

Task 1

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from scipy import stats
```

```
In [2]: cardio_data_train = pd.read_csv("cardio=train.csv", sep = ";")
cardio_data_train.head()
```

	id	age	gender	height	weight	ap_hi	ap_lo	cholesterol	gluc	smoke	alco	active	cardio	
0	92160	105360	NaN	163.0	75.0	120.0	NaN	NaN	Normal	0.0	0.0	1.0	0	
1	88880	16725.0	Men	168.0	68.0	110.0	70.0	NaN	NaN	NaN	0.0	0.0	NaN	0
2	10483	19761.0	Men	170.0	75.0	120.0	NaN	NaN	Normal	0.0	0.0	1.0	1	
3	22798	19035.0	Men	NaN	89.0	NaN	80.0	Normal	NaN	0.0	NaN	NaN	0	
4	85542	NaN	NaN	158.0	74.0	NaN	93.0	High	NaN	0.0	1.0	1.0	1	

```
print("Number of features in the cardio data: ",len(cardio_data_train.columns))
print('Features: ')
print(n*10)
features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

for i in range(len(cardio_data_train.columns)):
    print(cardio_data_train.columns[i]+"-----> "+features_desc_from_kaggle[i])

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Systolic blood pressure","Diastolic blood pressure","Cholesterol","Glucose","Smoking","Alcohol intake","Physical activity","Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]

features_desc_from_kaggle = ["id","Age(days)","Gender","Height(cm)","Weight(kg)","Syst
```

```
In [3]: print("Number of features in the cardio data: ", len(cardio_data_train.columns))
print("Features: ")
print("----")
print("id")
features_desc_from_kaggle = ["id", "Age(days)", "Gender", "Height(cm)", "Weight(kg)", "Systolic blood pressure",
                             "Diastolic blood pressure", "Cholesterol", "Glucose", "Smoking", "\n",
                             "Alcohol intake", "Physical activity", "Presence(1) or absence(0) of cardiovascular disease (Target Variable)"]
for i in range(len(cardio_data_train.columns)):
    print(cardio_data_train.columns[i])
    print("-----")
    print(features_desc_from_kaggle[i])
```

Number of features in the cardio data: 13
Features:

id -----> id
age -----> Age (days)
gender -----> Gender
height -----> Height (cm)
weight -----> Weight (kg)
ap_hi -----> Systolic blood pressure
ap_lo -----> Diastolic blood pressure
cholesterol -----> Cholesterol
gluc -----> Glucose
smoke -----> Smoking
alco -----> Alcohol intake
active -----> Physical activity
cardio -----> Presence (1) or absence (0) of cardiovascular disease (Target Variable)

1. Identify the dataset columns into nominal, categorical, continues etc. categories

- **Nominal (Categorical) features:** gender, cholesterol, gluc, smoke, alco, active
- **Numeric (Continuous) Features:** age, height, weight, ap_hi, ap_lo
- **Target variable:** cardio

2. Use dataframe.info and dataframe.describe to get the insights about the data.

```
In [4]: cardio_data_train.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 500 entries, 0 to 499
Data columns (total 13 columns):
 #   Column              Non-Null Count  Dtype
---  --
 0   id                   500 non-null    int64
 1   age                  335 non-null    float64
 2   gender               329 non-null    object
 3   height              198 non-null    float64
 4   weight              336 non-null    float64
 5   ap_hi               347 non-null    float64
 6   ap_lo               332 non-null    float64
 7   cholesterol         333 non-null    object
 8   gluc                333 non-null    object
 9   smoke               326 non-null    float64
10   alco                335 non-null    float64
11   active              343 non-null    float64
12   cardio              500 non-null    int64
dtypes: float64(8), int64(2), object(3)
memory usage: 50.9+ KB

Observations:

• We can see presence of NaN(missing) values in all columns except 'id' and 'cardio'.
• Most NaN values are in 'height' column, as non-null count is very low.
• Float values are present in 'age', 'height', 'weight', 'ap_hi', 'ap_lo', 'smoke', 'alco', and 'active'. These columns are numeric and continuous. Integer values are present in 'id' and 'cardio' columns.
• For columns 'gender', 'cholesterol' and 'gluc' we see data type as object because these are categorical features with string values.
```

```
In [5]: cardio_data_train.describe()

      id          age      height      weight      ap_hi      ap_lo      smoke      alco      active      cardio
count  500.000000    335.000000    198.000000    336.000000    347.000000    332.000000    326.000000    335.000000    343.000000    500.000000
mean   50279.946000   1940.886667   163.034343   74.347321   128.686879   90.060241   0.090205   0.065672   0.813411   0.50
std    29913.628731   2466.702487   8.265559   14.35964   18.490176   87.396945   0.289505   0.248078   0.390150   0.50
min     38.000000    14334.000000    120.000000    45.000000    60.000000    60.000000    0.000000    0.000000    0.000000    0.00
25%    23446.500000   1798.500000    159.250000    66.000000    120.000000    120.000000    0.000000    0.000000    1.000000    0.00
50%    51913.500000   1979.500000    165.000000    72.000000    120.000000    120.000000    0.000000    0.000000    1.000000    1.00
75%    78656.000000   2197.500000    168.000000    82.000000    140.000000    140.000000    0.000000    0.000000    1.000000    1.00
max    99662.000000   23479.000000    187.000000    155.000000    190.000000    1000.000000    1.000000    1.000000    1.000000    1.00
```

Observations:

- In the output of describe function, we can see basic stats like min, max, 25%, 50%, 75%, mean and std values for numerical columns in our dataset.
- Count represents number of present(non-null) values in each column.
- We can see that range of each column is very different. I.e. 'age' has very high values compared to 'active' or 'smoke' (0/1 values). We can guess that 'age' is given in days. 'ap_lo' has max value which is very high compared to other values of same column, which suggests presence of outlier(s).

3. Find the number of null values for each columns

```
In [6]: for col in cardio_data_train.columns:
        print("Number of null values in column " + col + " : ", cardio_data_train[col].isnull().sum())

number of null values in column id: 0
number of null values in column age: 165
number of null values in column gender: 171
number of null values in column height: 302
number of null values in column weight: 164
number of null values in column ap_hi: 164
number of null values in column ap_lo: 168
number of null values in column cholesterol: 167
number of null values in column gluc: 167
number of null values in column smoke: 174
number of null values in column alco: 165
number of null values in column active: 157
number of null values in column cardio: 0
```

4. Know about the patients (Example of analysis for ages)

```
In [7]: #age is given in days, we can convert it to years by age = age/365
cardio_data_train['age'] = (cardio_data_train['age']/365).round()
data_age = cardio_data_train['age'].round()
print("A. Oldest person in the data is : %0.1f years" % data_age['max'])
print("B. Youngest person in the data is : %0.1f years" % data_age['min'])
print("C. Average age of a person in the data is : %0.1f years" % data_age['mean'])
print("D. Median age of a person in the data is : %0.1f years" % data_age['50%'])

a. Oldest person in the data is : 64.0 years
b. Youngest person in the data is : 39.0 years
c. Average age of a person in the data is : 53.4 years
d. Median age of a person in the data is : 54.0 years
```

e. Find the relationship between the cardio and ages (the cardio column is your prediction variable)

```
In [8]: age_non_null = pd.DataFrame(cardio_data_train[cardio_data_train.age.notnull()])
sns.kdeplot(
    data=age_non_null.loc[(age_non_null['cardio'] == 1), 'age'],
    color = "darkturquoise",
    shade = True
)
sns.kdeplot(
    data=age_non_null.loc[(age_non_null['cardio'] == 0), 'age'],
    color = "lightcoral",
    shade = True
)
plt.legend(['cardio', 'No cardio'])
plt.title('age vs cardio')
plt.xlim(0, 100)
plt.show()
```



Observation:

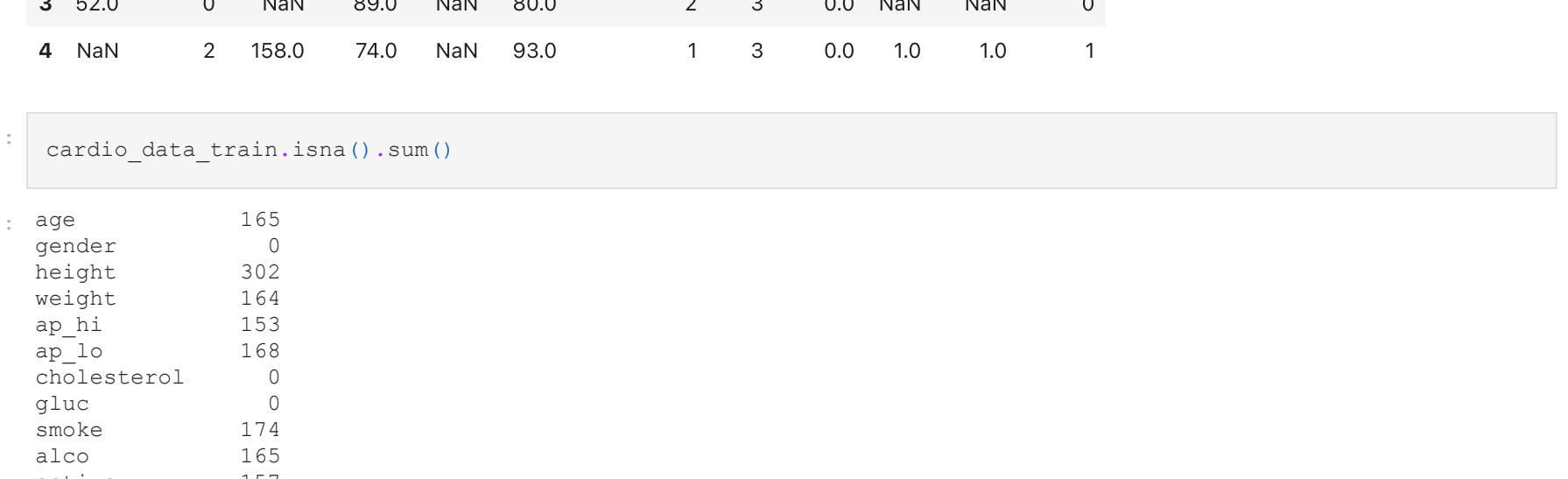
- As we can see from the density plot, age density distribution of people with cardio disease is on the right side of density distribution of people with no cardio disease.
- This suggests that people with higher age(>55) have more risk for cardio disease than people with lesser age(<45).

f. Find the age groups whose survival rate is the largest

```
In [9]: # using only the records that have value for 'age' column
bins = ["30-35", "35-40", "40-45", "45-50", "50-55", "55-60", "60-65", "65-70", "70-75"]
lo = [30, 35, 40, 45, 50, 55, 60, 65, 70]
hi = [35, 40, 45, 50, 55, 60, 65, 70, 75]
counts = [0] * len(bins)
survivors = [0] * len(bins)
ages = pd.DataFrame(age_non_null['age'])
cardio = pd.DataFrame(age_non_null['cardio'])
bins_age = [''] * len(ages)
print("Testing bins of ", str(len(ages)), "people.")
for i in range(len(ages)):
    age = ages.values[i]
    cardio = cardio.values[i]
    for j in range(len(bins)):
        if age >= lo[j] and age < hi[j]:
            counts[j] += 1
            bins_age[i] = bins[j]
            if cardio == 0:
                survivors[j] += 1
            break
for i in range(len(bins)):
    sur_rate = 0
    if counts[i] != 0:
        sur_rate = survivors[i]/counts[i]
    print("Age group: ", bins[i], ", total people: ", counts[i], ", people with no cardio: ", survivors[i], ", survival rate: ", sur_rate)
```

creating bins of 333 people.
In age group 30-35 : total people: 0 , people with no cardio: 0 , survival rate: 0
In age group 35-40 : total people: 1 , people with no cardio: 1 , survival rate: 1.0
In age group 40-45 : total people: 48 , people with no cardio: 36 , survival rate: 0.75
In age group 45-50 : total people: 37 , people with no cardio: 18 , survival rate: 0.49
In age group 50-55 : total people: 68 , people with no cardio: 51 , survival rate: 0.53
In age group 55-60 : total people: 73 , people with no cardio: 34 , survival rate: 0.47
In age group 60-65 : total people: 80 , people with no cardio: 22 , survival rate: 0.28
In age group 65-70 : total people: 0 , people with no cardio: 0 , survival rate: 0
In age group 70-75 : total people: 0 , people with no cardio: 0 , survival rate: 0
Age group 40-45 has the most survival rate.

Note: Ignoring the 1 person in age group 35-40 as we cannot generalize survival rate for that age group.



Observations:

- From countplot above, we can see the distribution of people in age groups of 5 years.
- People can be divided into bins from age 35-40 to 60-65 as we have people from age 39 to 64 years in our data.
- We have highest count of people from age group 60-65.
- 75% of people from age group 40-45 (36 out of 48) have cardio=0, which makes it age group with highest survival rate.

Let's replace categorical features having non-numerical values with integer values.

```
In [11]: print("cholesterol unique values: ", cardio_data_train["cholesterol"].unique())
print("gluc unique values: ", cardio_data_train["gluc"].unique())
print("gender unique values: ", cardio_data_train["gender"].unique())

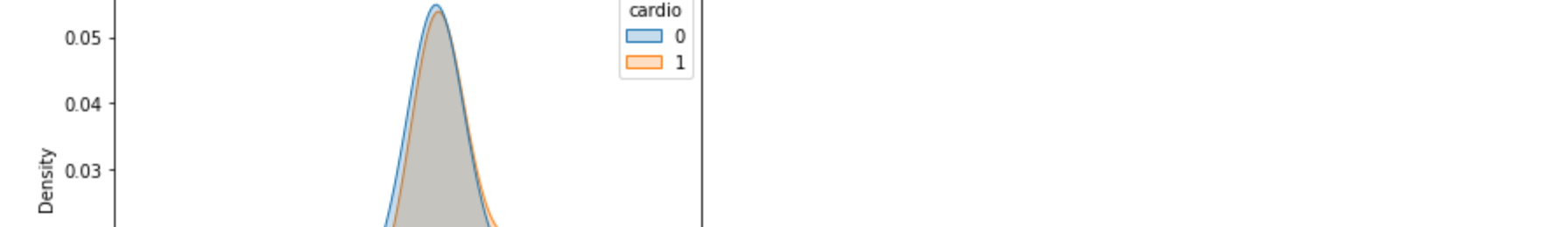
cholesterol unique values: [nan 'Normal' 'High' 'Above Normal']
glucose unique values: [nan 'Normal' nan 'High' 'Above Normal']
gender unique values: [nan 'Men' 'Women']
```

```
In [12]: from sklearn.preprocessing import LabelEncoder
cardio_data_train["cholesterol"] = LabelEncoder().fit_transform(cardio_data_train["cholesterol"])
cardio_data_train["gluc"] = LabelEncoder().fit_transform(cardio_data_train["gluc"])
cardio_data_train["gender"] = LabelEncoder().fit_transform(cardio_data_train["gender"])
```

```
In [13]: #setting the number of points each class has
class_0 = len(cardio_data_train[cardio_data_train["cardio"] == 0])
class_1 = len(cardio_data_train[cardio_data_train["cardio"] == 1])

print("Train data: ")
print("Number of data points for class_0: %d" % class_0)
print("Number of data points for class_1: %d" % class_1)

cardio_data_train.groupby("cardio")["id"].count().plot.bar()
plt.ylabel("Count")
plt.title("Count per class in Train data")
plt.show()
```



```
In [14]: #dropping the "id" column
cardio_data_train.drop("id", axis=1, inplace=True)
```

```
In [15]: cardio_data_train.head()
```

	age	gender	height	weight	ap_hi	ap_lo	cholesterol	gluc	smoke	alco	active	cardio
0	56.0	2	163.0	75.0	120.0	NaN	3	2	0.0	0.0	1.0	0
1	44.0	0	168.0	68.0	110.0	70.0	3	3	0.0	0.0	NaN	0
2	54.0	2	170.0	75.0	120.0	NaN	3	2	0.0	0.0	1.0	1
3	52.0	0	NaN	89.0	NaN	80.0	2	3	0.0	NaN	NaN	0
4	NaN	2	158.0	74.0	NaN	93.0	1	3	0.0	1.0	1.0	1

```
In [16]: cardio_data_train.isna().sum()

[]
```

```
Out [16]: age          0
gender         0
height        162
weight        164
ap_hi         153
ap_lo         168
cholesterol    0
gluc           0
smoke         174
alco          165
active        157
cardio         0
dtype: int64
```

```
In [17]: #replacing NaN values with mean in numerical before finding relationships between features and target variable
#replacing NaN values with mode in categorical attributes

from scipy import stats

numerical_columns = ["age", "ap_hi", "ap_lo", "smoke", "alco", "active", "height", "weight"]
for col in numerical_columns:
    cardio_data_train[col] = cardio_data_train[col].fillna(stats.mode(cardio_data_train[col])[0][0], inplace=True)

categorical_columns = ["cholesterol", "gluc", "gender"]
for col in categorical_columns:
    cardio_data_train[col] = cardio_data_train[col].fillna(stats.mode(cardio_data_train[col])[0][0], inplace=True)

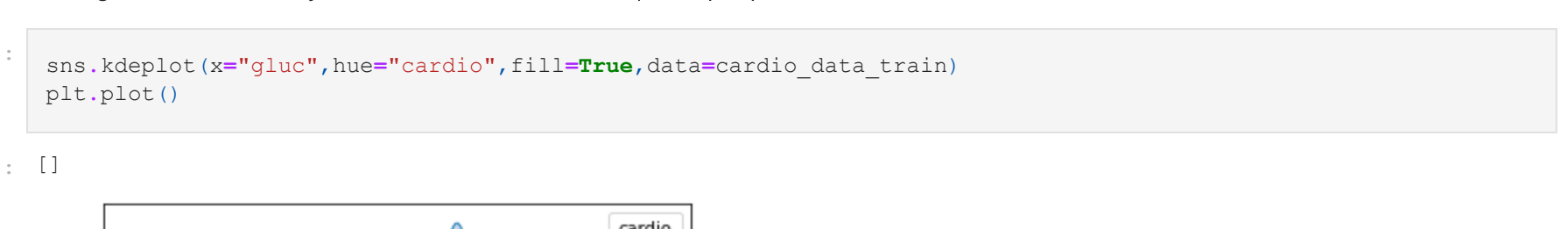
cardio_data_train.isna().sum()
```

```
Out [17]: age          0
gender         0
height         0
weight         0
ap_hi          0
ap_lo          0
cholesterol     0
gluc            0
smoke           0
alco            0
active          0
cardio          0
dtype: int64
```

g. Find similar relationships for at least 3-4 columns that you think can play a role in prediction (For example, systolic BP, cholesterol etc.)

Using density plots to find relationships between few features and target variable. Note: All the NaN values have been replaced with mean(numeric features)/mode(categorical features) in previous cell. So now, in following plots, we will use entire training set (500 rows).

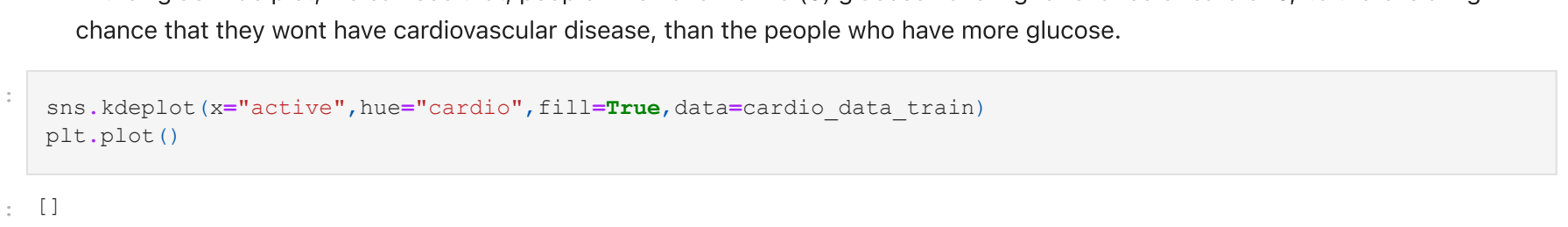
```
In [18]: sns.kdeplot(x="age", hue="cardio", fill=True, data=cardio_data_train)
plt.show()
```



Observations:

- We can see in the kde plot that people around the age 35-45 have higher chance of cardio=0.
- People around the age 50-65 have high survival rate as well as high fatality rate.
- People with age 57-70 have high chance of having cardio=1.

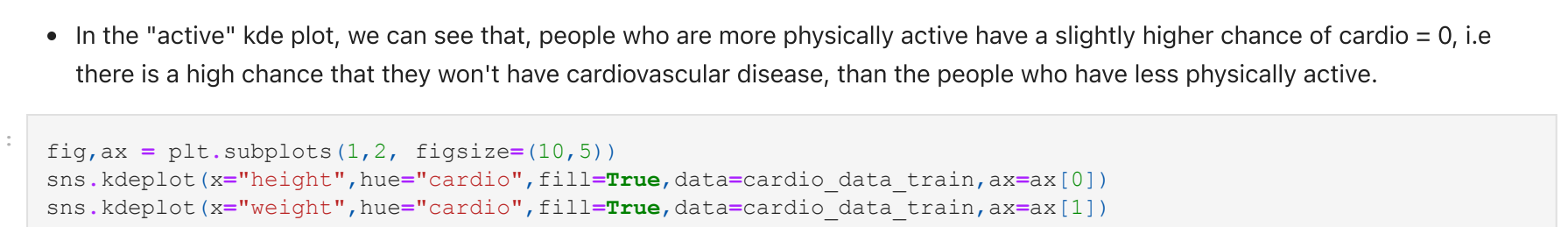
```
In [19]: sns.kdeplot(x="ap_hi", hue="cardio", fill=True, data=cardio_data_train)
plt.plot()
```



Observations:

- People with average Systolic blood pressure(ap_hi) have higher chance of having cardio=0, i.e. there is higher chance that they don't have cardio disease.

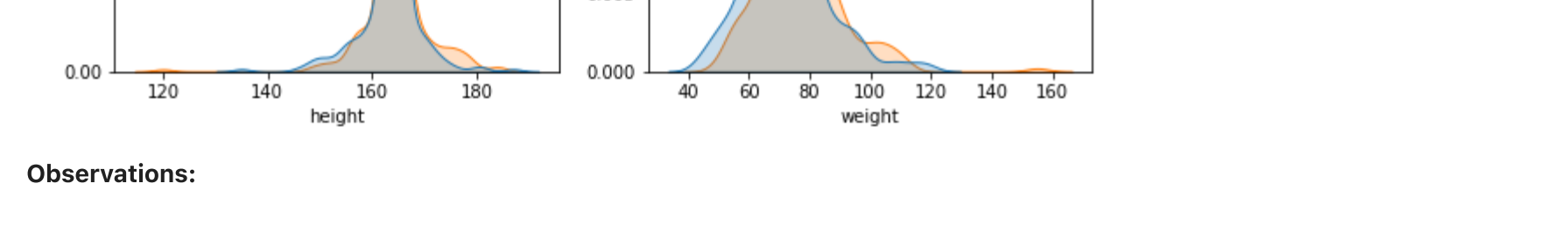
```
In [20]: sns.kdeplot(x="ap_lo", hue="cardio", fill=True, data=cardio_data_train)
plt.plot()
```



Observations:

- We can see from density plot of ap_lo that people with average Diastolic blood pressure(ap_lo) have very high probability of not having cardio disease.

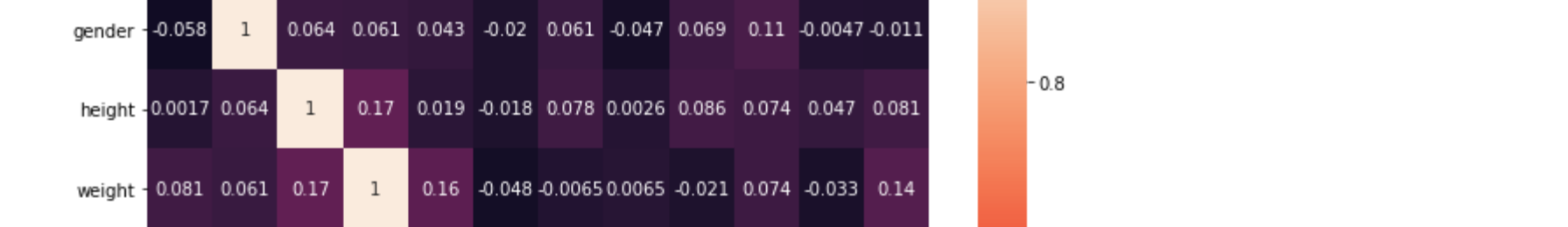
```
In [21]: sns.kdeplot(x="cholesterol", hue="cardio", fill=True, data=cardio_data_train)
plt.plot()
```



Observations:

- From density plot of cholesterol, we can see that people with Normal(0) cholesterol have higher chance of cardio=0, i.e. there is high chance that they will not have cardio disease, than people who have more cholesterol.

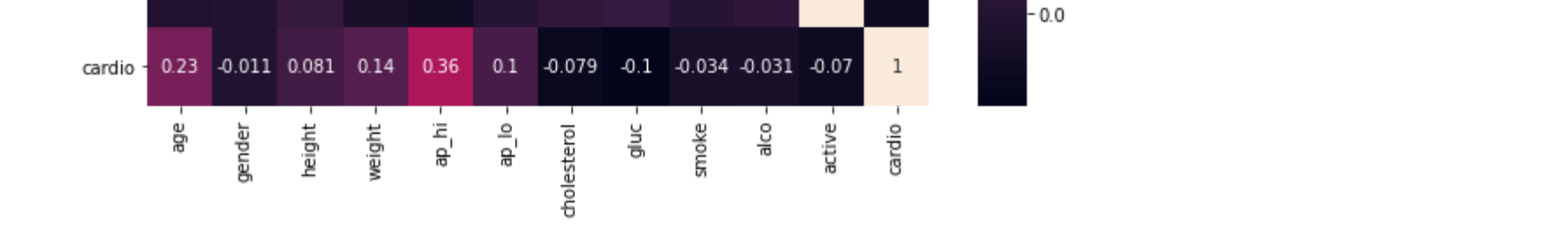
```
In [22]: sns.kdeplot(x="gluc", hue="cardio", fill=True, data=cardio_data_train)
plt.plot()
```



Observations:

- In the "gluc" kde plot, we can see that, people who have Normal(0) glucose have higher chance of cardio=0, i.e. there is a high chance that they won't have cardiovascular disease, than the people who have more glucose.

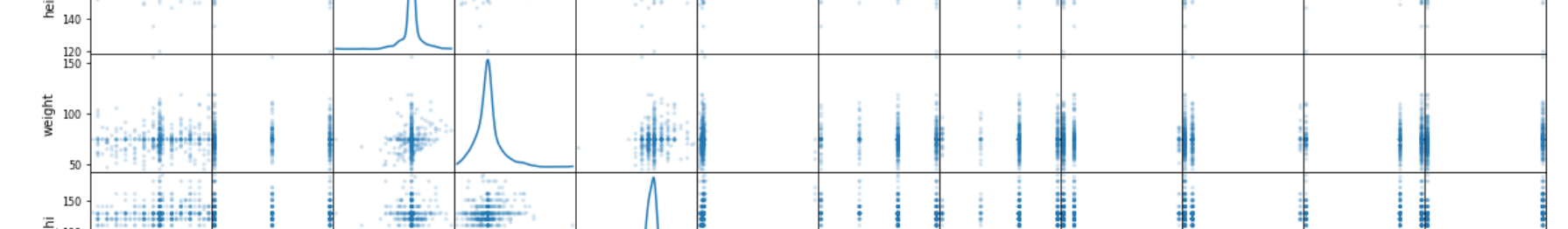
```
In [23]: sns.kdeplot(x="active", hue="cardio", fill=True, data=cardio_data_train)
plt.plot()
```



Observations:

- In the "active" kde plot, we can see that, people who are more physically active have a slightly higher chance of cardio = 0, i.e. there is a high chance that they won't have cardiovascular disease, than the people who have less physical active.

```
In [24]: fig, ax = plt.subplots(1, 2, figsize=(10, 5))
sns.kdeplot(x="height", hue="cardio", fill=True, data=cardio_data_train, ax=ax[0])
sns.kdeplot(x="weight", hue="cardio", fill=True, data=cardio_data_train, ax=ax[1])
plt.plot()
```

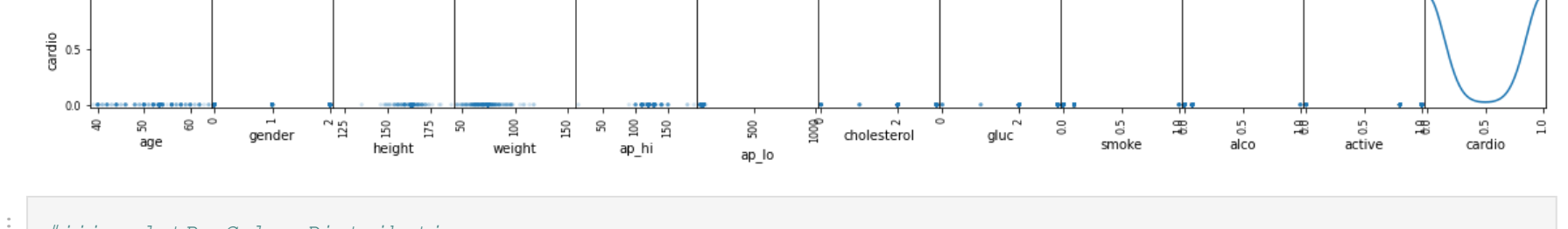


Observations:

- In kde plot of height, the plot of cardio=0 and cardio=1 is roughly the same and from this we can also tell height column does not add much importance, so we can drop height column when we are building our model.
- In kde plot of weight, the plot of cardio=0 and cardio=1 is roughly the same and from this we can also tell weight column does not add much importance, so we can drop weight column when we are building our model.

h. Get more visuals on data distributions

```
In [25]: #1. Use plotCorrelationMatrix
corr_matrix = cardio_data_train.corr()
sns.heatmap(corr_matrix, annot=True)
plt.show()
```



Observations:

- Above plot shows the heatmap of null values per column.
- We can clearly see that 'height' column has 8 most null values as it has most light color.
- 'cardio' and 'id' has entire column with black color, which indicates 0 missing values for these columns.

```
In [26]: # custom label encoder methods for categorical columns
#High, Above Normal, Normal=0
def replace02(x):
    if x == "High":
        return 1
    elif x == "Above Normal":
        return 1
    elif x == "Normal":
        return 0
    else:
        return 0
#Women=0, Men=1
def replace01(x):
    if x == "Women":
        return 0
    elif x == "Men":
        return 1
    else:
        return 0
```

```
In [27]: # turning categorical string values into numerical values
cardio_data_train.replace("cholesterol", cardio_data_train.replace("cholesterol").apply(replace02))
cardio_data_train.replace("gluc", cardio_data_train.replace("gluc").apply(replace02))
cardio_data_train.replace("gender", cardio_data_train.replace("gender").apply(replace01))
cardio_data_train.replace("id", cardio_data_train.replace("id").apply(replace01))
```

i. Applying a different technique to handle missing values (For each technique verify your prediction results)


```
n_neighbors = 18, Accuracy: 0.716, F1-score = 0.684
n_neighbors = 20, Accuracy: 0.712, F1-score = 0.700
n_neighbors = 21, Accuracy: 0.724, F1-score = 0.693
n_neighbors = 22, Accuracy: 0.720, F1-score = 0.685
n_neighbors = 23, Accuracy: 0.716, F1-score = 0.684
n_neighbors = 24, Accuracy: 0.724, F1-score = 0.691
n_neighbors = 27, Accuracy: 0.712, F1-score = 0.691
n_neighbors = 28, Accuracy: 0.720, F1-score = 0.685
n_neighbors = 32, Accuracy: 0.720, F1-score = 0.688
n_neighbors = 33, Accuracy: 0.728, F1-score = 0.693
n_neighbors = 34, Accuracy: 0.724, F1-score = 0.693
n_neighbors = 35, Accuracy: 0.728, F1-score = 0.699
n_neighbors = 36, Accuracy: 0.724, F1-score = 0.696
n_neighbors = 37, Accuracy: 0.724, F1-score = 0.693
n_neighbors = 38, Accuracy: 0.720, F1-score = 0.691
n_neighbors = 39, Accuracy: 0.720, F1-score = 0.688
n_neighbors = 40, Accuracy: 0.724, F1-score = 0.688
n_neighbors = 41, Accuracy: 0.728, F1-score = 0.679
n_neighbors = 42, Accuracy: 0.732, F1-score = 0.693
n_neighbors = 43, Accuracy: 0.720, F1-score = 0.693
n_neighbors = 44, Accuracy: 0.716, F1-score = 0.688
n_neighbors = 45, Accuracy: 0.732, F1-score = 0.687
n_neighbors = 46, Accuracy: 0.728, F1-score = 0.691
n_neighbors = 47, Accuracy: 0.728, F1-score = 0.694
n_neighbors = 48, Accuracy: 0.732, F1-score = 0.697
n_neighbors = 49, Accuracy: 0.720, F1-score = 0.683
n_neighbors = 50, Accuracy: 0.724, F1-score = 0.691
n_neighbors = 51, Accuracy: 0.728, F1-score = 0.696
n_neighbors = 52, Accuracy: 0.732, F1-score = 0.699
n_neighbors = 53, Accuracy: 0.724, F1-score = 0.691
n_neighbors = 54, Accuracy: 0.720, F1-score = 0.705
n_neighbors = 55, Accuracy: 0.732, F1-score = 0.700
n_neighbors = 56, Accuracy: 0.736, F1-score = 0.705
n_neighbors = 57, Accuracy: 0.720, F1-score = 0.700
n_neighbors = 58, Accuracy: 0.728, F1-score = 0.696
n_neighbors = 59, Accuracy: 0.724, F1-score = 0.691
n_neighbors = 60, Accuracy: 0.728, F1-score = 0.711
n_neighbors = 61, Accuracy: 0.724, F1-score = 0.691
n_neighbors = 62, Accuracy: 0.720, F1-score = 0.697
n_neighbors = 63, Accuracy: 0.728, F1-score = 0.699
n_neighbors = 64, Accuracy: 0.728, F1-score = 0.696
n_neighbors = 65, Accuracy: 0.728, F1-score = 0.708
n_neighbors = 66, Accuracy: 0.720, F1-score = 0.688
n_neighbors = 67, Accuracy: 0.728, F1-score = 0.696
n_neighbors = 68, Accuracy: 0.720, F1-score = 0.711
n_neighbors = 69, Accuracy: 0.720, F1-score = 0.688
n_neighbors = 70, Accuracy: 0.728, F1-score = 0.700
n_neighbors = 71, Accuracy: 0.728, F1-score = 0.694
n_neighbors = 72, Accuracy: 0.728, F1-score = 0.694
n_neighbors = 73, Accuracy: 0.720, F1-score = 0.700
n_neighbors = 74, Accuracy: 0.736, F1-score = 0.705
n_neighbors = 75, Accuracy: 0.740, F1-score = 0.709
n_neighbors = 76, Accuracy: 0.732, F1-score = 0.700
n_neighbors = 77, Accuracy: 0.732, F1-score = 0.696
n_neighbors = 78, Accuracy: 0.740, F1-score = 0.728
n_neighbors = 79, Accuracy: 0.732, F1-score = 0.673
n_neighbors = 80, Accuracy: 0.720, F1-score = 0.679
n_neighbors = 81, Accuracy: 0.720, F1-score = 0.679
n_neighbors = 82, Accuracy: 0.720, F1-score = 0.679
n_neighbors = 83, Accuracy: 0.724, F1-score = 0.682
n_neighbors = 84, Accuracy: 0.720, F1-score = 0.679
n_neighbors = 85, Accuracy: 0.724, F1-score = 0.685
n_neighbors = 86, Accuracy: 0.724, F1-score = 0.688
n_neighbors = 87, Accuracy: 0.724, F1-score = 0.691
n_neighbors = 88, Accuracy: 0.724, F1-score = 0.685
n_neighbors = 89, Accuracy: 0.728, F1-score = 0.691
n_neighbors = 90, Accuracy: 0.724, F1-score = 0.688
n_neighbors = 91, Accuracy: 0.720, F1-score = 0.682
n_neighbors = 92, Accuracy: 0.724, F1-score = 0.688
n_neighbors = 93, Accuracy: 0.732, F1-score = 0.697
n_neighbors = 94, Accuracy: 0.724, F1-score = 0.688
n_neighbors = 95, Accuracy: 0.728, F1-score = 0.735
n_neighbors = 96, Accuracy: 0.732, F1-score = 0.697
n_neighbors = 97, Accuracy: 0.736, F1-score = 0.700
n_neighbors = 98, Accuracy: 0.728, F1-score = 0.688
n_neighbors = 99, Accuracy: 0.724, F1-score = 0.682
```

```
In [42]: model_knn = KNeighborsClassifier(n_neighbors=54, weights="distance")
model_knn.fit(x_train,y_train)
y_pred_test = model_knn.predict(x_test)
acc = accuracy_score(y_test, y_pred_test)
f1 = f1_score(y_test, y_pred_test)
print("KNN Accuracy = %.3f, F1-score = %.3f%" (acc, f1))

KNN Accuracy = 0.736, F1-score = 0.705
```

Bagging Classifier

```
In [43]: from sklearn.ensemble import BaggingClassifier
from numpy import arange

base_estimators = [10, 50, 100, 200, 400]
for n in base_estimators:
    for i in arange(0.1, 1.1, 0.1):
        model_bagging = BaggingClassifier(n_estimators = n, max_samples = 1)
        model_bagging.fit(x_train,y_train)
        y_pred_test = model_bagging.predict(x_test)
        acc = accuracy_score(y_test, y_pred_test)
        f1 = f1_score(y_test, y_pred_test)
        if acc > 0.72 and f1 > 0.72:
            print("n_estimators = %d, max_samples = %.2f, Accuracy: %.3f, F1-score = %.3f" % (n, i, acc, f1))

n_estimators = 10 max_samples = 0.40, Accuracy: 0.724, F1-score = 0.721
n_estimators = 30 max_samples = 0.20, Accuracy: 0.732, F1-score = 0.729
n_estimators = 50 max_samples = 0.30, Accuracy: 0.740, F1-score = 0.737
n_estimators = 80 max_samples = 0.80, Accuracy: 0.728, F1-score = 0.726
n_estimators = 100 max_samples = 1.00, Accuracy: 0.728, F1-score = 0.724
n_estimators = 100 max_samples = 0.10, Accuracy: 0.736, F1-score = 0.725
n_estimators = 100 max_samples = 0.20, Accuracy: 0.736, F1-score = 0.747
n_estimators = 100 max_samples = 0.30, Accuracy: 0.732, F1-score = 0.724
n_estimators = 100 max_samples = 0.40, Accuracy: 0.736, F1-score = 0.749
n_estimators = 100 max_samples = 1.00, Accuracy: 0.732, F1-score = 0.737
n_estimators = 200 max_samples = 0.10, Accuracy: 0.744, F1-score = 0.728
n_estimators = 200 max_samples = 0.20, Accuracy: 0.732, F1-score = 0.744
n_estimators = 200 max_samples = 0.30, Accuracy: 0.732, F1-score = 0.727
n_estimators = 200 max_samples = 0.40, Accuracy: 0.732, F1-score = 0.729
n_estimators = 200 max_samples = 0.60, Accuracy: 0.732, F1-score = 0.722
n_estimators = 200 max_samples = 0.80, Accuracy: 0.728, F1-score = 0.728
n_estimators = 200 max_samples = 0.90, Accuracy: 0.740, F1-score = 0.735
n_estimators = 400 max_samples = 0.10, Accuracy: 0.744, F1-score = 0.746
n_estimators = 400 max_samples = 0.30, Accuracy: 0.740, F1-score = 0.737
n_estimators = 400 max_samples = 0.40, Accuracy: 0.732, F1-score = 0.734
n_estimators = 400 max_samples = 0.70, Accuracy: 0.736, F1-score = 0.732
n_estimators = 400 max_samples = 0.80, Accuracy: 0.748, F1-score = 0.739
```

```
In [45]: model_bagging = BaggingClassifier(n_estimators = 100, max_samples = 0.2)
model_bagging.fit(x_train,y_train)
y_pred_test = model_bagging.predict(x_test)
acc = accuracy_score(y_test, y_pred_test)
f1 = f1_score(y_test, y_pred_test)
print("Bagging Accuracy = %.3f, F1-score = %.3f%" (acc, f1))
# kaggle score: 0.680 (bagging(50, 0.5))

Bagging Accuracy = 0.752, F1-score = 0.742
```

Voting Classifier

```
In [46]: from sklearn.ensemble import VotingClassifier
model_voting = VotingClassifier('rf', model_rf1, ('knn', model_knn), ('svc', model_svm), voting='hard')
model_voting.fit(x_train,y_train)
y_pred_test = model_voting.predict(x_test)
acc = accuracy_score(y_test, y_pred_test)
f1 = f1_score(y_test, y_pred_test)
print("Voting Classifier Accuracy = %.3f, F1-score = %.3f%" (acc, f1))

Voting Classifier Accuracy = 0.756, F1-score = 0.734
```

Evaluate models

```
In [47]: test_data = pd.read_csv("cardio-test.csv", sep=";")
print(test_data.shape)
test_data.head()
```

o. Use the cardio-validation.csv and cardio-train.csv as well to make your final prediction.

```

: # REPLACE CLF
: final_model = model_rf
:
: train the model on training + validation data
: final_model.fit(x_train.append(x_test), y_train.append(y_test))
:
: generate submission file
: y_pred_submission = final_model.predict(test_data)
: df_result = pd.DataFrame(y_pred_submission, columns=["cardio"])
: sample = pd.read_csv("sample-submission.csv")
: df_result_final = pd.concat([sample.loc[:, "id"], df_result], axis=1)
: df_result_final.to_csv("submissions/submit-final-1f-4.csv", index=False)
:
:
: df_result_final
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:
:

```

```
In [48]: test_data.isna().sum()
```

```
Out[48]: id          0
age          0
gender       0
height       0
weight       0
ap_hi        0
ap_lo        0
cholesterol  0
gluc         0
smoke        0
alco         0
active       0
dtype: int64
```

```
In [49]: from sklearn.preprocessing import StandardScaler, MinMaxScaler
test_data = preprocess(test_data)
scaler = StandardScaler()
test_data = pd.DataFrame(scaler.fit_transform(test_data))
```

```
In [50]: test_data.head()
```

		0	1	2	3	4	5	6	7	8	9	10
Out[50]:	0	0.217526	0.743502	-1.047999	-1.126201	-0.967672	-0.107793	-0.485965	-0.352689	-0.318311	3.702146	-1.760216
	1	-0.533600	0.743502	0.552257	-0.917645	-0.373133	-0.088149	-0.485965	-0.352689	-0.318311	-0.270114	0.568112
	2	0.387751	-1.344487	0.687099	0.477278	0.815946	-0.109758	-0.485965	-0.352689	-0.318311	-0.270114	0.568112
	3	1.569553	0.743502	-0.801028	0.477278	2.005024	-0.068504	2.551318	-0.352689	-0.318311	-0.270114	0.568112
	4	0.387751	0.743502	-0.801028	1.306499	0.815946	-0.109758	-0.485965	1.651228	-0.318311	-0.270114	-1.760216

n. Upload your test data predictions to Kaggle competition in the correct submission format.

o. Use the cardio-validation.csv and cardio-train.csv as well to make your final prediction.

```
In [52]: # REF:KAGGLE CUP
final_model = model_rf

#train the model on training + validation data
final_model.fit(x_train.append(x_test), y_train.append(y_test))

# generate submission file
df_result = pd.DataFrame(y_pred_submission,columns=["cardio"])
sample = pd.read_csv("sample_submission.csv")
df_result_final = pd.concat([sample.loc[:,["id"]],df_result], axis=1)
df_result_final.to_csv("submissions/submit-final-rf-4.csv", index=False)
```

```
In [53]: df_result_final
```

```
11 active      1000 non-null int64
12 cardio      1000 non-null int64
dtypes: float64(1), int64(9), object(3)
memory usage: 101.7+ KB
```

```
cardio_complete_data.describe()
```

	id	age	height	weight	ap_hi	ap_lo	smoke	alco	active
--	----	-----	--------	--------	-------	-------	-------	------	--------

250 rows x 2 columns

Task 2

```
In [101]: from sklearn.metrics import accuracy_score, f1_score
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression

cardio_complete_data = pd.read_csv("cardio-complete .csv")
print(cardio_complete_data.shape)
cardio_complete_data.head()
```

	id	age	gender	height	weight	ap_hi	ap_lo	cholesterol	gluc	smoke	alco	active	cardio
Out[101]:	0	66667	20252	Women	161	62.0	140	90	High	High	0	0	1
	1	22956	21129	Men	166	66.0	125	70	Normal	Normal	1	0	1
	2	40536	16602	Men	160	74.0	140	90	Normal	Normal	0	0	1
	3	39712	15172	Men	167	77.0	120	80	Normal	Normal	0	0	1
	4	82165	19858	Women	176	93.0	140	90	Above Normal	Normal	0	1	1

```
In [102]: #checking presence of Null values
cardio_complete_data.isna().sum()
```

```
Out[102]: id          0
age          0
gender       0
height       0
weight       0
ap_hi        0
ap_lo        0
cholesterol  0
gluc         0
smoke        0
alco         0
active       0
cardio       0
dtype: int64
```

```
In [103]: cardio_complete_data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 13 columns):
 #   Column        Non-Null Count  Dtype  
---  -
 0   id            1000 non-null    int64  
 1   age           1000 non-null    int64  
 2   gender        1000 non-null    object  
 3   height        1000 non-null    int64  
 4   weight        1000 non-null    float64 
 5   ap_hi         1000 non-null    int64  
 6   ap_lo         1000 non-null    int64  
 7   cholesterol    1000 non-null    object  
 8   gluc          1000 non-null    object  
 9   smoke         1000 non-null    int64  
10  alco          1000 non-null    int64  
11  active         1000 non-null    int64  
12  cardio        1000 non-null    int64  
dtypes: float64(1), int64(9), object(3)
memory usage: 101.7+ KB
```

	id	age	height	weight	ap_hi	ap_lo	smoke	alco	active
Out[104]:	count	1000.000000	1000.000000	1000.000000	1000.000000	1000.000000	1000.000000	1000.000000	1000.000000
	mean	49937.141000	149430.11200	646.288000	74.256200	126.221000	88.654000	0.080000	0.045000
	std	28845.044455	29491.98155	8.415811	14.141323	19.316969	83.606641	0.271429	0.207048
	min	135.000000	14344.00000	109.000000	36.000000	70.000000	40.000000	0.000000	0.000000
	25%	25325.500000	17822.50000	159.000000	64.000000	120.000000	80.000000	0.000000	0.000000
	50%	4957.000000	19716.00000	164.000000	72.000000	120.000000	80.000000	0.000000	0.000000
	75%	74524.000000	21311.75000	170.000000	83.000000	140.000000	90.000000	0.000000	0.000000
	max	99699.000000	23655.00000	194.000000	140.000000	215.000000	1100.000000	1.000000	1.000000

Observations:

- The cardio complete data has 0 null values. Hence, we will require less preprocessing for this data.
- Complete data has 1000 rows while in task 1, training data had 500 rows and validation data had 500 rows. So in this part, we have more data to fit the model.

```
In [105]: # same preprocessing as task 1 (except missing data imputation is not needed)
cardio_complete_data = preprocess(cardio_complete_data)
cardio_complete_data.head()
```

	age	gender	height	weight	ap_hi	ap_lo	cholesterol	gluc	smoke	alco	active	cardio
Out[105]:	0	65.0	1	161	62.0	140	90	1	1	0	0	1
	1	58.0	0	166	66.0	125	70	2	2	1	0	1
	2	45.0	0	160	74.0	140	90	2	2	0	0	1
	3	42.0	0	167	77.0	120	80	2	2	0	0	1
	4	54.0	1	176	93.0	140	90	0	2	0	1	1

```
In [106]: y = cardio_complete_data.loc[:, "cardio"]
X = cardio_complete_data.drop("cardio", 1)
```

```
In [107]: from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test = train_test_split(X,y,stratify=y,test_size=0.2,random_state=1)
y_train_task2 = y_train.copy()
```

```
In [108]: c_values = [100, 10, 1.0, 0.1, 0.01]
for c in c_values:
    model_log_reg2 = LogisticRegression(C = c, max_iter = 5000)
    model_log_reg2.fit(x_train,y_train)
    y_pred_test = model_log_reg2.predict(x_test)
    acc = accuracy_score(y_test, y_pred_test)
    f1 = f1_score(y_test, y_pred_test)
    print("C = %.2f, Accuracy = %.3f, F1-score = %.3f%" (c, acc, f1))

C = 100.00, Accuracy = 0.680, F1-score = 0.680
C = 10.00, Accuracy = 0.680, F1-score = 0.680
C = 1.00, Accuracy = 0.685, F1-score = 0.687
C = 0.10, Accuracy = 0.695, F1-score = 0.700
C = 0.01, Accuracy = 0.705, F1-score = 0.704
```

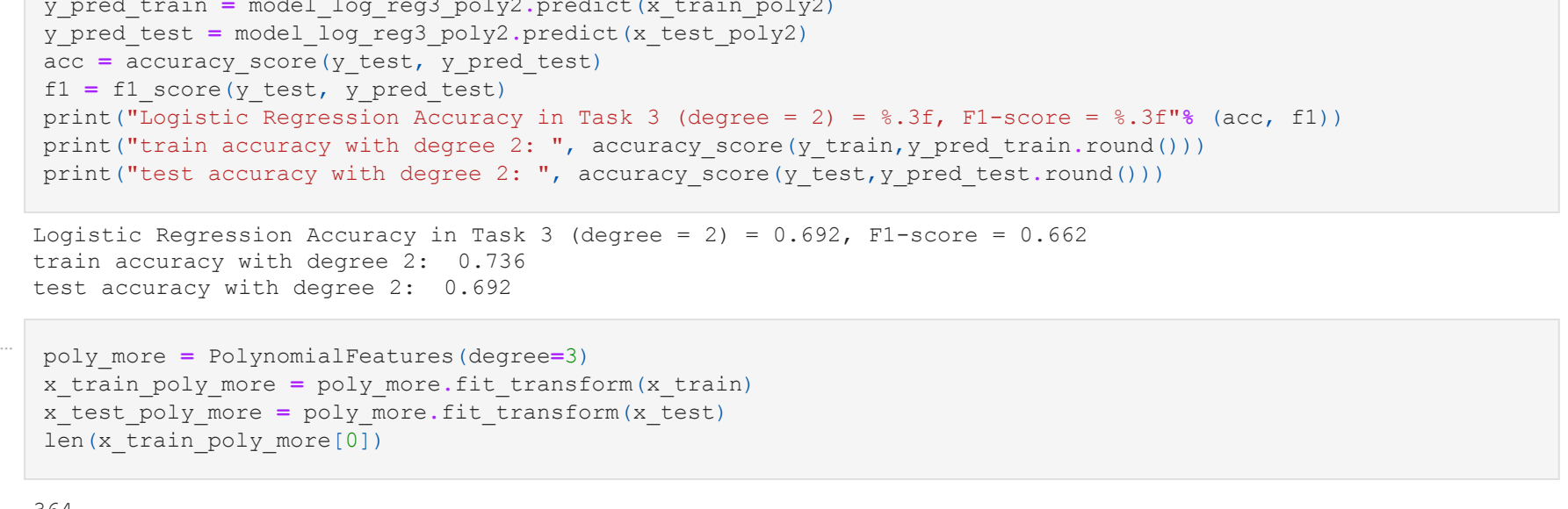
```
In [109]: model_log_reg2 = LogisticRegression(C = 0.01, max_iter = 5000)
model_log_reg2.fit(x_train,y_train)
y_pred_test = model_log_reg2.predict(x_test)
acc = accuracy_score(y_test, y_pred_test)
f1 = f1_score(y_test, y_pred_test)
print("Logistic Regression Accuracy in Task 2 = %.3f, F1-score = %.3f%" (acc, f1))
y_pred_task2 = y_pred_test.copy()
```

Logistic Regression Accuracy in Task 2 = 0.705, F1-score = 0.704

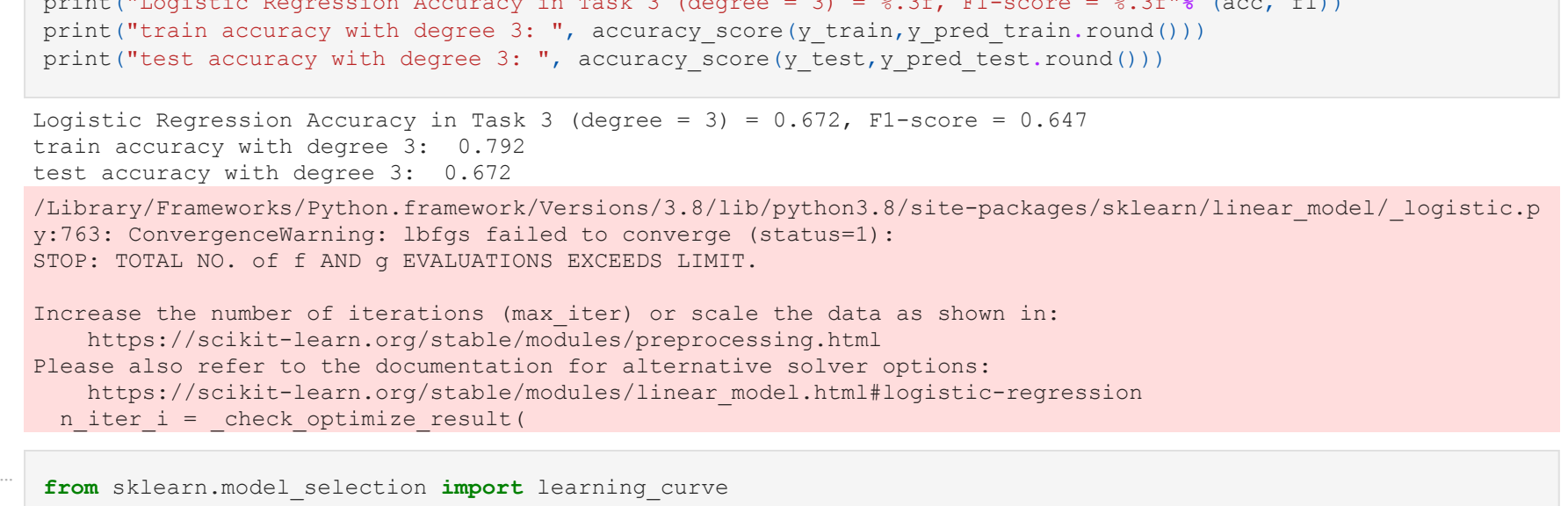
```
In [110]: from sklearn.metrics import precision_score, recall_score, classification_report, f1_score, confusion_matrix, roc_auc_score
import seaborn as sns
import matplotlib.pyplot as plt

print("Model = Logistic Regression")
print("For Task 1: Accuracy = 0.716, Precision = 0.720, Recall = 0.691, F1-score = 0.705")
print("For Task 2: Accuracy = 0.705, Precision = 0.722, Recall = 0.686, F1-score = 0.704")
print("For Task 3: Accuracy = 0.672, Precision = 0.672, Recall = 0.672, F1-score = 0.672")
```

```
In [111]: conf_matrix_task1 = confusion_matrix(y_test_task1, y_pred_task1)
plt.title("Confusion Matrix for Task 1")
sns.heatmap(conf_matrix_task1, annot=True, fmt="")
plt.show()
```



```
In [112]: conf_matrix_task2 = confusion_matrix(y_test_task2, y_pred_task2)
plt.title("Confusion Matrix for Task 2")
sns.heatmap(conf_matrix_task2, annot=True, fmt="")
plt.show()
```



Task 3

```
In [113]: # Repeating same data as Task 1
from sklearn.metrics import accuracy_score, f1_score
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression

train_data = pd.read_csv("cardio-train.csv", sep=";")
train_data = preprocess(train_data)
validation_data = preprocess(validation_data)
y_train = train_data.loc[:, "cardio"]
x_train = train_data.drop("cardio", 1)
y_test = validation_data.loc[:, "cardio"]
x_test = validation_data.drop("cardio", 1)
```

```
Out[114]: 78
```

```
In [115]: model_log_reg3_poly2 = LogisticRegression(max_iter = 50000)
model_log_reg3_poly2.fit(x_train_poly2,y_train)
y_pred_train = model_log_reg3_poly2.predict(x_train_poly2)
y_pred_test = model_log_reg3_poly2.predict(x_test_poly2)
acc = accuracy_score(y_test, y_pred_test)
f1 = f1_score(y_test, y_pred_test)
print("Logistic Regression Accuracy in Task 3 (degree = 2) = %.3f, F1-score = %.3f%" (acc, f1))
print("train accuracy with degree 2: ", accuracy_score(y_train,y_pred_train.round()))
print("test accuracy with degree 2: ", accuracy_score(y_test,y_pred_test.round()))

Logistic Regression Accuracy in Task 3 (degree = 2) = 0.692, F1-score = 0.662
train accuracy with degree 2: 0.736
test accuracy with degree 2: 0.692
```

```
In [116]: poly_more = PolynomialFeatures(degree=3)
x_train_poly_more = poly_more.fit_transform(x_train)
x_test_poly_more = poly_more.fit_transform(x_test)
len(x_train_poly_more[0])
```

```
Out[116]: 364
```

```
In [117]: model_log_reg3_poly_more = LogisticRegression(max_iter = 50000)
model_log_reg3_poly_more.fit(x_train_poly_more,y_train)
y_pred_train = model_log_reg3_poly_more.predict(x_train_poly_more)
y_pred_test = model_log_reg3_poly_more.predict(x_test_poly_more)
acc = accuracy_score(y_test, y_pred_test)
f1 = f1_score(y_test, y_pred_test)
print("Logistic Regression Accuracy in Task 3 (degree = 3) = %.3f, F1-score = %.3f%" (acc, f1))
print("train accuracy with degree 3: ", accuracy_score(y_train,y_pred_train.round()))
print("test accuracy with degree 3: ", accuracy_score(y_test,y_pred_test.round()))

Logistic Regression Accuracy in Task 3 (degree = 3) = 0.672, F1-score = 0.647
train accuracy with degree 3: 0.752
test accuracy with degree 3: 0.672
```

```
Library/Frameworks/Python.framework/Versions/3.8/lib/python3.8/site-packages/sklearn/linear_model/_logistic.py:763: ConvergenceWarning: lbfgs failed to converge (status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
n_iter_ = 10000
```

```
In [119]: from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, f1_score, confusion_matrix, roc_auc_score
def plot_learning_curve(estimator, title, X, y, axes_names, n_jobs=None, cv=None,
                        n_jobs=None, train_size=np.linspace(0.1, 1.0, 5)):
    if axes is None:
        _, axes = plt.subplots(1, 1, figsize=(10, 5))
    axes.set_title(title)
    axes.set_xlabel(axes_names[0])
    axes.set_ylabel("Score")
    train_sizes, train_scores, test_scores, fit_times, _ = \
        learning_curve(estimator, X, y, cv=cv, n_jobs=n_jobs,
                        return_times=True)
    train_scores_mean = np.mean(train_scores, axis=1)
    train_scores_std = np.std(train_scores, axis=1)
    test_scores_mean = np.mean(test_scores, axis=1)
    test_scores_std = np.std(test_scores, axis=1)
    fit_times_mean = np.mean(fit_times, axis=1)
    fit_times_std = np.std(fit_times, axis=1)
    # Plot learning curve
    axes.grid()
    axes
```