

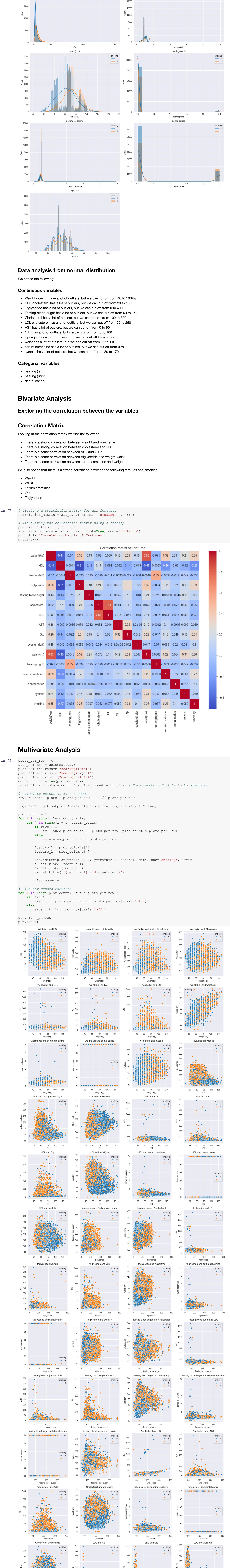
```
In [2]: import pandas as pd
data_path = "tahkine.csv"
all_data = pd.read_csv(data_path)

# nicotine dependency, carbon monoxide levels, daily cigarette consumption, age of
# smoking initiation, previous quit attempts, emotional well-being, personality traits, and motivation
# to
# quit

columns = all_data.columns.tolist()[1:]
columns.remove("smoking")
columns

Out[2]: ['weight(kg)', 'HDL', 'triglyceride', 'fasting blood sugar', 'Cholesterol', 'LDL', 'AST', 'Gtp', 'eyesight(left)', 'waist(cm)', 'hearing(right)', 'serum creatinine', 'dental caries', 'systolic']
```

Univariate Analysis



Data analysis from normal distribution

We notice the following:

Continuous variables

- Weight doesn't have a lot of outliers, but we can cut off from 40 to 100KG

- HDL cholesterol has a lot of outliers, but we can cut off from 20 to 100

- Triglyceride has a lot of outliers, but we can cut off from 0 to 400

- Fasting blood sugar has a lot of outliers, but we can cut off from 60 to 150

- Cholesterol has a lot of outliers, but we can cut off from 100 to 300

- LDL cholesterol has a lot of outliers, but we can cut off from 20 to 250

- AST has a lot of outliers, but we can cut off from 0 to 90

- GTP has a lot of outliers, but we can cut off from 0 to 180

- Eyesight has a lot of outliers, but we can cut off from 0 to 2

- Waist has a lot of outliers, but we can cut off from 55 to 110

- Serum creatinine has a lot of outliers, but we can cut off from 0 to 2

- Systolic has a lot of outliers, but we can cut off from 80 to 170

Categorical variables

- hearing (left)
- hearing (right)
- dental caries

Bivariate Analysis

Exploring the correlation between the variables

Correlation Matrix

Looking at the correlation matrix we find the following:

- There is a strong correlation between weight and waist size
- There is a strong correlation between cholesterol and LDL
- There is a some correlation between AST and GTP
- There is a some correlation between triglyceride and weight-waist
- There is a some correlation between serum creatinine and weight

We also notice that there is a strong correlation between the following features and smoking:

- Weight
- Waist
- Serum creatinine
- Gtp
- Triglyceride

Multivariate Analysis

