

NGIYEG /
Scientific-Computing

🔍

👤

▼

📄



<> Code

🕒 Issues

🔗 Pull requests

🎬 Actions

📁 Projects

📖 Wiki

🛡 Security

📈 In

📄

🔗

 master ▼

⋮

Scientific-Computing / pandas_intro / pandas_assignment.ipynb 

 NGIYEG Pandas intro

f9c84c4 · 5 minutes ago 

1510 lines (1510 loc) · 276 KB

Preview

Code

Blame

Raw





 ▼

```
In [25]: import pandas as pd
import matplotlib
import seaborn as sns
import matplotlib.pyplot as plt
```

```
In [26]: # Load data
dataaw = pd.read_csv("G:\\CS NOTES\\2.2\\SCIENTIFIC COMPUTING\\Assignment\\Diabet
dataaw
```

```
Out[26]:
```

	Pregnant	Glucose	Diastolic_BP	Skin_Fold	Serum_Insulin	BMI	Diabetes_Pedigree
0	6	148.0	72.0	35.0	NaN	33.6	0.627
1	1	85.0	66.0	29.0	NaN	26.6	0.351
2	8	183.0	64.0	NaN	NaN	23.3	0.672
3	1	89.0	66.0	23.0	94.0	28.1	0.167
4	0	137.0	40.0	35.0	168.0	43.1	2.288
...
763	10	101.0	76.0	48.0	180.0	32.9	0.171
764	2	122.0	70.0	27.0	NaN	36.8	0.340
765	5	121.0	72.0	23.0	112.0	26.2	0.245
766	1	126.0	60.0	NaN	NaN	30.1	0.349
767	1	93.0	70.0	31.0	NaN	30.4	0.315

768 rows × 9 columns

```
In [27]: # getting the dataset information
dataaw.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 [28]: *# checking the first 5 rows of the dataset*
`dataw.head()`

Out[28]:

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

In [29]: *# /checking the last 5 rows of the dataset*
`dataw.tail()`

Out[29]:

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

In [30]: *# checking the size of the data*
`dataw.shape`

Out[30]: (768, 9)

In [31]: *# checking the data types for each column of ther dataset*
`dataw.dtypes`

Out[31]:

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

Out[32]:

	Pregnant	Glucose	Diastolic_BP	Skin_Fold	Serum_Insulin	BMI	Diab
count	768.000000	763.000000	733.000000	541.000000	394.000000	757.000000	
mean	3.845052	121.686763	72.405184	29.153420	155.548223	32.457464	
std	3.369578	30.535641	12.382158	10.476982	118.775855	6.924988	
min	0.000000	44.000000	24.000000	7.000000	14.000000	18.200000	
25%	1.000000	99.000000	64.000000	22.000000	76.250000	27.500000	
50%	3.000000	117.000000	72.000000	29.000000	125.000000	32.300000	
75%	6.000000	141.000000	80.000000	36.000000	190.000000	36.600000	
max	17.000000	199.000000	122.000000	99.000000	846.000000	67.100000	

In [33]:

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

Out[33]:

Pregnant	0
Glucose	5
Diastolic_BP	35
Skin_Fold	227
Serum_Insulin	374
BMI	11
Diabetes_Pedigree	0
Age	0
Class	0
dtype:	int64

In [34]:

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

Out[34]:

	Pregnant	Glucose	Diastolic_BP	Skin_Fold	Serum_Insulin	BMI	Diabetes_Pedigree
0	6	148.0	72.0	35.00000	155.548223	33.6	0.627
1	1	85.0	66.0	29.00000	155.548223	26.6	0.351
2	8	183.0	64.0	29.15342	155.548223	23.3	0.672
3	1	89.0	66.0	23.00000	94.000000	28.1	0.167
4	0	137.0	40.0	35.00000	168.000000	43.1	2.288
...
763	10	101.0	76.0	48.00000	180.000000	32.9	0.171
764	2	122.0	70.0	27.00000	155.548223	36.8	0.340
765	5	121.0	72.0	23.00000	112.000000	26.2	0.245
766	1	126.0	60.0	29.15342	155.548223	30.1	0.349

767 1 93.0 70.0 31.00000 155.548223 30.4 0.315

768 rows × 9 columns



In [35]: *# calculation of key statistics*
dataw.mean()

Out[35]: Pregnant 3.845052
Glucose 121.686763
Diastolic_BP 72.405184
Skin_Fold 29.153420
Serum_Insulin 155.548223
BMI 32.457464
Diabetes_Pedigree 0.471876
Age 33.240885
Class 0.348958
dtype: float64

In [36]: dataw.median()

Out[36]: Pregnant 3.000000
Glucose 117.000000
Diastolic_BP 72.202592
Skin_Fold 29.153420
Serum_Insulin 155.548223
BMI 32.400000
Diabetes_Pedigree 0.372500
Age 29.000000
Class 0.000000
dtype: float64

In [37]: dataw.mode()

Out[37]:

	Pregnant	Glucose	Diastolic_BP	Skin_Fold	Serum_Insulin	BMI	Diabetes_Pedigree	Age
0	1.0	99.0	70.0	29.15342	155.548223	32.0	0.254	2
1	NaN	100.0	NaN	NaN	NaN	NaN	0.258	N

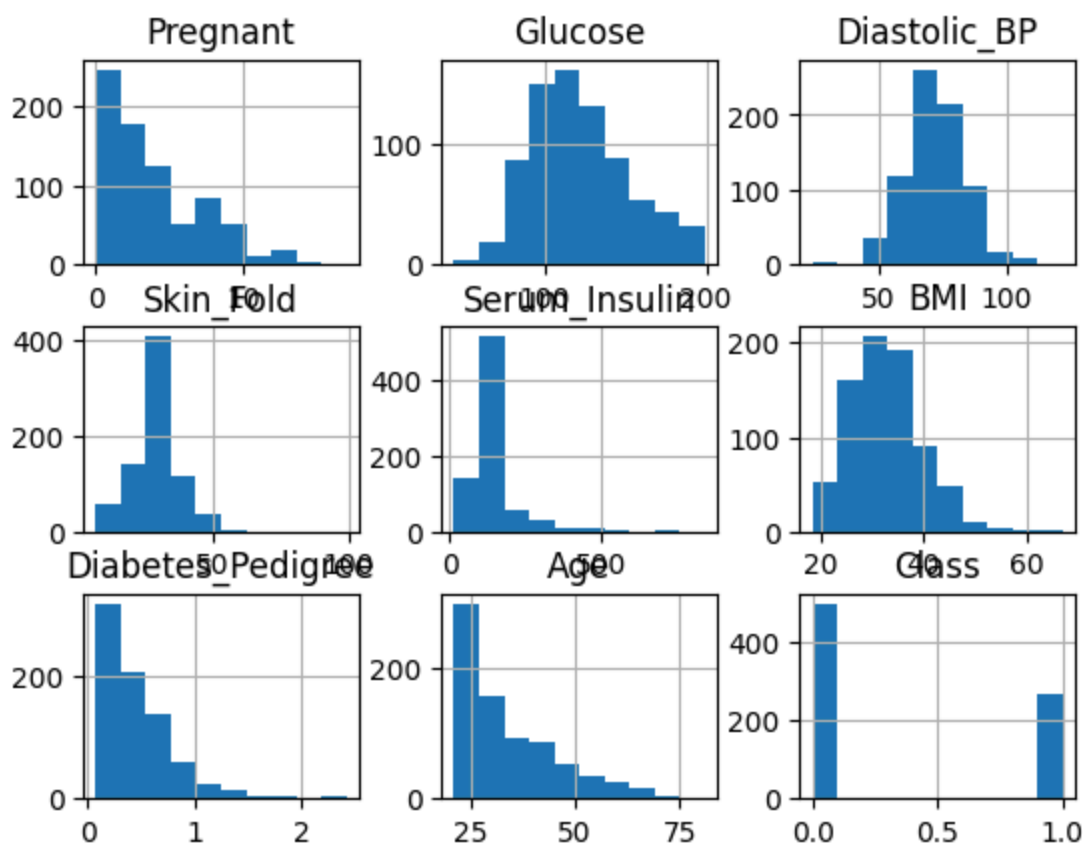


In [38]: dataw.std()

Out[38]: 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

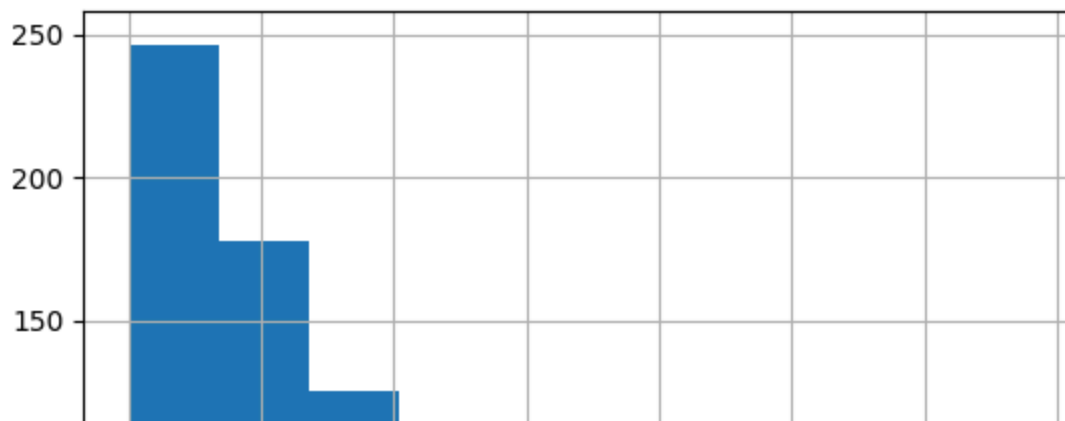
```
In [39]: # Dataset's Histogram representation for every column separately
dataw.hist()
```

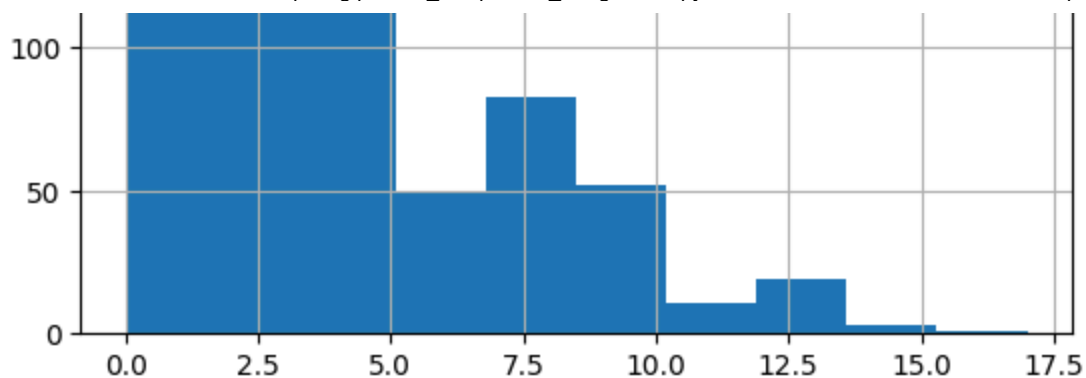
```
Out[39]: array([[<Axes: title={'center': 'Pregnant'}>,
<Axes: title={'center': 'Glucose'}>,
<Axes: title={'center': 'Diastolic_BP'}>],
[<Axes: title={'center': 'Skin_Fold'}>,
<Axes: title={'center': 'Serum_Insulin'}>,
<Axes: title={'center': 'BMI'}>],
[<Axes: title={'center': 'Diabetes_Pedigree'}>,
<Axes: title={'center': 'Age'}>,
<Axes: title={'center': 'Class'}>]], dtype=object)
```



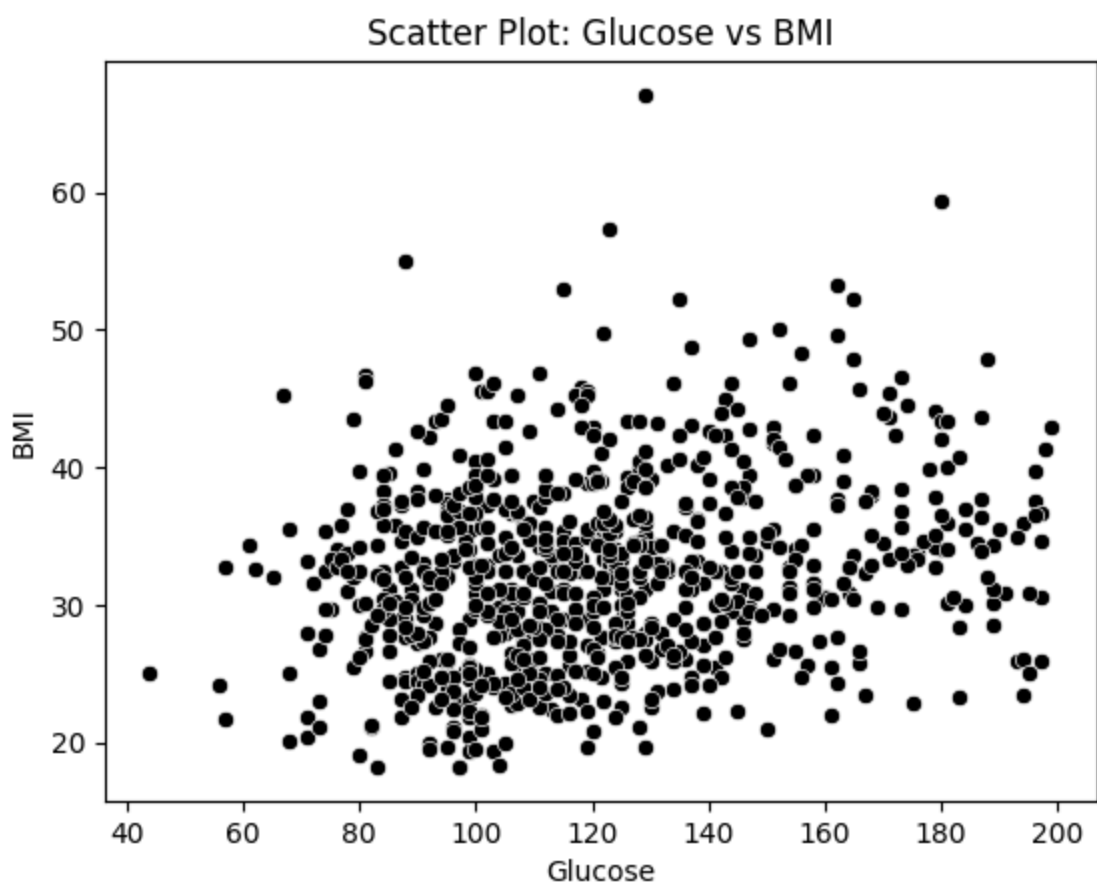
```
In [46]: # Pregnant column data visulization
dataw['Pregnant'].hist()
```

```
Out[46]: <Axes: >
```

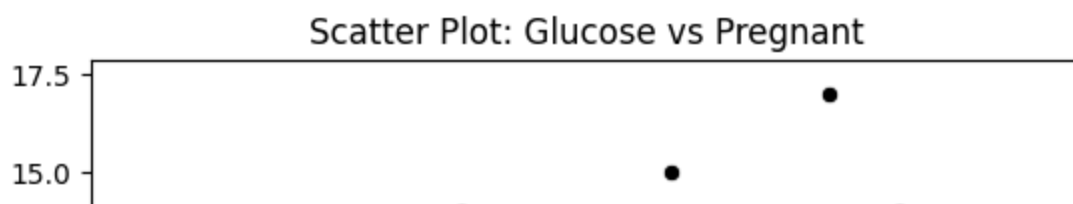


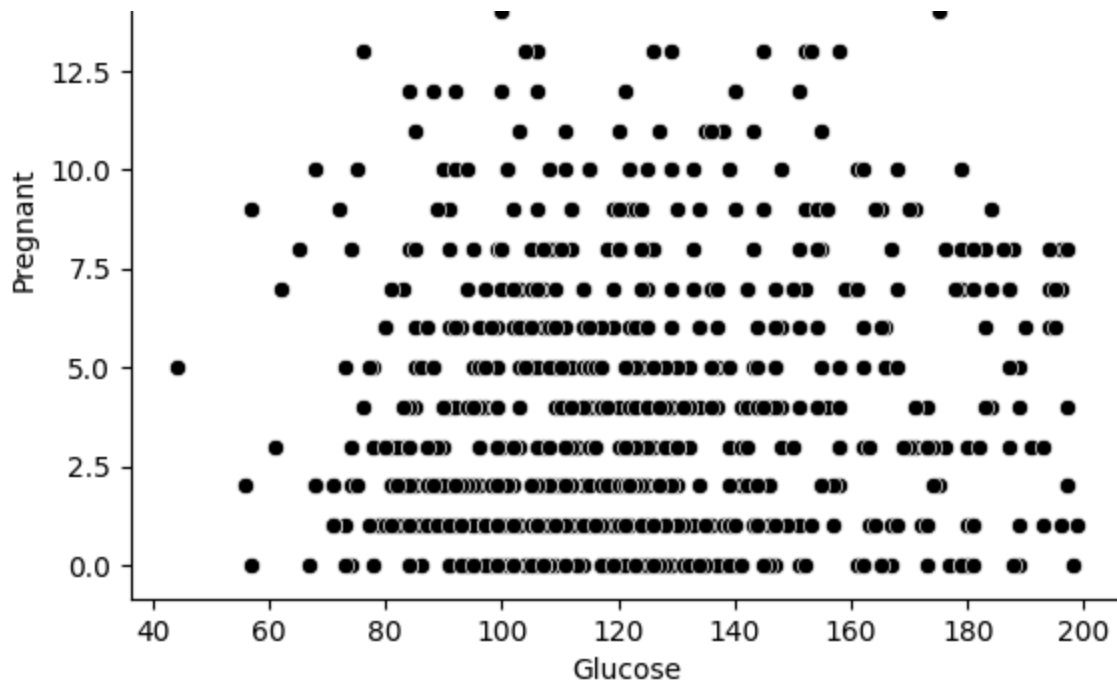


```
In [41]: sns.scatterplot(data=dataw, x="Glucose", y="BMI", color="black")  
  
plt.title("Scatter Plot: Glucose vs BMI")  
plt.show()
```



```
In [42]: sns.scatterplot(data=dataw, x="Glucose", y="Pregnant", color="black")  
  
plt.title("Scatter Plot: Glucose vs Pregnant")  
plt.show()
```





```
In [43]: # Compute the correlation matrix
corr_matrix = dataw.corr()

corr_matrix
```

```
Out[43]:
```

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

```
In [44]: plt.figure(figsize=(10, 6))
sns.heatmap(corr_matrix, annot=True, cmap="coolwarm", fmt=".2f", linewidths=0.5)

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

