

```

import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import matplotlib inline

data_df=pd.read_csv('C:\Users\Hohd\Shadab\Downloads\Projects\CAPSTONE PROJECTS\Project 2\Healthcare - Diabetes\Health care diabetes.csv')
data_df

Out[2]:
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
...
763     10     101           76           48     180  32.9              0.171  63      0
764     2     122           70           27      0  36.8              0.340  27      0
765     5     121           72           23     112  26.2              0.245  30      0
766     1     126           60           0      0  30.1              0.349  47      1
767     1      93           70           31      0  30.4              0.315  23      0
768 rows x 9 columns

In [3]:
data_df.shape

Out[3]:
(768, 9)

In [4]:
data_df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 768 entries, 0 to 767
Data columns (total 9 columns):
 #   Column                Non-Null Count  Dtype
---  --
 0   Pregnancies           768 non-null   int64
 1   Glucose               768 non-null   int64
 2   BloodPressure         768 non-null   int64
 3   SkinThickness         768 non-null   int64
 4   Insulin               768 non-null   int64
 5   BMI                  768 non-null   float64
 6   DiabetesPedigreeFunction  768 non-null   float64
 7   Age                  768 non-null   int64
 8   Outcome               768 non-null   int64
dtypes: float64(2), int64(7)
memory usage: 54.1 KB

In [5]:
data_df.describe()

Pregnancies      0
Glucose          35
BloodPressure    227
SkinThickness    374
Insulin         11
BMI              0
DiabetesPedigreeFunction  0
Age              0
Outcome          0
dtypes: int64(8)

In [6]:
data_df.isnull().sum()*100/len(data_df) # Percentage of null values

Pregnancies      0.000000
Glucose          0.000000
BloodPressure    0.000000
SkinThickness    0.000000
Insulin          0.000000
BMI              0.000000
DiabetesPedigreeFunction  0.000000
Age              0.000000
Outcome          0.000000
dtypes: float64(2)

In [7]:
#Since the data is skewed, we will use median values to treat null values

data_df['Glucose'].fillna(data_df['Glucose'].median(), inplace=True)
data_df['BloodPressure'].fillna(data_df['BloodPressure'].median(), inplace=True)
data_df['SkinThickness'].fillna(data_df['SkinThickness'].median(), inplace=True)
data_df['Insulin'].fillna(data_df['Insulin'].median(), inplace=True)
data_df['BMI'].fillna(data_df['BMI'].median(), inplace=True)

In [8]:
data_df.isnull().sum() #When Null values have been treated

Pregnancies      0
Glucose          0
BloodPressure    0
SkinThickness    0
Insulin          0
BMI              0
DiabetesPedigreeFunction  0
Age              0
Outcome          0
dtypes: int64(8)

In [9]:
data_df.isnull().sum()*100/len(data_df) # Percentage of null values

Pregnancies      0.000000
Glucose          0.000000
BloodPressure    0.000000
SkinThickness    0.000000
Insulin          0.000000
BMI              0.000000
DiabetesPedigreeFunction  0.000000
Age              0.000000
Outcome          0.000000
dtypes: float64(2)

In [10]:
#Since the data is skewed, we will use median values to treat null values

data_df['Glucose'].fillna(data_df['Glucose'].median(), inplace=True)
data_df['BloodPressure'].fillna(data_df['BloodPressure'].median(), inplace=True)
data_df['SkinThickness'].fillna(data_df['SkinThickness'].median(), inplace=True)
data_df['Insulin'].fillna(data_df['Insulin'].median(), inplace=True)
data_df['BMI'].fillna(data_df['BMI'].median(), inplace=True)

In [11]:
data_df.isnull().sum() #When Null values have been treated

Pregnancies      0
Glucose          0
BloodPressure    0
SkinThickness    0
Insulin          0
BMI              0
DiabetesPedigreeFunction  0
Age              0
Outcome          0
dtypes: int64(8)

In [12]:
data_df.isnull().sum()*100/len(data_df) # Percentage of null values

Pregnancies      0.000000
Glucose          0.000000
BloodPressure    0.000000
SkinThickness    0.000000
Insulin          0.000000
BMI              0.000000
DiabetesPedigreeFunction  0.000000
Age              0.000000
Outcome          0.000000
dtypes: float64(2)

In [13]:
#Since the data is skewed, we will use median values to treat null values

data_df['Glucose'].fillna(data_df['Glucose'].median(), inplace=True)
data_df['BloodPressure'].fillna(data_df['BloodPressure'].median(), inplace=True)
data_df['SkinThickness'].fillna(data_df['SkinThickness'].median(), inplace=True)
data_df['Insulin'].fillna(data_df['Insulin'].median(), inplace=True)
data_df['BMI'].fillna(data_df['BMI'].median(), inplace=True)

In [14]:
data_df.isnull().sum() #When Null values have been treated

Pregnancies      0
Glucose          0
BloodPressure    0
SkinThickness    0
Insulin          0
BMI              0
DiabetesPedigreeFunction  0
Age              0
Outcome          0
dtypes: int64(8)

In [15]:
data_df.isnull().sum()*100/len(data_df) # Percentage of null values

Pregnancies      0.000000
Glucose          0.000000
BloodPressure    0.000000
SkinThickness    0.000000
Insulin          0.000000
BMI              0.000000
DiabetesPedigreeFunction  0.000000
Age              0.000000
Outcome          0.000000
dtypes: float64(2)

In [16]:
sns.pairplot(data_df)

In [17]:
data_df.dtypes.value_counts().plot.bar(label='Count', xlabel='Datatype');

In [18]:
plt.figure(figsize=(18,16));
sns.pairplot(data_df)

In [19]:
data_df['Outcome'].value_counts()*100/len(data_df)

Outcome
0    65.104167
1    34.895833
Name: Outcome, dtype: float64

In [20]:
#The above percentage shows that the data is not imbalanced.

In [21]:
data_df.columns

Index(['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin',
      'BMI', 'DiabetesPedigreeFunction', 'Age', 'Outcome'],
      dtype='object')

In [22]:
#Q2. Create scatter charts between the pair of variables to understand the relationships. Describe your findings.

sns.pairplot(data_df)

In [23]:
sns.heatmap(data_df)

In [24]:
sns.heatmap(data_df)

In [25]:
sns.heatmap(data_df)

In [26]:
sns.heatmap(data_df)

In [27]:
sns.heatmap(data_df)

In [28]:
sns.heatmap(data_df)

In [29]:
sns.heatmap(data_df)

In [30]:
sns.heatmap(data_df)

In [31]:
sns.heatmap(data_df)

In [32]:
sns.heatmap(data_df)

In [33]:
sns.heatmap(data_df)

In [34]:
sns.heatmap(data_df)

In [35]:
sns.heatmap(data_df)

In [36]:
sns.heatmap(data_df)

In [37]:
sns.heatmap(data_df)

In [38]:
sns.heatmap(data_df)

In [39]:
sns.heatmap(data_df)

In [40]:
sns.heatmap(data_df)

In [41]:
sns.heatmap(data_df)

In [42]:
sns.heatmap(data_df)

In [43]:
sns.heatmap(data_df)

In [44]:
sns.heatmap(data_df)

In [45]:
sns.heatmap(data_df)

In [46]:
sns.heatmap(data_df)

In [47]:
sns.heatmap(data_df)

In [48]:
sns.heatmap(data_df)

In [49]:
sns.heatmap(data_df)

In [50]:
sns.heatmap(data_df)

In [51]:
sns.heatmap(data_df)

In [52]:
sns.heatmap(data_df)

In [53]:
sns.heatmap(data_df)

In [54]:
sns.heatmap(data_df)

In [55]:
sns.heatmap(data_df)

In [56]:
sns.heatmap(data_df)

In [57]:
sns.heatmap(data_df)

In [58]:
sns.heatmap(data_df)

In [59]:
sns.heatmap(data_df)

In [60]:
sns.heatmap(data_df)

In [61]:
sns.heatmap(data_df)

In [62]:
sns.heatmap(data_df)

In [63]:
sns.heatmap(data_df)

In [64]:
sns.heatmap(data_df)

In [65]:
sns.heatmap(data_df)

In [66]:
sns.heatmap(data_df)

In [67]:
sns.heatmap(data_df)

In [68]:
sns.heatmap(data_df)

In [69]:
sns.heatmap(data_df)

In [70]:
sns.heatmap(data_df)

In [71]:
sns.heatmap(data_df)

In [72]:
sns.heatmap(data_df)

In [73]:
sns.heatmap(data_df)

In [74]:
sns.heatmap(data_df)

In [75]:
sns.heatmap(data_df)

In [76]:
sns.heatmap(data_df)

In [77]:
sns.heatmap(data_df)

In [78]:
sns.heatmap(data_df)

In [79]:
sns.heatmap(data_df)

In [80]:
sns.heatmap(data_df)

In [81]:
sns.heatmap(data_df)

In [82]:
sns.heatmap(data_df)

In [83]:
sns.heatmap(data_df)

In [84]:
sns.heatmap(data_df)

In [85]:
sns.heatmap(data_df)

In [86]:
sns.heatmap(data_df)

In [87]:
sns.heatmap(data_df)

In [88]:
sns.heatmap(data_df)

In [89]:
sns.heatmap(data_df)

In [90]:
sns.heatmap(data_df)

In [91]:
sns.heatmap(data_df)

In [92]:
sns.heatmap(data_df)

In [93]:
sns.heatmap(data_df)

In [94]:
sns.heatmap(data_df)

In [95]:
sns.heatmap(data_df)

In [96]:
sns.heatmap(data_df)

In [97]:
sns.heatmap(data_df)

In [98]:
sns.heatmap(data_df)

In [99]:
sns.heatmap(data_df)

In [100]:
sns.heatmap(data_df)

In [101]:
sns.heatmap(data_df)

In [102]:
sns.heatmap(data_df)

In [103]:
sns.heatmap(data_df)

In [104]:
sns.heatmap(data_df)

In [105]:
sns.heatmap(data_df)

In [106]:
sns.heatmap(data_df)

In [107]:
sns.heatmap(data_df)

In [108]:
sns.heatmap(data_df)

In [109]:
sns.heatmap(data_df)

In [110]:
sns.heatmap(data_df)

In [111]:
sns.heatmap(data_df)

In [112]:
sns.heatmap(data_df)

In [113]:
sns.heatmap(data_df)

In [114]:
sns.heatmap(data_df)

In [115]:
sns.heatmap(data_df)

In [116]:
sns.heatmap(data_df)

In [117]:
sns.heatmap(data_df)

In [118]:
sns.heatmap(data_df)

In [119]:
sns.heatmap(data_df)

In [120]:
sns.heatmap(data_df)

In [121]:
sns.heatmap(data_df)

In [122]:
sns.heatmap(data_df)

In [123]:
sns.heatmap(data_df)

In [124]:
sns.heatmap(data_df)

In [125]:
sns.heatmap(data_df)

In [126]:
sns.heatmap(data_df)

In [127]:
sns.heatmap(data_df)

In [128]:
sns.heatmap(data_df)

In [129]:
sns.heatmap(data_df)

In [130]:
sns.heatmap(data_df)

In [131]:
sns.heatmap(data_df)

In [132]:
sns.heatmap(data_df)

In [133]:
sns.heatmap(data_df)

In [134]:
sns.heatmap(data_df)

In [135]:
sns.heatmap(data_df)

In [136]:
sns.heatmap(data_df)

In [137]:
sns.heatmap(data_df)

In [138]:
sns.heatmap(data_df)

In [139]:
sns.heatmap(data_df)

In [140]:
sns.heatmap(data_df)

In [141]:
sns.heatmap(data_df)

In [142]:
sns.heatmap(data_df)

In [143]:
sns.heatmap(data_df)

In [144]:
sns.heatmap(data_df)

In [145]:
sns.heatmap(data_df)

In [146]:
sns.heatmap(data_df)

In [147]:
sns.heatmap(data_df)

In [148]:
sns.heatmap(data_df)

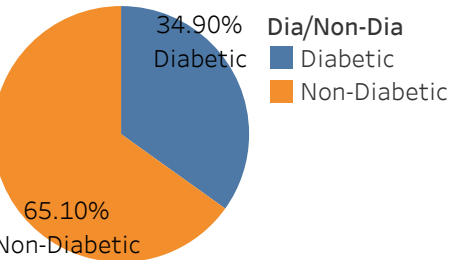
In [149]:
sns.heatmap(data_df)

In [150]:
sns.heatmap(data_df)

In [
```

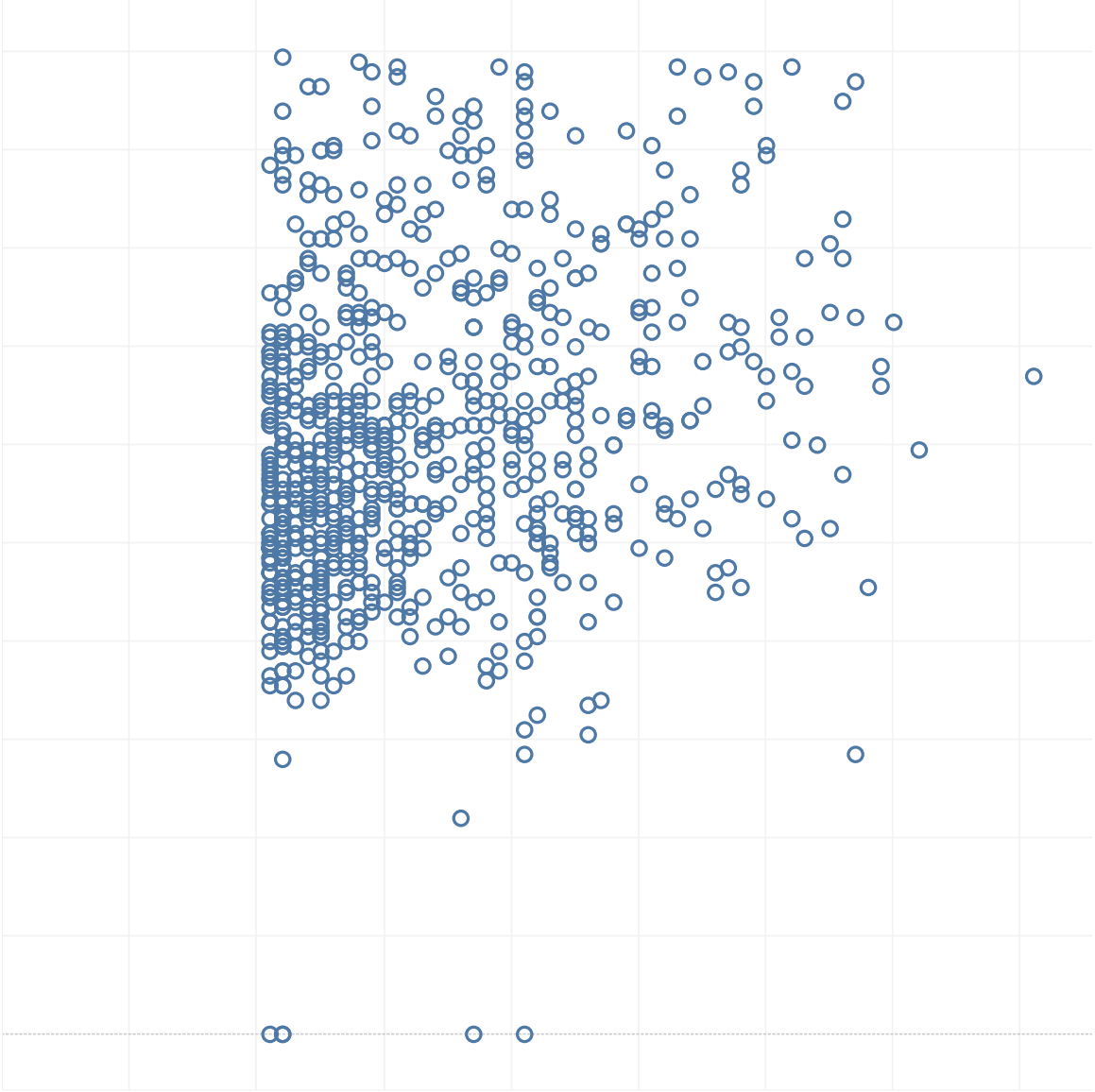


Diabetic vs  
Non-Diabet-  
ic



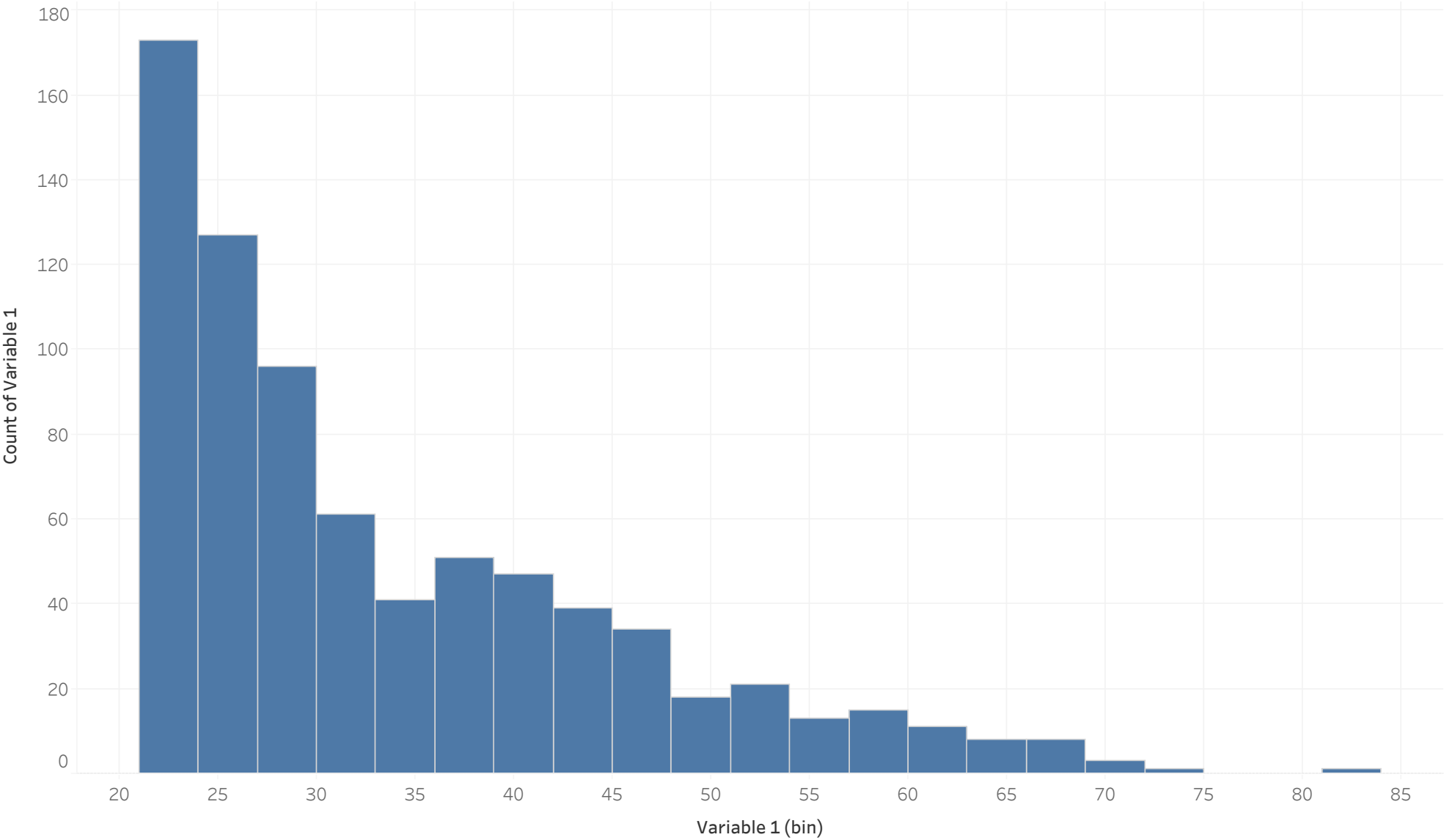
% of Total Count of  
Outcome and  
Dia/Non-Dia. Color  
shows details about  
Dia/Non-Dia. The  
marks are labeled by  
% of Total Count of  
Outcome and  
Dia/Non-Dia.

Age vs Glucose -Scatter Plots



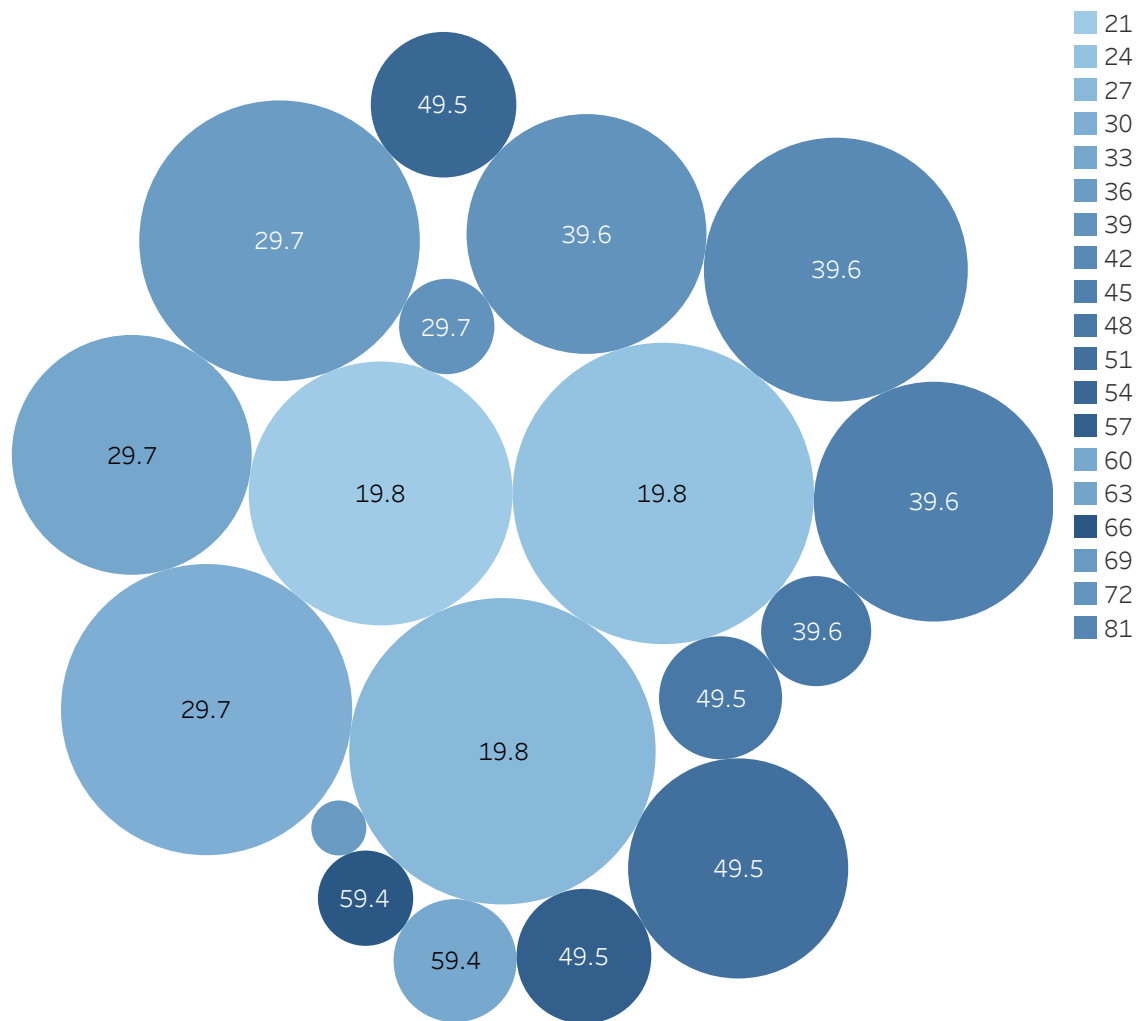
Variable 1 vs. Variable 2.

Age

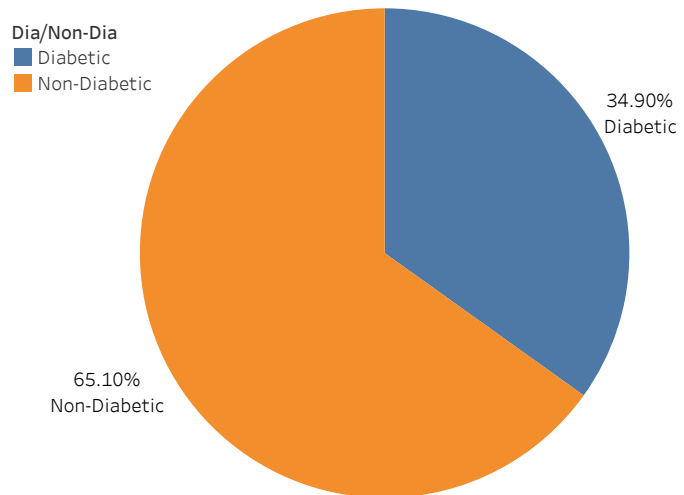


The trend of count of Variable 1 for Variable 1 (bin).

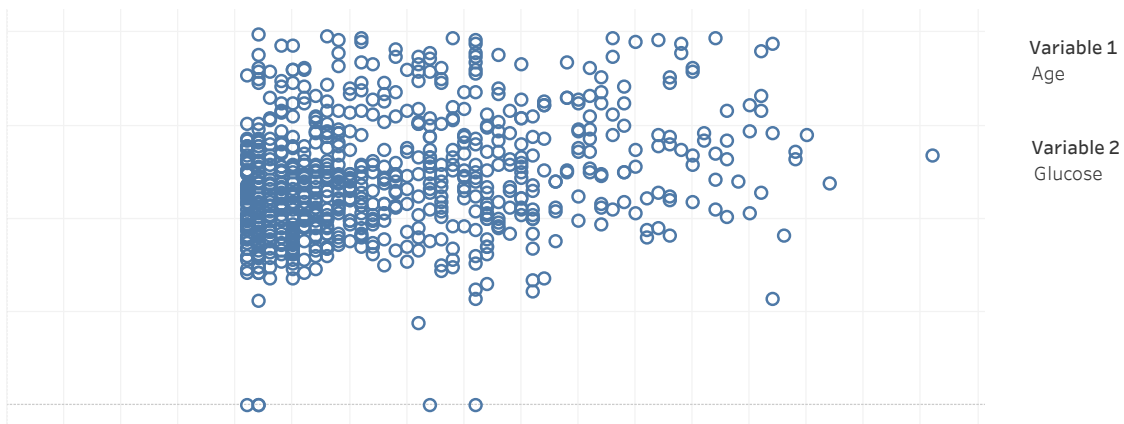
## Age vs Outcome - Bubble Charts



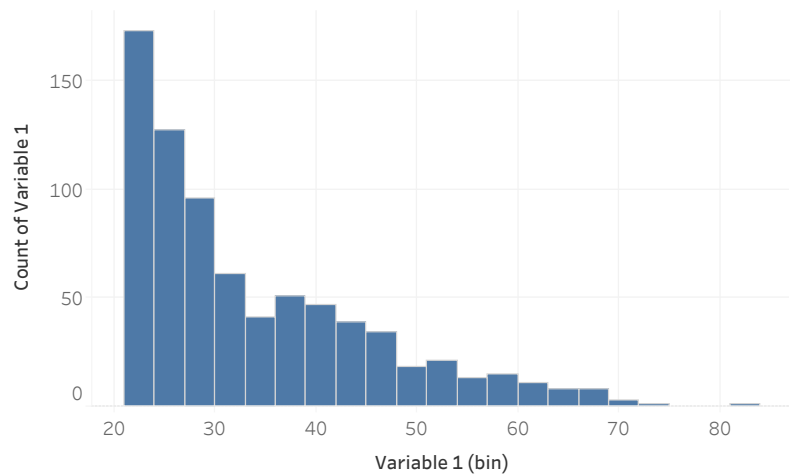
Variable 1 (bin) 2. Color shows details about Variable 1 (bin). Size shows sum of Outcome.  
The marks are labeled by Variable 1 (bin) 2.



Age vs Glucose - Scatter Plots



Age



Age vs Outcome - Bubble Charts

