

100

--- class_0

--- class_1

80

50

40

Seed

0.6

synthetic

0.02

0.01

0.8

0.6

0.4

0.2

0.0

0.2

0.4

20

100

--- class_0

— class_1

0.00

3.5

3.0

1.0

0.06

0.05

0.04

Density ©00

0.0

200 100

Outcome

<u>Insulin</u>

400

0.175

0.150

0.125

∰ 0.100 ·

_ი.075 -

0.050

0.025

0.000

0.08

0.07

0.06

0.05

0.04

0.02

0.01

0.00

2.00

1.75

1.50

≥ 1.25

1.00

0.75

0.50

0.25

0.00

0.014

0.012

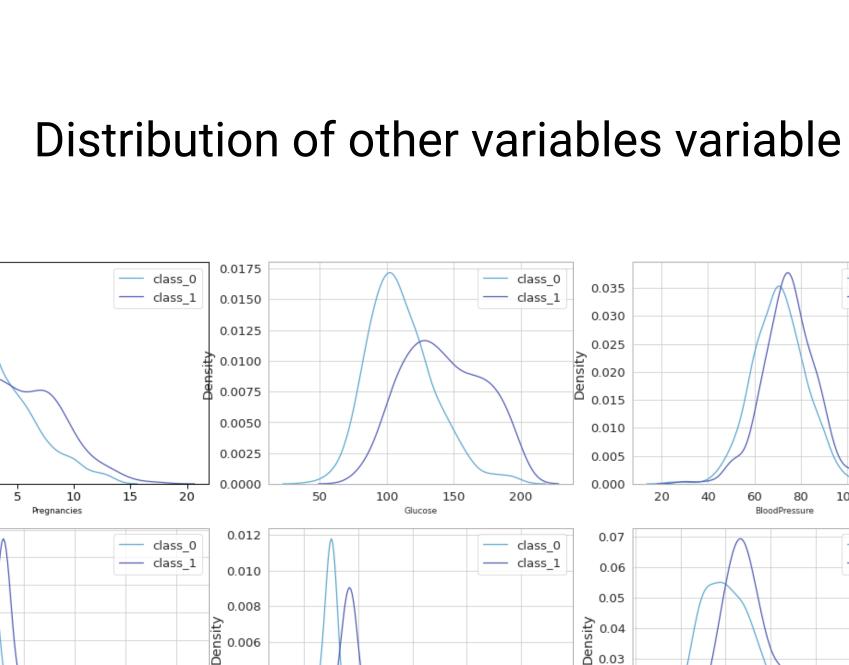
0.010

0.006

0.000

_{0.03} ,

<u>SkinThickness</u>



200

400

600

800

— class_0

— class_1

0.004

0.002

0.04

0.03

0.02

0.01

20

60

0.5 1.0 1.5 2.0 2.5

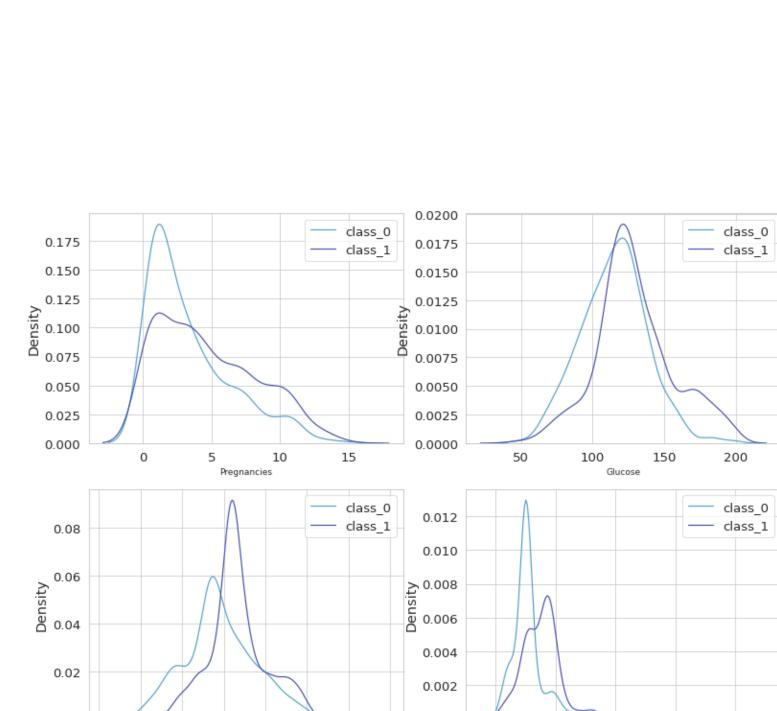
40

80

100

--- class_0

--- class_1



0.000

0.03

0.02

0.01

400

600

800

--- class_0

--- class_1

50

60

--- class_0

--- class_1

40

DiabetesPedigreeFunction

30

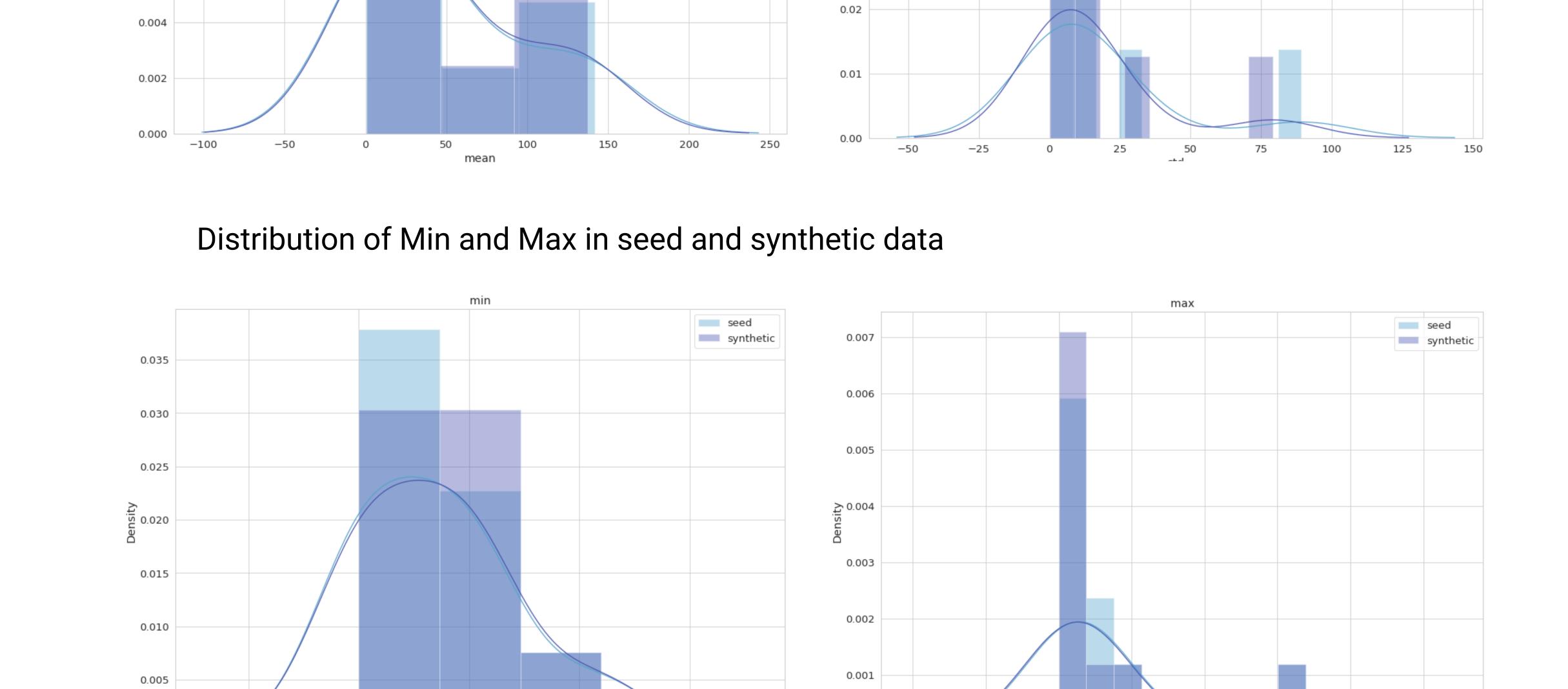
Outcome

Distribution of Mean and Std. deviation in seed and synthetic data

20

Distribution of 25%, 50% and 75% of data in seed and synthetic data

60



0.000

-250

250

max

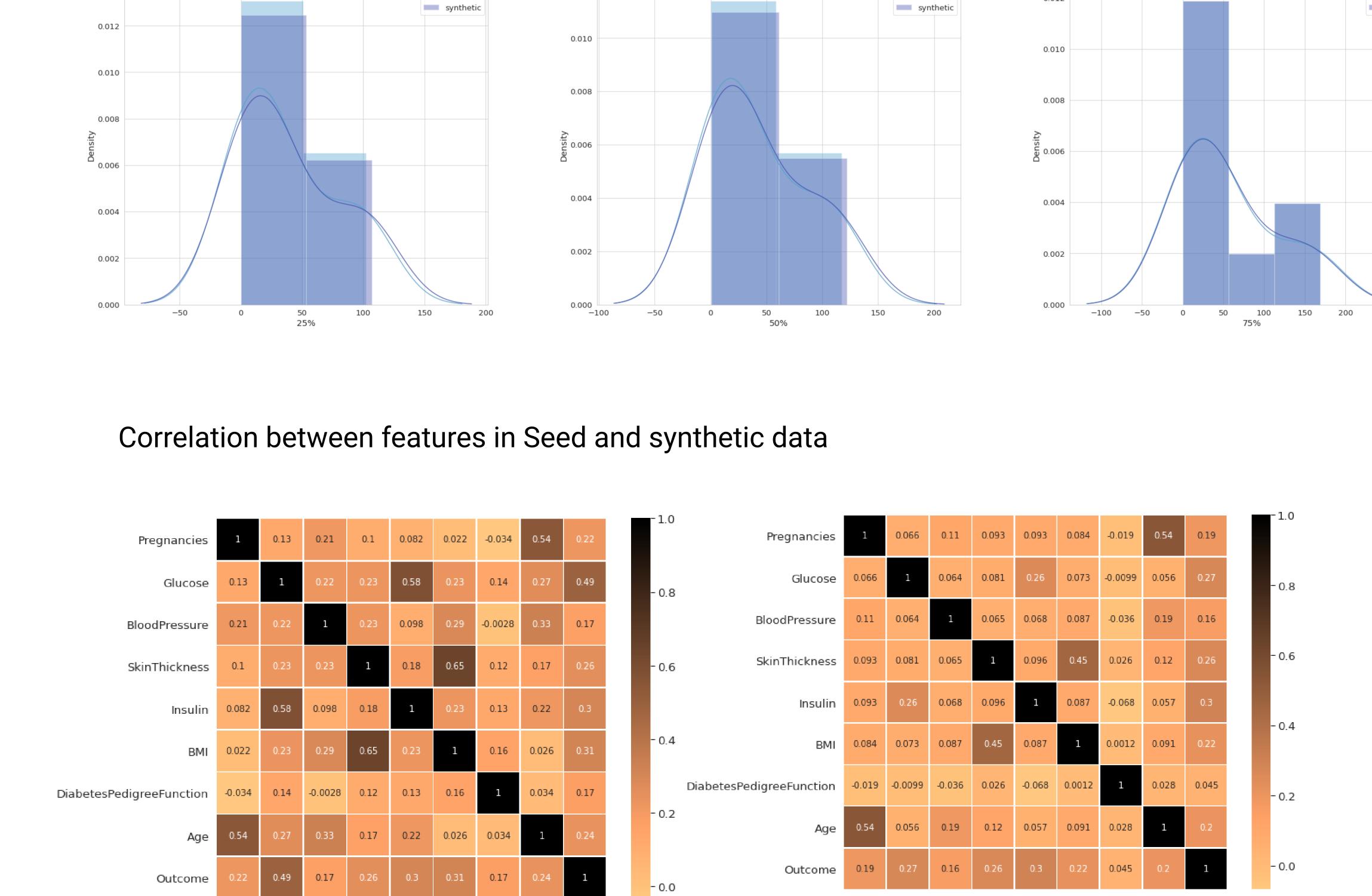
750

Synthetic

1000

1250

250



Seed Distribution of correlation values of all features in seed and synthetic data

in seed data but in the synthetic data is significantly low.

Insulin

each other.



Distribution of Mean, std. deviation, min and max over all features is quite same, as they overlap

#Distribution of 25%, 50% and 75% of data over all features in seed and synthetic data is also quite same.

There are some diffrences pin correlation between features of two seed and synthetic data but for

more part its is quite comparable, for e.g. there is a correlation of +0.58 between insulin and glucose

Jensen-Shannon score(0.1827) is also low which means there's less diffrence between seed and synthetic data.