

```
In [80]: import pandas as pd
import matplotlib
import seaborn as sns
import matplotlib.pyplot as plt

In [ ]: ## |Load data
dataw = pd.read_csv("G:\\CS NOTES\\2.2\\SCIENTIFIC COMPUTING\\Assignment\\pandas_into\\Diabetes Data.csv")
dataw

In [82]: # getting the dataset information
dataw.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 768 entries, 0 to 767
Data columns (total 9 columns):
#   Column                Non-Null Count  Dtype  
---  -
0   Pregnant              768 non-null   int64   
1   Glucose               763 non-null   float64  
2   Diastolic_BP         733 non-null   float64  
3   Skin_Fold             541 non-null   float64  
4   Serum_Insulin         394 non-null   float64  
5   BMI                  757 non-null   float64  
6   Diabetes_Pedigree     768 non-null   float64  
7   Age                  768 non-null   int64   
8   Class                768 non-null   int64   
dtypes: float64(6), int64(3)
memory usage: 54.1 KB

In [83]: # checking the first 5 rows of the dataset
dataw.head()

Out[83]:
```

	Pregnant	Glucose	Diastolic_BP	Skin_Fold	Serum_Insulin	BMI	Diabetes_Pedigree	Age	Class
0	6	148.0	72.0	35.0	NaN	33.6	0.627	50	1
1	1	85.0	66.0	29.0	NaN	26.6	0.351	31	0
2	8	183.0	64.0	NaN	NaN	23.3	0.672	32	1
3	1	89.0	66.0	23.0	94.0	28.1	0.167	21	0
4	0	137.0	40.0	35.0	168.0	43.1	2.288	33	1

```
In [84]: # |checking the Last 5 rows of the dataset
dataw.tail()

Out[84]:
```

	Pregnant	Glucose	Diastolic_BP	Skin_Fold	Serum_Insulin	BMI	Diabetes_Pedigree	Age	Class
763	10	101.0	76.0	48.0	180.0	32.9	0.171	63	0
764	2	122.0	70.0	27.0	NaN	36.8	0.340	27	0
765	5	121.0	72.0	23.0	112.0	26.2	0.245	30	0
766	1	126.0	60.0	NaN	NaN	30.1	0.349	47	1
767	1	93.0	70.0	31.0	NaN	30.4	0.315	23	0

```
In [85]: # checking the size of the data
dataw.shape

Out[85]: (768, 9)

In [86]: # checking the data types for each column of ther dataset
dataw.dtypes

Out[86]: Pregnant                int64
Glucose                float64
Diastolic_BP          float64
Skin_Fold              float64
Serum_Insulin          float64
BMI                   float64
Diabetes_Pedigree     float64
Age                   int64
Class                 int64
dtype: object

In [87]: # getting summary Statistics for each columns
dataw.describe(include='all')
```

Out[87]:

	Pregnant	Glucose	Diastolic_BP	Skin_Fold	Serum_Insulin	BMI	Diabetes_Pedigree	Age	Class
count	768.000000	763.000000	733.000000	541.000000	394.000000	757.000000	768.000000	768.000000	768.000000
mean	3.845052	121.686763	72.405184	29.153420	155.548223	32.457464	0.471876	33.240885	0.348958
std	3.369578	30.535641	12.382158	10.476982	118.775855	6.924988	0.331329	11.760232	0.476951
min	0.000000	44.000000	24.000000	7.000000	14.000000	18.200000	0.078000	21.000000	0.000000
25%	1.000000	99.000000	64.000000	22.000000	76.250000	27.500000	0.243750	24.000000	0.000000
50%	3.000000	117.000000	72.000000	29.000000	125.000000	32.300000	0.372500	29.000000	0.000000
75%	6.000000	141.000000	80.000000	36.000000	190.000000	36.600000	0.626250	41.000000	1.000000
max	17.000000	199.000000	122.000000	99.000000	846.000000	67.100000	2.420000	81.000000	1.000000

In [88]:

```
# checking count of missing values per column
dataw.isnull().sum()
```

Out[88]:

Pregnant0
Glucose5
Diastolic_BP35
Skin_Fold227
Serum_Insulin374
BMI11
Diabetes_Pedigree0
Age0
Class0
dtype: int64

In [89]:

```
# filling the missing values in the dataset with the mean
dataw.fillna(dataw.mean(numeric_only=True), inplace=True)
dataw
```

Out[89]:

	Pregnant	Glucose	Diastolic_BP	Skin_Fold	Serum_Insulin	BMI	Diabetes_Pedigree	Age	Class
0	6	148.0	72.0	35.00000	155.548223	33.6	0.627	50	1
1	1	85.0	66.0	29.00000	155.548223	26.6	0.351	31	0
2	8	183.0	64.0	29.15342	155.548223	23.3	0.672	32	1
3	1	89.0	66.0	23.00000	94.000000	28.1	0.167	21	0
4	0	137.0	40.0	35.00000	168.000000	43.1	2.288	33	1
...
763	10	101.0	76.0	48.00000	180.000000	32.9	0.171	63	0
764	2	122.0	70.0	27.00000	155.548223	36.8	0.340	27	0
765	5	121.0	72.0	23.00000	112.000000	26.2	0.245	30	0
766	1	126.0	60.0	29.15342	155.548223	30.1	0.349	47	1
767	1	93.0	70.0	31.00000	155.548223	30.4	0.315	23	0

768 rows × 9 columns

In [90]:

```
# calculation of key statistics
dataw.mean()
```

Out[90]:

Pregnant3.845052
Glucose121.686763
Diastolic_BP72.405184
Skin_Fold29.153420
Serum_Insulin155.548223
BMI32.457464
Diabetes_Pedigree0.471876
Age33.240885
Class0.348958
dtype: float64

In [91]:

```
dataw.median()
```

Out[91]:

Pregnant3.000000
Glucose117.000000
Diastolic_BP72.202592
Skin_Fold29.153420
Serum_Insulin155.548223
BMI32.400000
Diabetes_Pedigree0.372500
Age29.000000
Class0.000000
dtype: float64

```
In [92]: dataw.mode()
```

Out[92]:

	Pregnant	Glucose	Diastolic_BP	Skin_Fold	Serum_Insulin	BMI	Diabetes_Pedigree	Age	Class
0	1.0	99.0	70.0	29.15342	155.548223	32.0	0.254	22.0	0.0
1	NaN	100.0	NaN	NaN	NaN	NaN	0.258	NaN	NaN

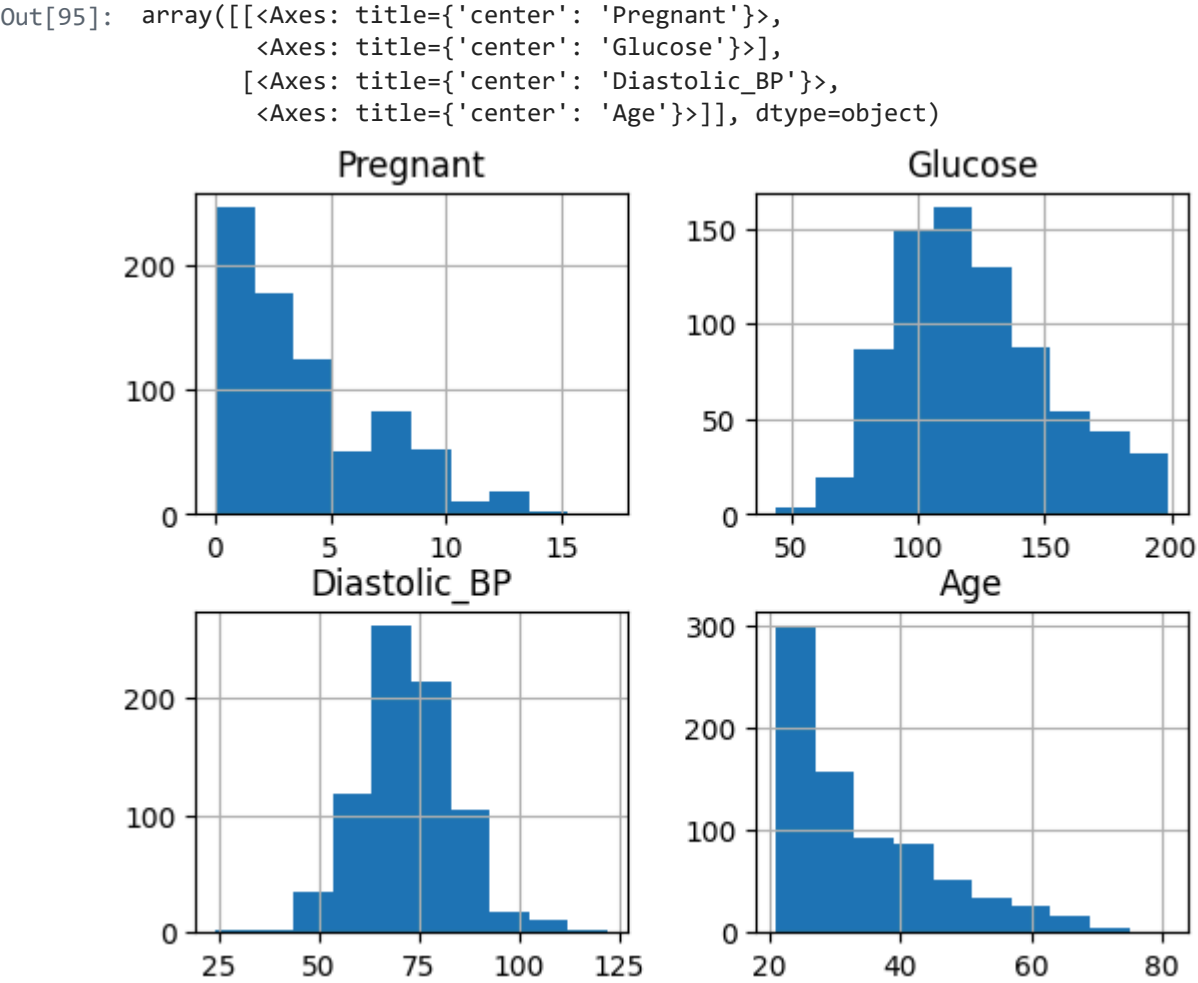
```
In [93]: dataw.std()
```

Out[93]:

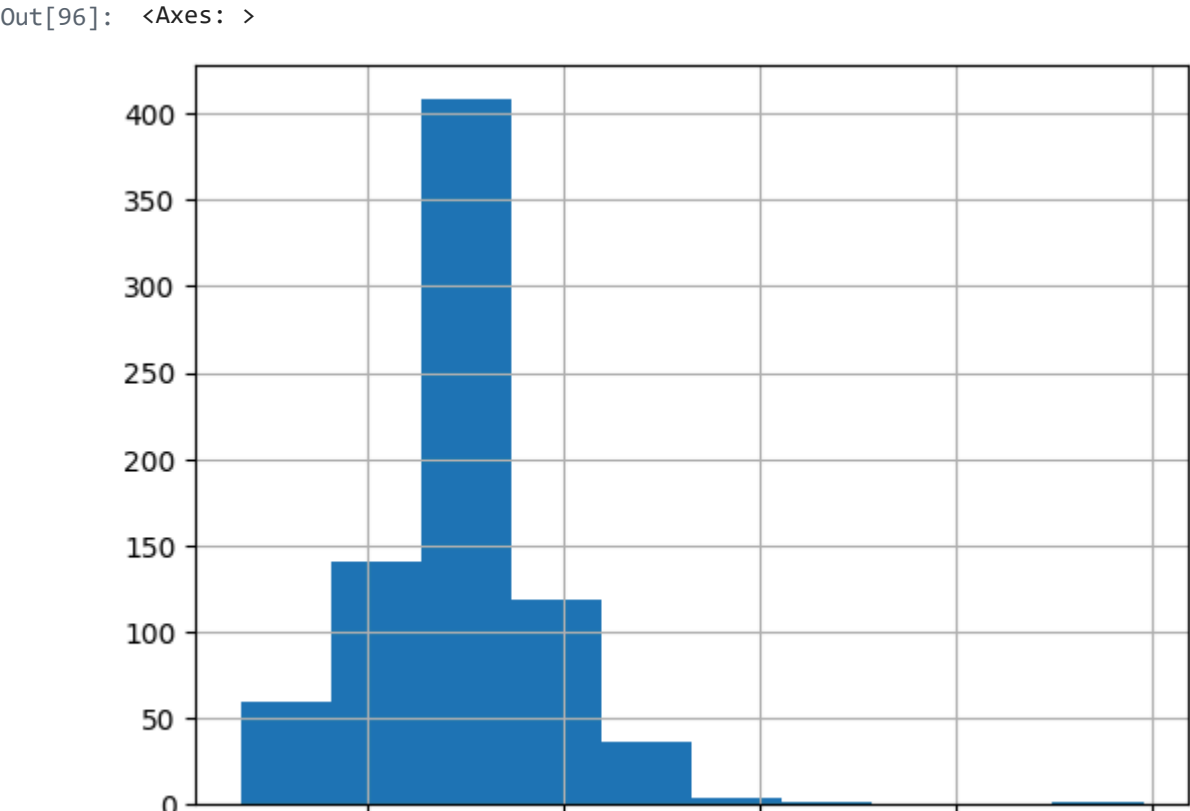
```
Pregnant      3.369578
Glucose       30.435949
Diastolic_BP  12.096346
Skin_Fold      8.790942
Serum_Insulin  85.021108
BMI            6.875151
Diabetes_Pedigree  0.331329
Age           11.760232
Class         0.476951
dtype: float64
```

```
In [94]: # Dataset's Histogram representation for every column separately
# dataw.hist()
```

```
In [95]: # Selecting multiple columns of the dataset and visualizing their data using a histogram
dataw[['Pregnant', 'Glucose', 'Diastolic_BP',"Age"]].hist()
```

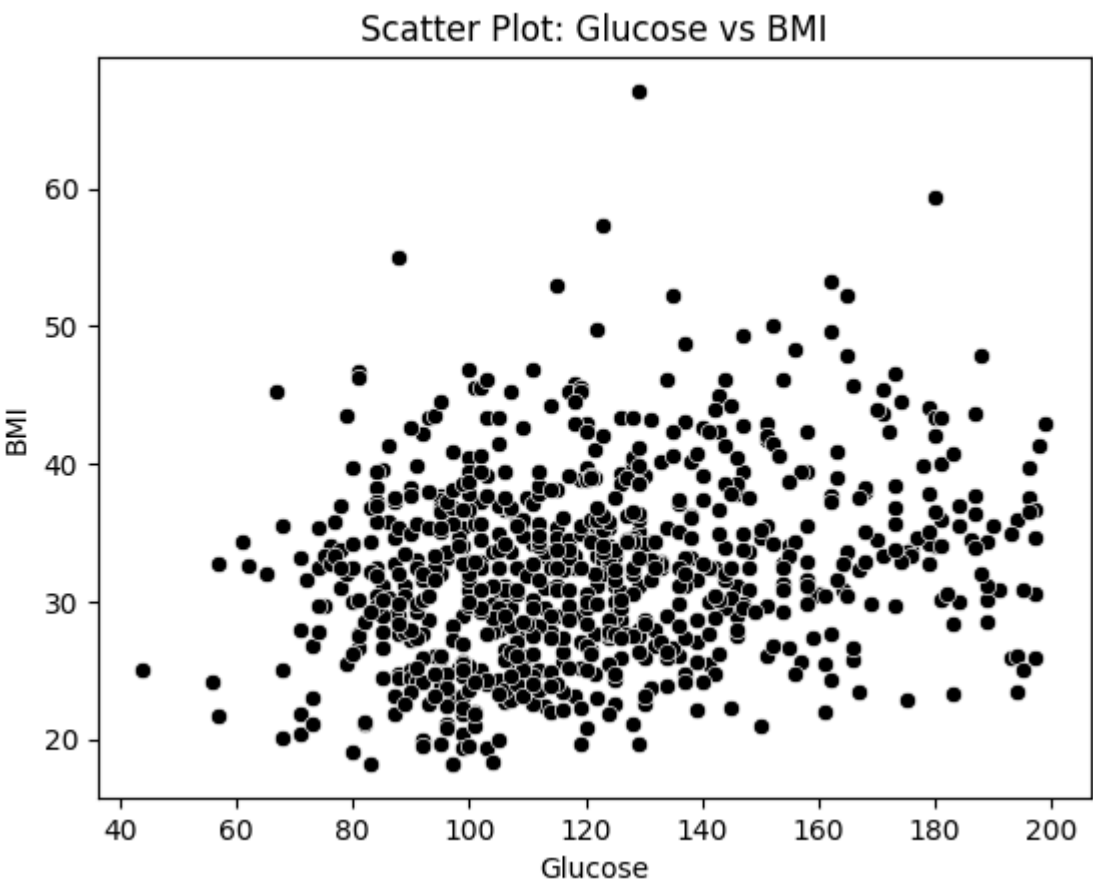


```
In [96]: # Selecting a single column (Skin_Fold column) and visualise it's data
dataw['Skin_Fold'].hist()
```



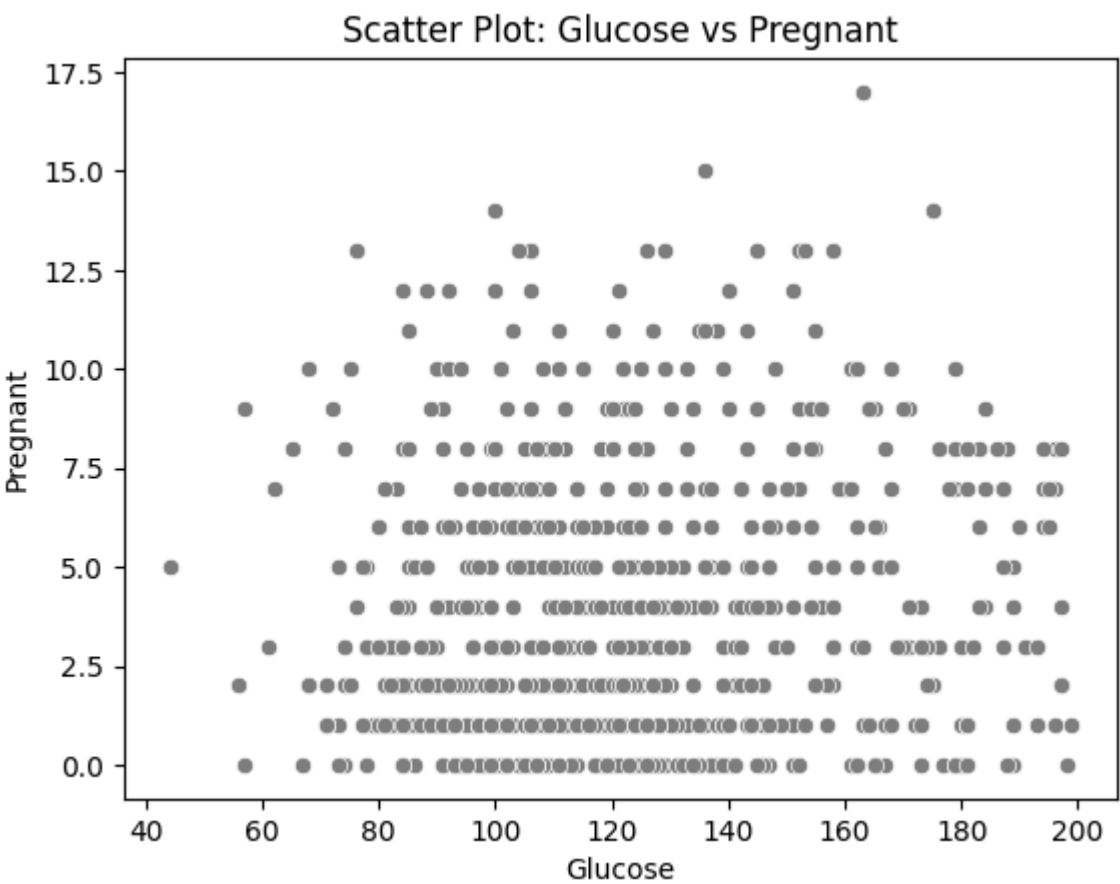
```
In [97]: # scatterplot1
sns.scatterplot(data=dataaw, x="Glucose", y="BMI", color="black")

plt.title("Scatter Plot: Glucose vs BMI")
plt.show()
```



```
In [98]: # scatterplot2
sns.scatterplot(data=dataaw, x="Glucose", y="Pregnant", color="grey")

plt.title("Scatter Plot: Glucose vs Pregnant")
plt.show()
```



```
In [99]: # Compute the correlation matrix
corr_matrix = dataaw.corr()

corr_matrix
```

Out[99]:

	Pregnant	Glucose	Diastolic_BP	Skin_Fold	Serum_Insulin	BMI	Diabetes_Pedigree	Age	Class
Pregnant	1.000000	0.127911	0.208522	0.082989	0.056027	0.021565	-0.033523	0.544341	0.221898
Glucose	0.127911	1.000000	0.218367	0.192991	0.420157	0.230941	0.137060	0.266534	0.492928
Diastolic_BP	0.208522	0.218367	1.000000	0.192816	0.072517	0.281268	-0.002763	0.324595	0.166074
Skin_Fold	0.082989	0.192991	0.192816	1.000000	0.158139	0.542398	0.100966	0.127872	0.215299
Serum_Insulin	0.056027	0.420157	0.072517	0.158139	1.000000	0.166586	0.098634	0.136734	0.214411
BMI	0.021565	0.230941	0.281268	0.542398	0.166586	1.000000	0.153400	0.025519	0.311924

Diabetes_Pedigree	-0.033523	0.137060	-0.002763	0.100966	0.098634	0.153400	1.000000	0.033561	0.173844
Age	0.544341	0.266534	0.324595	0.127872	0.136734	0.025519	0.033561	1.000000	0.238356
Class	0.221898	0.492928	0.166074	0.215299	0.214411	0.311924	0.173844	0.238356	1.000000

In [100...

```
# correlation heatmap
plt.figure(figsize=(10, 6))
sns.heatmap(corr_matrix, annot=True, cmap="coolwarm", fmt=".2f", linewidths=0.5)

plt.title("Correlation Heatmap")
plt.show()
```

