| **Name: Salma Alarfaj** |
| --- |

| **Purpose** |
| --- |
| To get exposed to a very large data set that has more than 10 features and more than 10000 records and apply data cleaning, Machine learning Algorithms, and statistical processes to better study the data and practice data analysis and science. |
| **Picked Dataset** |
| **Cardiovascular Disease dataset**<br>The dataset consists of 70 000 records of patient's data, 11 features + target. |
| **Target Variable** |
| Presence or absence of cardiovascular disease (binary). |
| **Features** |
| 1. Age \| age \| int (days)<br>2. Height \| height \| int (cm)<br>3. Weight \| weight \| float (kg)<br>4. Gender \| gender \| categorical code<br>5. Systolic blood pressure \| ap_hi \| int<br>6. Diastolic blood pressure \| ap_lo \| int<br>7. Cholesterol \| cholesterol \| 1: normal, 2: above normal, 3: well above normal<br>8. Glucose \| gluc \| 1: normal, 2: above normal, 3: well above normal \|<br>9. Smoking \| smoke \| binary \|<br>10. Alcohol intake \| alco \| binary \|<br>11. Physical activity \| active \| binary \| |
| **Analysis** |
| The most common three factors that usually accompany **cardiovascular diseases.**<br>Run statistical description on the data.<br>The age-groups commonly carry **cardiovascular diseases.**<br>More analysis will be carried out after exploring the data closely. |