

Diabetes classification using KNN

```
In [1]: import pandas as pd
import numpy as np
```

```
In [3]: data = pd.read_csv("diabetes.csv")
data.head()
```

```
Out[3]:
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
0	6	148	72	35	0	33.6	0.627	50	1
1	1	85	66	29	0	26.6	0.351	31	0
2	8	183	64	0	0	23.3	0.672	32	1
3	1	89	66	23	94	28.1	0.167	21	0
4	0	137	40	35	168	43.1	2.288	33	1

```
In [5]: data.isnull().any()
```

```
Out[5]: Pregnancies      False
Glucose      False
BloodPressure  False
SkinThickness  False
Insulin      False
BMI          False
DiabetesPedigreeFunction  False
Age          False
Outcome      False
dtype: bool
```

```
In [7]: data.describe().T
```

```
In [7]: data.describe().T
```

```
Out[7]:
```

	count	mean	std	min	25%	50%	75%	max
Pregnancies	768.0	3.845052	3.369578	0.000	1.00000	3.0000	6.00000	17.00
Glucose	768.0	120.894531	31.972618	0.000	99.00000	117.0000	140.25000	199.00
BloodPressure	768.0	69.105469	19.355807	0.000	62.00000	72.0000	80.00000	122.00
SkinThickness	768.0	20.536458	15.952218	0.000	0.00000	23.0000	32.00000	99.00
Insulin	768.0	79.799479	115.244002	0.000	0.00000	30.5000	127.25000	846.00
BMI	768.0	31.992578	7.884160	0.000	27.30000	32.0000	36.60000	67.10
DiabetesPedigreeFunction	768.0	0.471876	0.331329	0.078	0.24375	0.3725	0.62625	2.42
Age	768.0	33.240885	11.760232	21.000	24.00000	29.0000	41.00000	81.00
Outcome	768.0	0.348958	0.476951	0.000	0.00000	0.0000	1.00000	1.00

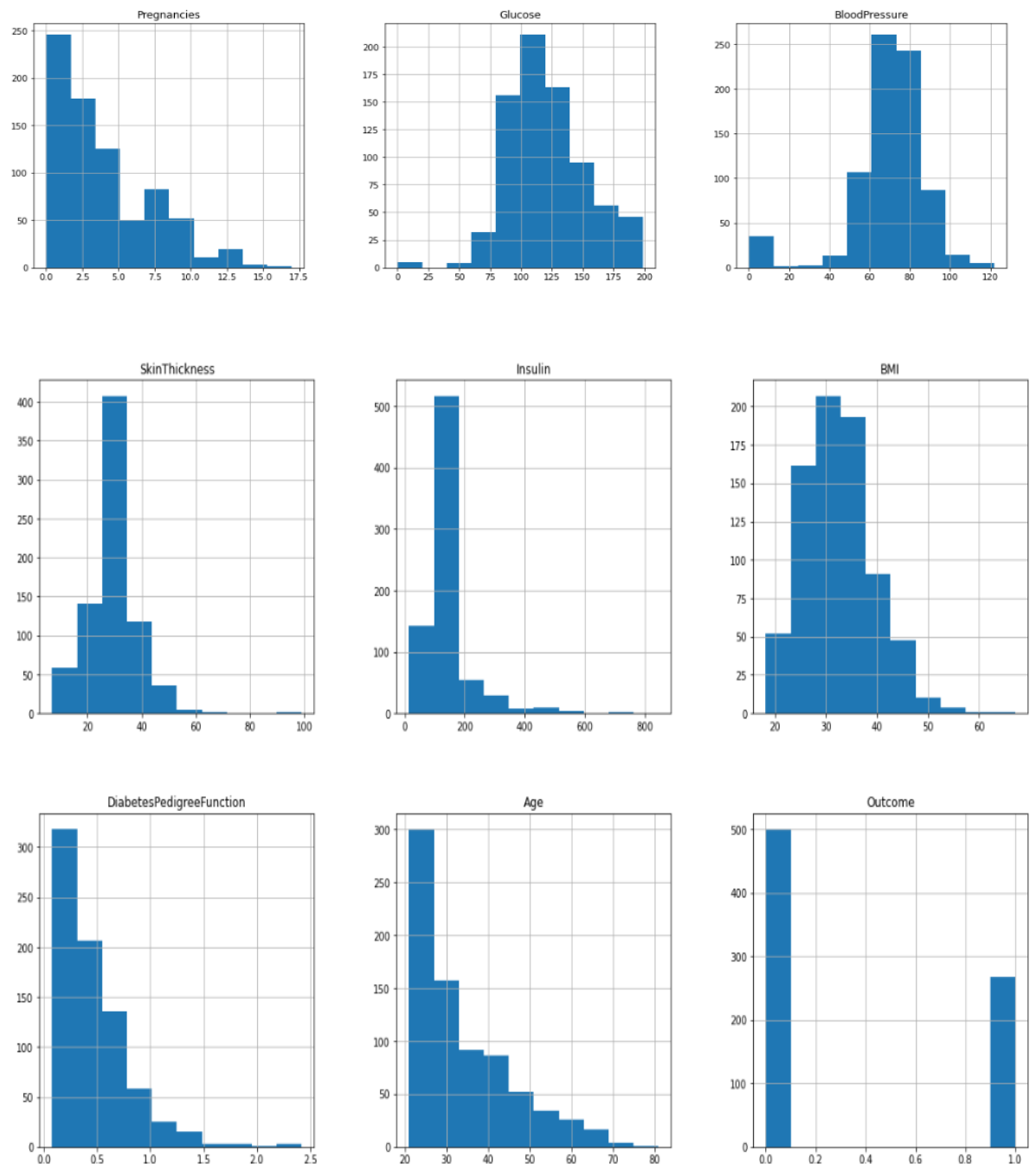
Glucose, BloodPressure, SkinThickness, Insulin, BMI
columns have values 0 which does not make sense, hence are missing values

```
In [13]: data_copy = data.copy(deep = True)
data_copy[['Glucose', 'BloodPressure', 'SkinThickness', 'Insulin', 'BMI']] = data_copy[['Glucose', 'BloodPressure', 'SkinThickness', 'Insulin', 'BMI']].where(data_copy.isnull().sum() > 0)
```

```
Out[13]: Pregnancies      0
          Glucose         5
          BloodPressure   35
          SkinThickness   227
          Insulin         374
          BMI            11
          DiabetesPedigreeFunction  0
          Age            0
          Outcome         0
          dtype: int64
```

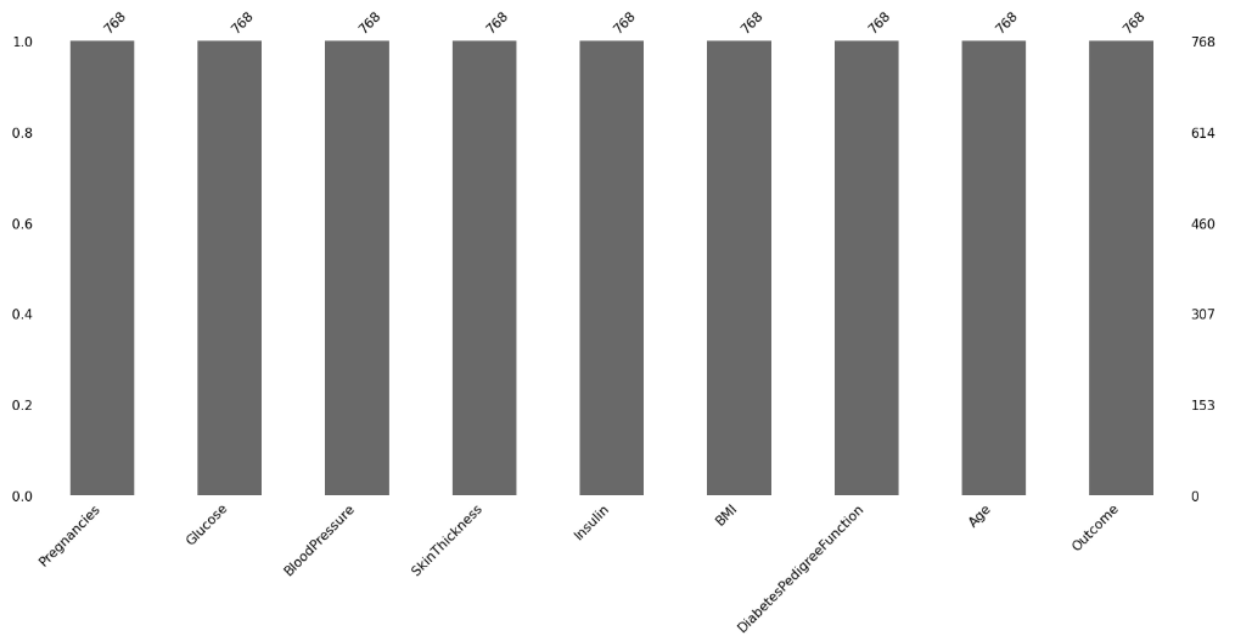
To fill these Nan values the data distribution needs to be understood

```
In [14]: p = data.hist(figsize = (20,20))
```



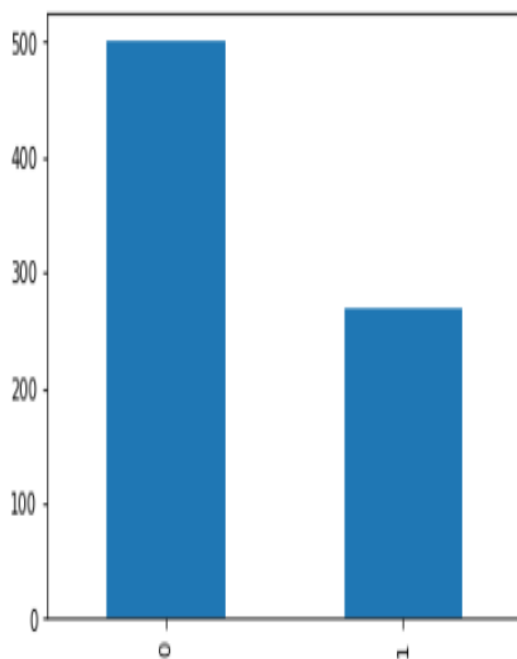
```
In [18]: import missingno as msno
          p = msno.bar(data)
```

```
In [18]: import missingno as msno
p = msno.bar(data)
```



```
In [24]: p=data.Outcome.value_counts().plot(kind="bar")
```

```
In [24]: p=data.Outcome.value_counts().plot(kind="bar")
```



The above graph shows that the data is biased towards datapoints having outcome value as 0 where it means that diabetes was not present actually. The number of non-diabetics is almost twice the number of diabetic patients