# technohacks-task-3-diabeties-1

January 27, 2024

```
[1]: # import the required libraries
     import pandas as pd
     import numpy as np
     import matplotlib.pyplot as plt
     import seaborn as sns
```

```
[2]: #loading the dataset
     df=pd.read_csv("diabetes.csv")
     df
```

```
[2]:      Pregnancies  Glucose  BloodPressure  SkinThickness  Insulin   BMI  \
     0              6      148             72             35        0  33.6
     1              1       85             66             29        0  26.6
     2              8      183             64              0        0  23.3
     3              1       89             66             23       94  28.1
     4              0      137             40             35      168  43.1
     ..           ...      ...            ...            ...      ...   ...
     763           10      101             76             48      180  32.9
     764            2      122             70             27        0  36.8
     765            5      121             72             23      112  26.2
     766            1      126             60              0        0  30.1
     767            1       93             70             31        0  30.4

          DiabetesPedigreeFunction  Age  Outcome
     0                       0.627   50        1
     1                       0.351   31        0
     2                       0.672   32        1
     3                       0.167   21        0
     4                       2.288   33        1
     ..                        ...  ...      ...
     763                     0.171   63        0
     764                     0.340   27        0
     765                     0.245   30        0
     766                     0.349   47        1
     767                     0.315   23        0

     [768 rows x 9 columns]
```

```
[3]: #checking any null values and missing values
     df.isnull().sum()
```

```
[3]: Pregnancies                 0
     Glucose                     0
     BloodPressure               0
     SkinThickness               0
     Insulin                     0
     BMI                         0
     DiabetesPedigreeFunction    0
     Age                         0
     Outcome                     0
     dtype: int64
```

```
[4]: #displaying the statistical information about the dataset
     df.describe()
```

```
[4]:        Pregnancies      Glucose  BloodPressure  SkinThickness      Insulin  \
     count   768.000000   768.000000     768.000000     768.000000   768.000000
     mean      3.845052   120.894531      69.105469      20.536458    79.799479
     std       3.369578    31.972618      19.355807      15.952218   115.244002
     min       0.000000     0.000000       0.000000       0.000000     0.000000
     25%       1.000000    99.000000      62.000000       0.000000     0.000000
     50%       3.000000   117.000000      72.000000      23.000000    30.500000
     75%       6.000000   140.250000      80.000000      32.000000   127.250000
     max      17.000000   199.000000     122.000000      99.000000   846.000000

                   BMI  DiabetesPedigreeFunction         Age     Outcome
     count  768.000000                768.000000  768.000000  768.000000
     mean    31.992578                  0.471876   33.240885    0.348958
     std      7.884160                  0.331329   11.760232    0.476951
     min      0.000000                  0.078000   21.000000    0.000000
     25%     27.300000                  0.243750   24.000000    0.000000
     50%     32.000000                  0.372500   29.000000    0.000000
     75%     36.600000                  0.626250   41.000000    1.000000
     max     67.100000                  2.420000   81.000000    1.000000
```

```
[5]: # checking the duplicates value
     df.duplicated()
```

```
[5]: 0      False
     1      False
     2      False
     3      False
     4      False
            …
     763    False
```

```
764     False
765     False
766     False
767     False
Length: 768, dtype: bool
```

[6]:
```python
correlation_matrix=df.corr()
correlation_matrix
```

[6]:
```
                          Pregnancies   Glucose  BloodPressure  SkinThickness  \
Pregnancies                  1.000000  0.129459       0.141282      -0.081672
Glucose                      0.129459  1.000000       0.152590       0.057328
BloodPressure                0.141282  0.152590       1.000000       0.207371
SkinThickness               -0.081672  0.057328       0.207371       1.000000
Insulin                     -0.073535  0.331357       0.088933       0.436783
BMI                          0.017683  0.221071       0.281805       0.392573
DiabetesPedigreeFunction    -0.033523  0.137337       0.041265       0.183928
Age                          0.544341  0.263514       0.239528      -0.113970
Outcome                      0.221898  0.466581       0.065068       0.074752

                            Insulin       BMI  DiabetesPedigreeFunction  \
Pregnancies               -0.073535  0.017683                 -0.033523
Glucose                    0.331357  0.221071                  0.137337
BloodPressure              0.088933  0.281805                  0.041265
SkinThickness              0.436783  0.392573                  0.183928
Insulin                    1.000000  0.197859                  0.185071
BMI                        0.197859  1.000000                  0.140647
DiabetesPedigreeFunction   0.185071  0.140647                  1.000000
Age                       -0.042163  0.036242                  0.033561
Outcome                    0.130548  0.292695                  0.173844

                               Age   Outcome
Pregnancies               0.544341  0.221898
Glucose                   0.263514  0.466581
BloodPressure             0.239528  0.065068
SkinThickness            -0.113970  0.074752
Insulin                  -0.042163  0.130548
BMI                       0.036242  0.292695
DiabetesPedigreeFunction  0.033561  0.173844
Age                       1.000000  0.238356
Outcome                   0.238356  1.000000
```
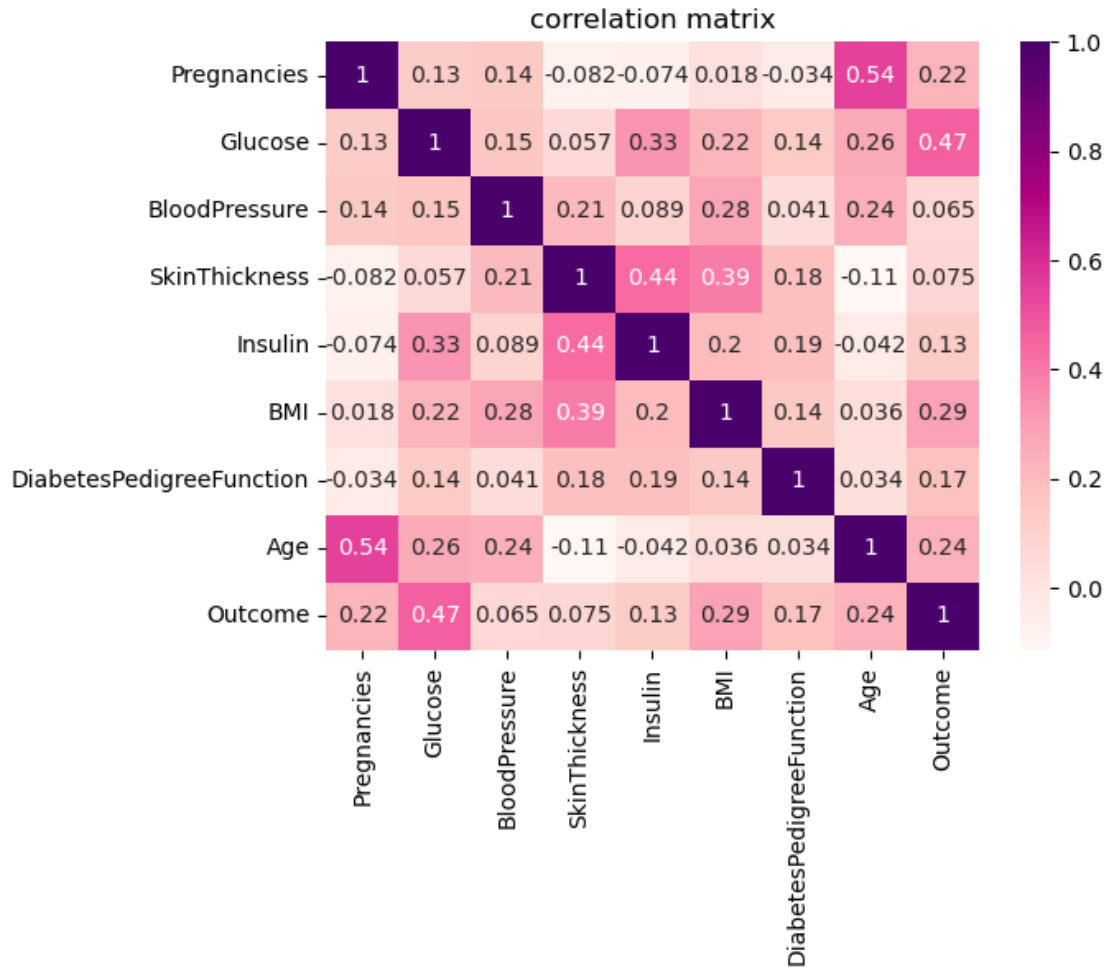
[7]:
```python
sns.heatmap(correlation_matrix,annot=True,cmap='RdPu')
plt.title('correlation matrix')
```

[7]: Text(0.5, 1.0, 'correlation matrix')

## correlation matrix

|  | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeFunction | Age | Outcome |
|---|---|---|---|---|---|---|---|---|---|
| Pregnancies | 1 | 0.13 | 0.14 | -0.082 | -0.074 | 0.018 | -0.034 | 0.54 | 0.22 |
| Glucose | 0.13 | 1 | 0.15 | 0.057 | 0.33 | 0.22 | 0.14 | 0.26 | 0.47 |
| BloodPressure | 0.14 | 0.15 | 1 | 0.21 | 0.089 | 0.28 | 0.041 | 0.24 | 0.065 |
| SkinThickness | -0.082 | 0.057 | 0.21 | 1 | 0.44 | 0.39 | 0.18 | -0.11 | 0.075 |
| Insulin | -0.074 | 0.33 | 0.089 | 0.44 | 1 | 0.2 | 0.19 | -0.042 | 0.13 |
| BMI | 0.018 | 0.22 | 0.28 | 0.39 | 0.2 | 1 | 0.14 | 0.036 | 0.29 |
| DiabetesPedigreeFunction | -0.034 | 0.14 | 0.041 | 0.18 | 0.19 | 0.14 | 1 | 0.034 | 0.17 |
| Age | 0.54 | 0.26 | 0.24 | -0.11 | -0.042 | 0.036 | 0.034 | 1 | 0.24 |
| Outcome | 0.22 | 0.47 | 0.065 | 0.075 | 0.13 | 0.29 | 0.17 | 0.24 | 1 |

```
[8]: # get the unique values for paticular column
     unique_value=df['Outcome'].unique()
     print(unique_value)
     #to get the value count
     df['Outcome'].value_counts()
```

```
[1 0]
```

```
[8]: 0    500
     1    268
     Name: Outcome, dtype: int64
```

```
[9]: # preparing dependent and target varaible
     X=df.drop(columns='Outcome',axis=1)
     y=df['Outcome']
     print(X)
     print(y)
```

```
     Pregnancies  Glucose  BloodPressure  SkinThickness  Insulin   BMI  \
0              6      148             72             35        0  33.6
1              1       85             66             29        0  26.6
2              8      183             64              0        0  23.3
3              1       89             66             23       94  28.1
4              0      137             40             35      168  43.1
..           ...      ...            ...            ...      ...   ...
763           10      101             76             48      180  32.9
764            2      122             70             27        0  36.8
765            5      121             72             23      112  26.2
766            1      126             60              0        0  30.1
767            1       93             70             31        0  30.4

     DiabetesPedigreeFunction  Age
0                       0.627   50
1                       0.351   31
2                       0.672   32
3                       0.167   21
4                       2.288   33
..                        ...  ...
763                     0.171   63
764                     0.340   27
765                     0.245   30
766                     0.349   47
767                     0.315   23

[768 rows x 8 columns]
0      1
1      0
2      1
3      0
4      1
      ..
763    0
764    0
765    0
766    1
767    0
Name: Outcome, Length: 768, dtype: int64
```

```python
#displaying the number of rows and columns in X dataset
X.shape
y.shape
```

```
(768,)
```

```
[12]: #displaying the number of rows and columns in X dataset
      X.shape
      y.shape
```

[12]: (768,)

[ ]: