

## Diabetes prediction

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline

In [2]: df = pd.read_csv('diabetes.csv')
df.head()
```

```
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

```
In [3]: df.shape
Out[3]: (768, 9)
```

```
In [4]: df.info()
<class 'pandas.core.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    float64
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]: df.describe()
Out[5]:
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
count	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000
mean	3.865262	120.924521	89.506489	20.536488	79.799079	31.992578	0.471976	33.248985	0.348998
std	3.359676	31.977504	19.355067	15.952218	115.244002	7.884160	0.331379	11.760232	0.476951
min	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.076900	21.000000	0.000000
25%	1.000000	99.000000	62.000000	0.000000	0.000000	27.300000	0.243750	24.000000	0.000000
50%	3.000000	117.000000	72.000000	23.000000	30.600000	32.000000	0.372500	29.000000	0.000000
75%	6.000000	140.250000	80.000000	32.000000	127.250000	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 [6]: df.isna().sum()
Out[6]:
Pregnancies    0
Glucose         0
BloodPressure   0
SkinThickness   0
Insulin         0
BMI             0
DiabetesPedigreeFunction  0
Age             0
Outcome         0
dtype: int64
```

```
In [7]: df['Pregnancies'].nunique()
Out[7]: 17
```

```
In [8]: plt.figure(figsize=(10, 5))
sns.set_style('whitegrid')
sns.countplot(data=df,
              x='Pregnancies',
              hue='Outcome')
plt.xlabel('No of Pregnancies')
plt.ylabel('Amount of People')
plt.legend(labels=['no diabetes', 'diabetes'])
plt.title('Diabetes Occurrences based on Pregnancies')
```

```
Out[8]:
```

Text(6.5, 1.8, 'Diabetes Occurrences based on Pregnancies')

```
In [9]: plt.figure(figsize=(10, 5))
plt.scatter(x=df['Glucose'],
            y=df['Outcome'],
            color='brown')
plt.xlabel('Glucose Level')
plt.ylabel('0 - no diabetes, 1 - diabetes')
plt.title('Diabetes Occurrences based on Glucose Level')
```

```
Out[9]:
```

Text(6.8, 1.8, 'Diabetes Occurrences based on Glucose Level')

```
In [10]: sns.set_style('whitegrid')
sns.countplot(data=df,
              x='Outcome',
              palette='RdBu_r')
plt.title('Amount of People w Diabetes & Without Diabetes')
plt.xlabel('0 - no diabetes, 1 - diabetes')
plt.ylabel('Number of People')
```

```
Out[10]:
```

Amount of People w Diabetes & Without Diabetes

```
In [11]: df.hist(figsize=(18, 10));
```

```
In [12]: import warnings
warnings.filterwarnings('ignore', category=FutureWarning, module='seaborn._oldcore')
```

```
In [13]: sns.pairplot(data=df,
                   hue='Outcome',
                   palette='winter');
```

```
Out[13]:
```

C:\Users\ushaj\Anaconda3\Lib\site-packages\seaborn\axisgrid.py:123: UserWarning: The figure layout has changed to tight

```
In [14]: df.head(18)
```

```
Out[14]:
```

	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
5	5	116	74	0	0	25.6	0.201	30	0
6	3	78	50	32	88	31.0	0.248	26	1
7	10	115	0	0	0	35.3	0.134	29	0
8	2	197	70	45	543	35.5	0.158	63	1
9	8	125	96	0	0	6.0	0.232	54	1

```
In [15]: df1 = df.copy()
```

```
Out[15]:
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
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 [16]: df1.fill_missing_with_group_mean(df, columns, group_col):
def fill_missing_with_group_mean(df1, columns, group_col):
    for column in columns:
        df1[column] = df1.groupby(group_col)[column].transform(lambda x: x.fillna(x.mean()))
    return df1
```

```
Out[16]:
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
0	6	148.0	72.0	35.0	106.500000	33.6	0.627	50	1
1	1	85.0	66.0	29.0	222.333333	26.6	0.351	31	0
2	8	183.0	64.0	29.1	95.500000	23.3	0.672	32	1
3	1	89.0	66.0	23.0	94.000000	28.1	0.167	21	0
4	0	137.0	40.0	35.0	168.000000	43.1	2.288	33	1

```
768 rows x 9 columns
```

```
In [17]: df1.isna().sum()
Out[17]:
Pregnancies    0
Glucose         0
BloodPressure   0
SkinThickness   1
Insulin         1
BMI             1
DiabetesPedigreeFunction  0
Age             0
Outcome         0
dtype: int64
```

```
In [18]: df1.head()
```

```
Out[18]:
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
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 [19]: df1.fill_missing_with_group_mean(df1, ["Glucose", "BloodPressure", "SkinThickness", "Insulin", "BMI", "Age"], 'Age')
def fill_missing_with_group_mean(df1, columns, group_col):
    for column in columns:
        df1[column] = df1.groupby(group_col)[column].transform(lambda x: x.fillna(x.mean()))
    return df1
```

```
Out[19]:
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
0	6	148.0	72.0	35.0	106.500000	33.6	0.627	50	1
1	1	85.0	66.0	29.0	222.333333	26.6	0.351	31	0
2	8	183.0	64.0	29.1	95.500000	23.3	0.672	32	1
3	1	89.0	66.0	23.0	94.000000	28.1	0.167	21	0
4	0	137.0	40.0	35.0	168.000000	43.1	2.288	33	1

```
768 rows x 9 columns
```

```
In [20]: df1.isna().sum()
Out[20]:
Pregnancies    0
Glucose         0
BloodPressure   0
SkinThickness   1
Insulin         1
BMI             1
DiabetesPedigreeFunction  0
Age             0
Outcome         0
dtype: int64
```

```
In [21]: df1[df1['Insulin'].isna()]
Out[21]:
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
123	5	132.0	80.0	NaN	NaN	26.8	0.196	69	0
129	0	195.0	84.0	58.0	NaN	27.9	0.741	62	1
148	5	147.0	78.0	37.0	NaN	33.7	0.218	65	0
221	2	158.0	90.0	NaN	NaN	31.6	0.805	66	1
294	0	161.0	50.0	37.0	NaN	21.9	0.254	65	0
362	5	103.0	108.0	37.0	NaN	39.2	0.305	65	0
363	4	146.0	78.0	NaN	NaN	26.8	0.520	67	1
453	2	119.0	NaN	NaN	NaN	19.6	0.832	72	0
495	1	135.0	54.0	58.0	NaN	26.7	0.687	62	0
489	8	194.0	80.0	NaN	NaN	26.1	0.551	67	0
495	6	166.0	74.0	NaN	NaN	26.6	0.304	66	0
509	8	120.0	78.0	NaN	NaN	25.0	0.409	64	0
537	0	57.0	60.0	NaN	NaN	21.7	0.735	67	0
552	6	114.0	88.0	NaN	NaN	27.8	0.247	66	1
579	2	119.0	70.0	99.0	NaN	34.7	0.575	62	1
582	12	121.0	78.0	17.0	NaN	26.5	0.259	62	0
666	4	145.0	82.0	18.0	NaN	32.5	0.235	70	1
674	8	91.0	82.0	NaN	NaN	39.6	0.587	68	0
684	5	136.0	82.0	NaN	NaN	28.8	0.640	69	0
759	6	190.0	92.0	NaN	NaN	35.5	0.278	66	1

```
In [22]: median_insulin_by_bmi = df1.groupby('BMI')['Insulin'].median()
def fill_missing_with_group_mean(df1, columns, group_col):
    for column in columns:
        df1[column] = df1.groupby(group_col)[column].transform(lambda x: x.fillna(x.mean()))
    return df1
```

```
Out[22]:
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
123	5	132.0	80.0	NaN	NaN	26.8	0.196	69	0
129	0	195.0	84.0	58.0	NaN	27.9	0.741	62	1
148	5	147.0	78.0	37.0	NaN	33.7	0.218	65	0
221	2	158.0	90.0	NaN	NaN	31.6	0.805	66	1
294	0	161.0	50.0	37.0	NaN	21.9	0.254	65	0
362	5	103.0	108.0	37.0	NaN	39.2	0.305	65	0
363	4	146.0	78.0	NaN	NaN	26.8	0.520	67	1
453	2	119.0	NaN	NaN	NaN	19.6	0.832	72	0
495	6	166.0	74.0	NaN	NaN	26.6	0.304	66	0
509	8	120.0	78.0	NaN	NaN	25.0	0.409	64	0
537	0	57.0	60.0	NaN	NaN	21.7	0.735	67	0
552	6	114.0	88.0	NaN	NaN	27.8	0.247	66	1
579	2	119.0	70.0	99.0	NaN	34.7	0.575	62	1
582	12	121.0	78.0	17.0	NaN	26.5	0.259	62	0
666	4	145.0	82.0	18.0	NaN	32.5	0.235	70	1
674	8	91.0	82.0	NaN	NaN	39.6	0.587	68	0
684	5	136.0	82.0	NaN	NaN	28.8	0.640	69	0
759	6	190.0	92.0	NaN	NaN	35.5	0.278	66	1

```
In [23]: median_insulin_by_bmi = df1.groupby('BMI')['Insulin'].median()
def fill_missing_with_group_mean(df1, columns, group_col):
    for column in columns:
        df1[column] = df1.groupby(group_col)[column].transform(lambda x: x.fillna(x.mean()))
    return df1
```

```
Out[23]:
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
123	5	132.0	80.0	NaN	NaN	26.8	0.196	69	0
129	0	195.0	84.0	58.0	NaN	27.9	0.741	62	1
148	5	147.0	78.0	37.0	NaN	33.7	0.218	65	0
221	2	158.0	90.0	NaN	NaN	31.6	0.805	66	1
294	0	161.0	50.0	37.0	NaN	21.9	0.254	65	0
362	5	103.0	108.0	37.0	NaN	39.2	0.305	65	0
363	4	146.0	78.0	NaN	NaN	26.8	0.520	67	1
453	2	119.0	NaN	NaN	NaN	19.6	0.832	72	0
495	6	166.0	74.0	NaN	NaN	26.6	0.304	66	0
509	8	120.0	78.0	NaN	NaN	25.0	0.409	64	0
537	0	57.0	60.0	NaN	NaN	21.7	0.735	67	0
552	6	114.0	88.0	NaN	NaN	27.8	0.247	66	1
579	2	119.0	70.0	99.0	NaN	34.7	0.575	62	1
582	12	121.0	78.0	17.0	NaN	26.5	0.259	62	0
666	4	145.0	82.0	18.0	NaN	32.5	0.235	70	1
674	8	91.0	82.0	NaN	NaN	39.6	0.587	68	0
684	5	136.0	82.0	NaN	NaN	28.8	0.640	69	0
759	6	190.0	92.0	NaN	NaN	35.5	0.278	66	1

```
In [24]: df1.isna().sum()
Out[24]:
Pregnancies    0
Glucose         0
BloodPressure   0
SkinThickness   1
Insulin         1
BMI             1
DiabetesPedigreeFunction  0
Age             0
Outcome         0
dtype: int64
```

```
In [25]: df1[df1['SkinThickness'].isna()]
Out[25]:
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
123	5	132.0	80.0	NaN	NaN	26.8	0.196	69	0
129	0	195.0	84.0	58.0	NaN	27.9	0.741	62	1
148	5	147.0	78.0	37.0	NaN	33.7	0.218	65	0
221	2	158.0	90.0	NaN	NaN	31.6	0.805	66	1
294	0	161.0	50.0	37.0	NaN	21.9	0.254	65	0
362	5	103.0	108.0	37.0	NaN	39.2	0.305	65	0
363	4	146.0	78.0	NaN	NaN	26.8	0.520	67	1
453	2	119.0	NaN	NaN	NaN	19.6	0.832		