

Correlation between Risk Factors and Coronary Heart Disease (CHD)

-Group N3-

Au Yat Sin Candice 20860834
Chan Chung Yin 20978198
Chui Yuen Tsun 20864165

Introduction

01

~37,000 deaths
were related to coronary
heart diseases in 2021

02

3rd most common
cause of death in Hong
Kong

03

Risk factors
like smoking, drinking,
unhealthy diet, physical
inactivity

04

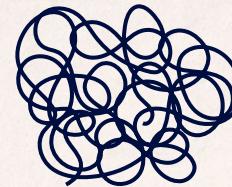
Unpredictable
Complications
even experienced
specialists may not always
foresee the development of
these complications

Questions we would like to answer

- (1) How different risk factors affect the prevalence rate of CHD?
- (2) Which factors are positively/negatively correlated with heart disease or attack?
- (3) Do other common complications share any similarities with CHD in the correlations between different risk factors?



Data collection and processing



Data collection

1. Dataset derived from a survey conducted from the Behavioural Risk Factor Surveillance System:

- Respondents were asked about:
 - Health status & living habit
 - Have CHD or not



Data cleaning

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 253680 entries, 0 to 253679
Data columns (total 20 columns):
 #   Column          Non-Null Count  Dtype  
--- 
 0   HeartDiseaseorAttack  253680 non-null   float64
 1   HighBP              253680 non-null   float64
 2   HighChol             253680 non-null   float64
 3   BMI                 253680 non-null   float64
 4   Smoker               253680 non-null   float64
 5   Stroke               253680 non-null   float64
 6   Diabetes              253680 non-null   float64
 7   PhysActivity          253680 non-null   float64
 8   Fruits                253680 non-null   float64
 9   Veggies               253680 non-null   float64
 10  HvyAlcoholConsump    253680 non-null   float64
 11  AnyHealthcare         253680 non-null   float64
 12  GenHlth               253680 non-null   float64
 13  MentHlth              253680 non-null   float64
 14  PhysHlth              253680 non-null   float64
 15  DiffWalk              253680 non-null   float64
 16  Sex                  253680 non-null   float64
 17  Age                  253680 non-null   float64
 18  Education             253680 non-null   float64
 19  Income                253680 non-null   float64
dtypes: float64(20)
memory usage: 38.7 MB
```

Data processing

```
HeartDiseaseorAttack  HighBP  HighChol  CholCheck  BMI  Smoker  Stroke \
0                   0.0     1.0      1.0       1.0  40.0    1.0     0.0
1                   0.0     0.0      0.0       0.0  25.0    1.0     0.0
2                   0.0     1.0      1.0       1.0  28.0    0.0     0.0
3                   0.0     1.0      0.0       1.0  27.0    0.0     0.0
4                   0.0     1.0      1.0       1.0  24.0    0.0     0.0

Diabetes  PhysActivity  Fruits  Veggies  HvyAlcoholConsump  AnyHealthcare \
0        0.0           0.0     0.0      1.0                  0.0            1.0
1        0.0           1.0     0.0      0.0                  0.0            0.0
2        0.0           0.0     1.0      0.0                  0.0            1.0
3        0.0           1.0     1.0      1.0                  0.0            1.0
4        0.0           1.0     1.0      1.0                  0.0            1.0

NoDocbcCost  GenHlth  MentHlth  PhysHlth  DiffWalk  Sex  Age  Education \
0        0.0      5.0     18.0     15.0      1.0  0.0  9.0     4.0
1        1.0      3.0     0.0      0.0      0.0  0.0  7.0     6.0
2        1.0      5.0     30.0     30.0      1.0  0.0  9.0     4.0
3        0.0      2.0     0.0      0.0      0.0  0.0 11.0     3.0
4        0.0      2.0     3.0      0.0      0.0  0.0 11.0     5.0

Income  HeartDiseaseorAttack  HighBP  HighChol  BMI  Smoker  Stroke  Diabetes \
0        3.0                  0.0     1.0      1.0  40.0    1.0     0.0     0.0
1        1.0                  0.0     0.0      0.0  25.0    1.0     0.0     0.0
2        8.0                  0.0     1.0      1.0  28.0    0.0     0.0     0.0
3        6.0                  0.0     1.0      0.0  27.0    0.0     0.0     0.0
4        4.0                  0.0     1.0      1.0  24.0    0.0     0.0     0.0

PhysActivity  Fruits  Veggies  HvyAlcoholConsump  AnyHealthcare  GenHlth \
0           0.0     0.0      1.0                  0.0            1.0     5.0
1           1.0     0.0      0.0                  0.0            0.0     3.0
2           0.0     1.0      0.0                  0.0            1.0     5.0
3           1.0     1.0      1.0                  0.0            1.0     2.0
4           1.0     1.0      1.0                  0.0            1.0     2.0

MentHlth  PhysHlth  DiffWalk  Sex  Age  Education  Income
0      -18.0    -15.0      1.0  0.0  9.0     4.0    3.0
1       -0.0     -0.0      0.0  0.0  7.0     6.0    1.0
2      -30.0    -30.0      1.0  0.0  9.0     4.0    8.0
3       -0.0     -0.0      0.0  0.0 11.0     3.0    6.0
4      -3.0     -0.0      0.0  0.0 11.0     5.0    4.0
```

Methodology

2. Multivariate dataset from the UCI Machine Learning Repository
 - survey where all respondents have a certain degree of chronic heart disease and related symptoms
 - list of complications, outcomes and any other disease that may result from coronary heart diseases

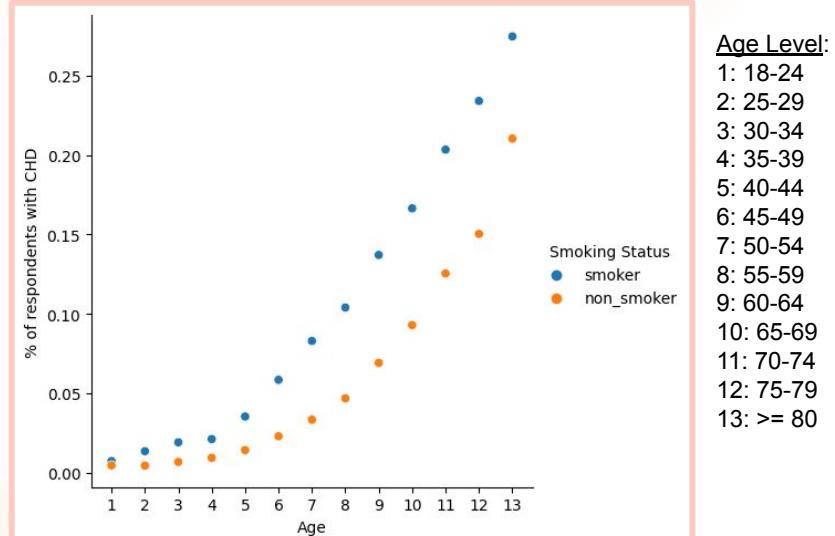
Data analysis and visualization

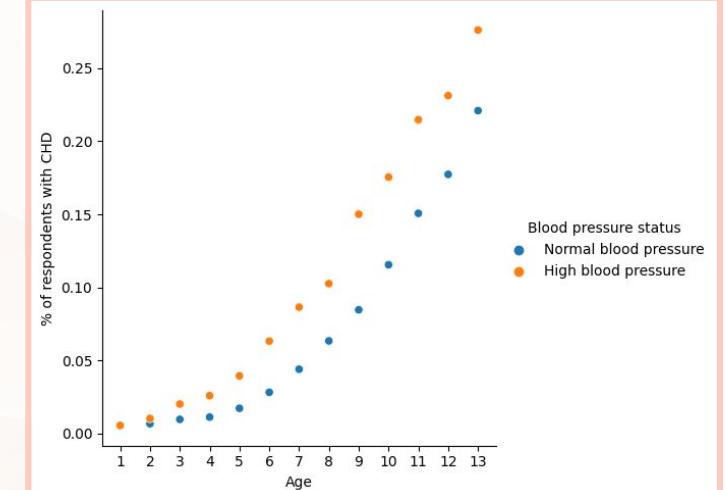
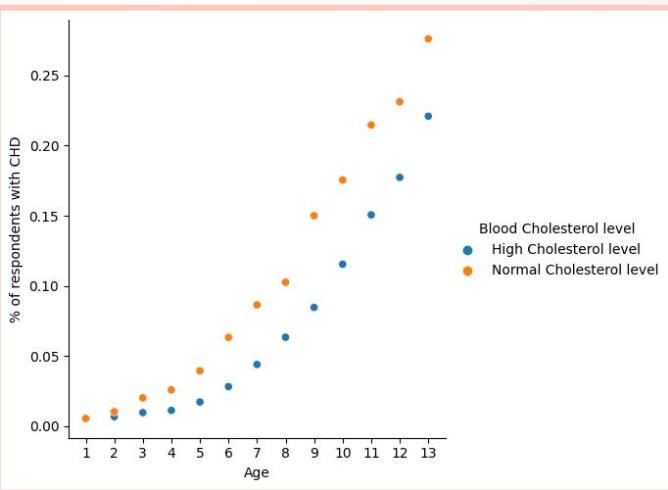
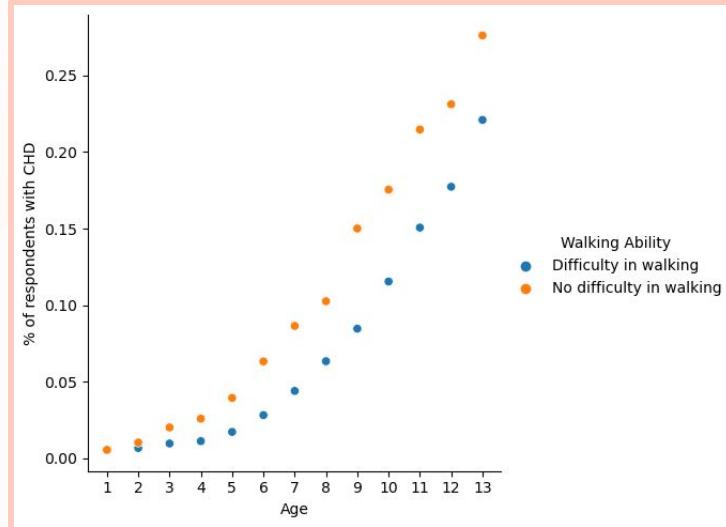
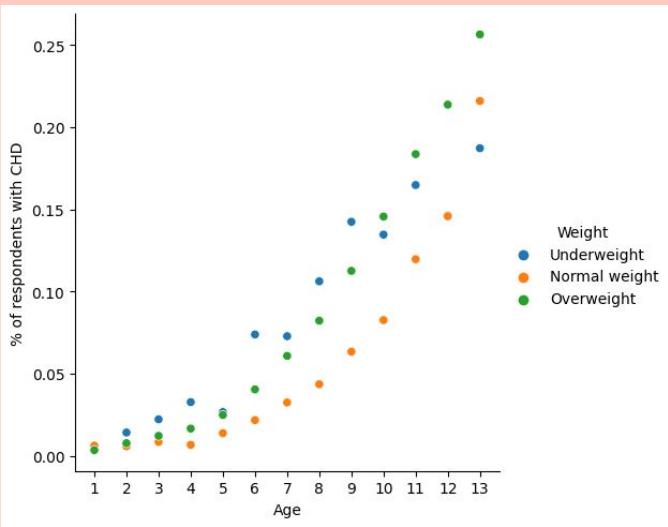


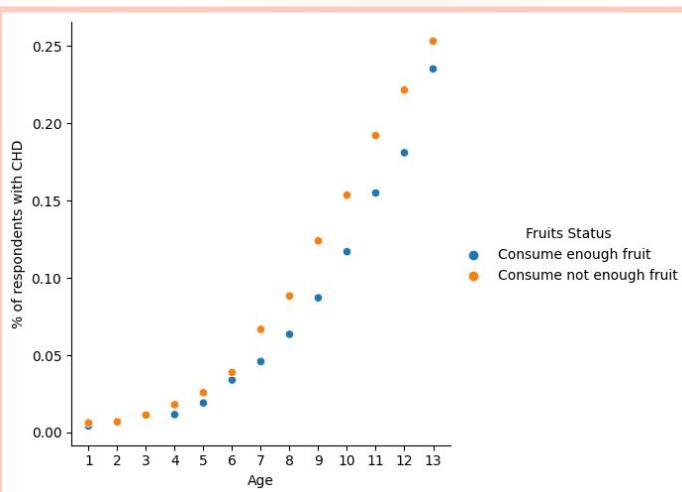
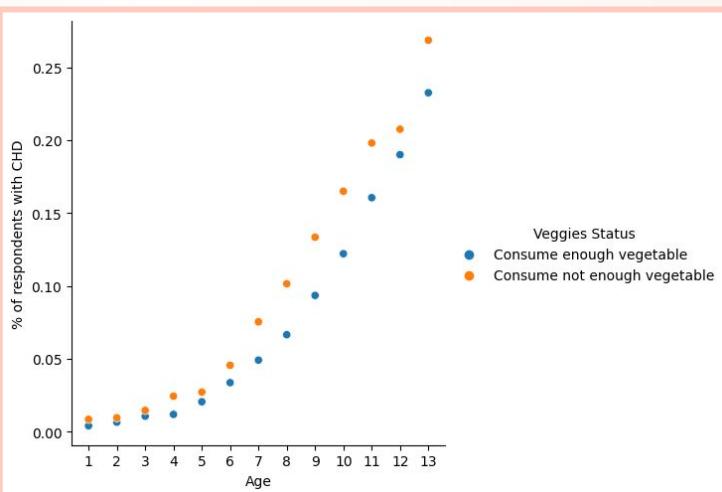
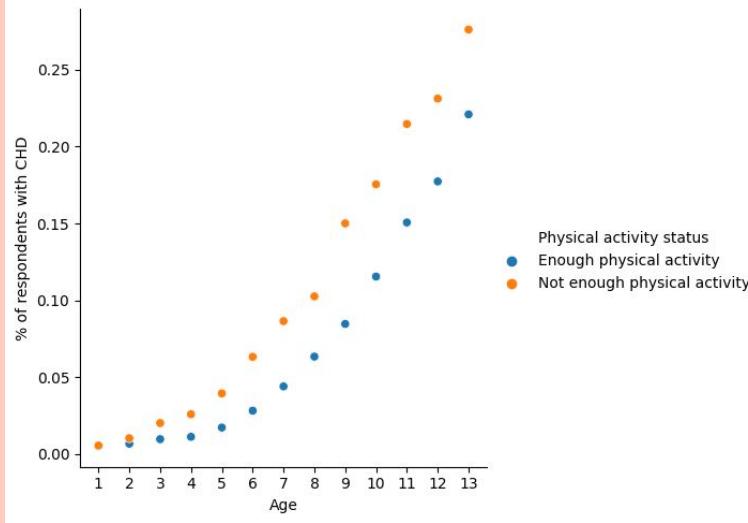
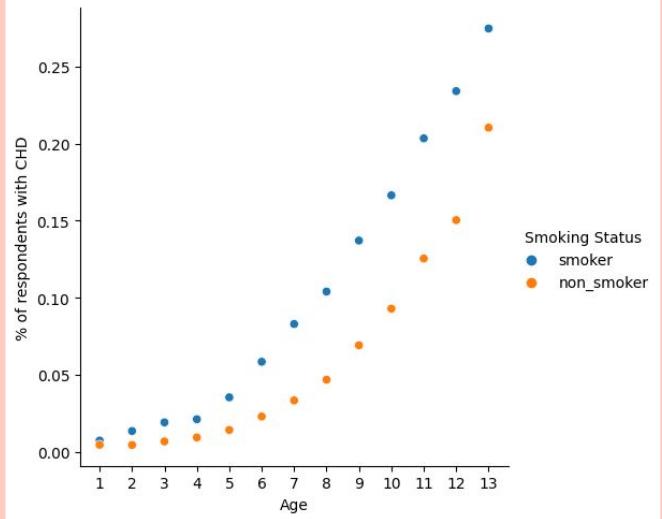
Methodology

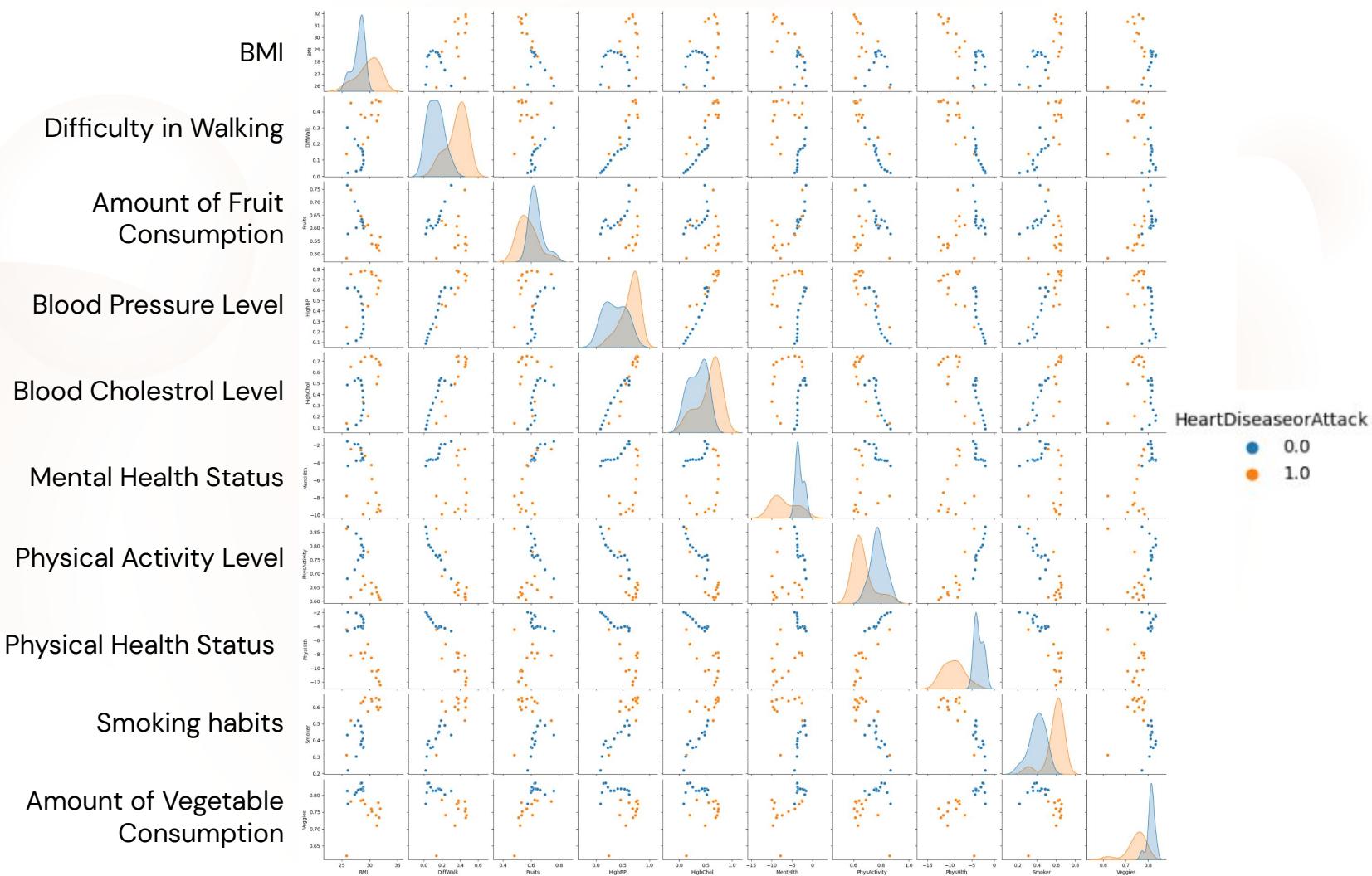
- Plot graphs of prevalence rate of CHD against different age levels
- Each graphs show variation in one risk factor
- Risk factors:
 - Smoking habit
 - Blood pressure level
 - Blood cholesterol level
 - Body Mass Index (BMI)
 - Walking difficulty
 - Physical activity
 - Fruit and Vegetable consumption

E.g. Variation in Smoking habits





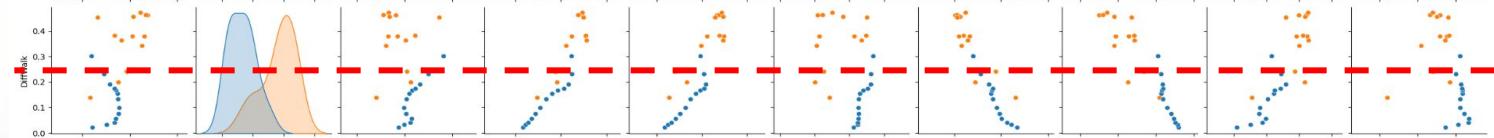






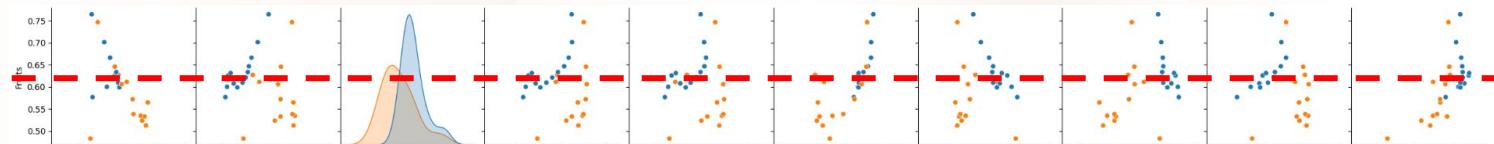
More obvious correlation with CHD

Difficulty in Walking

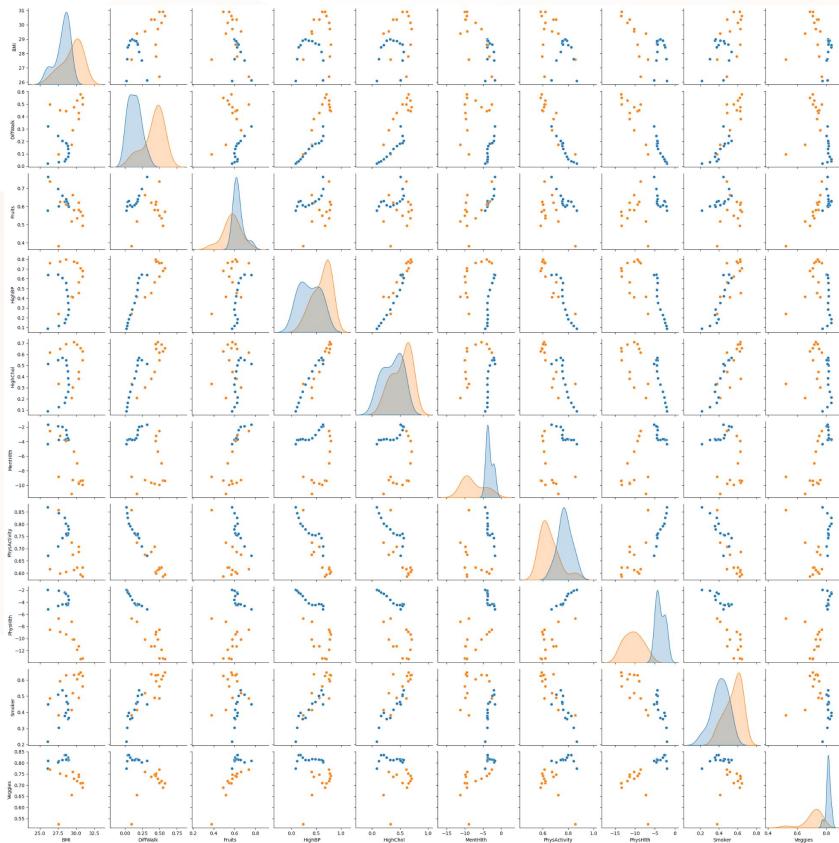


Less obvious correlation with CHD

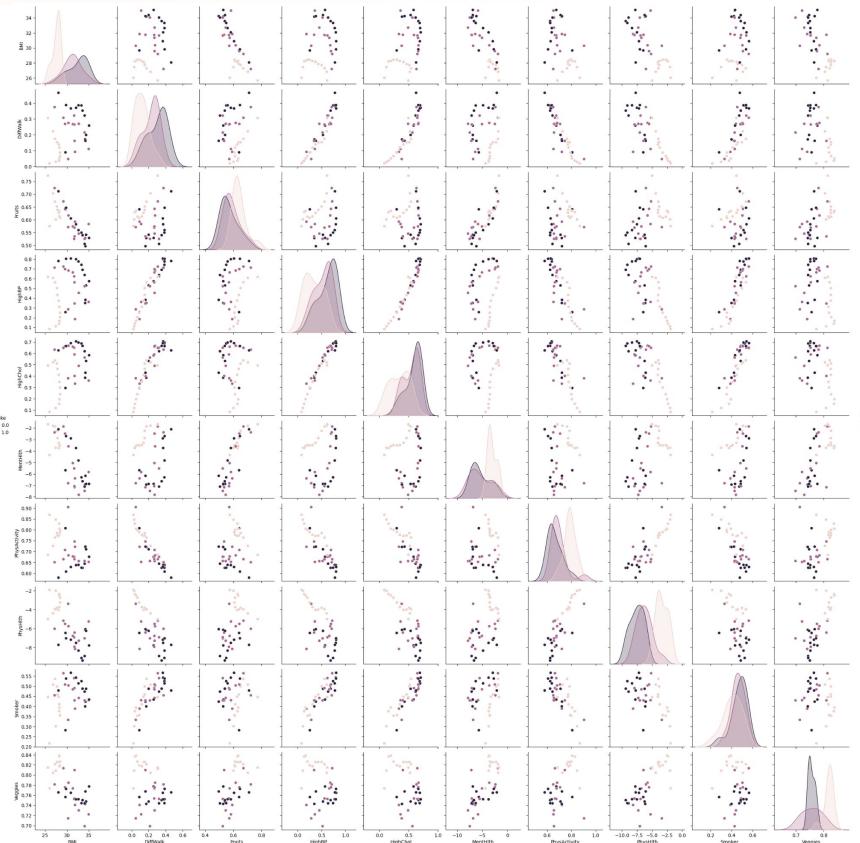
Amount of Fruit Consumption



Stroke

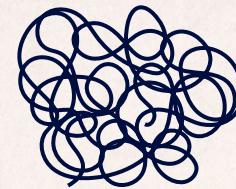


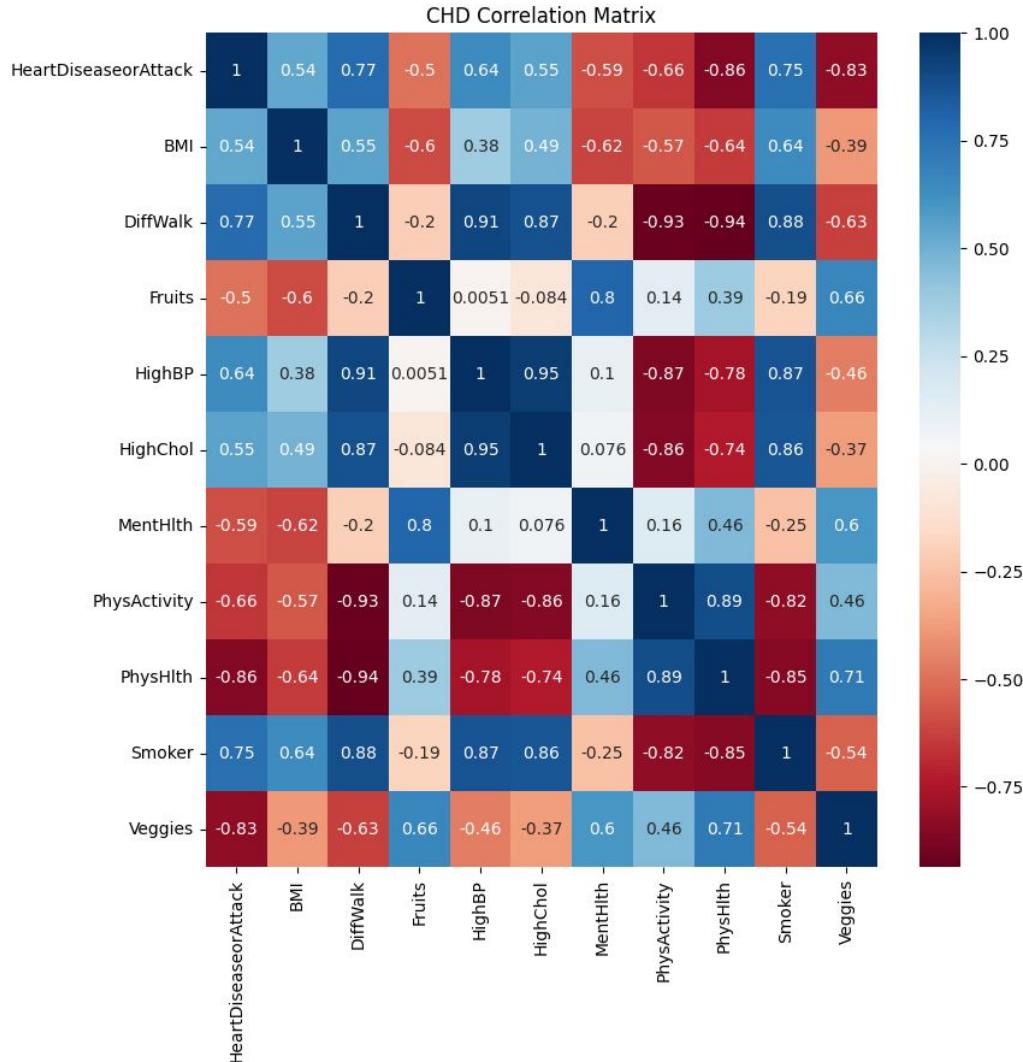
Diabetes



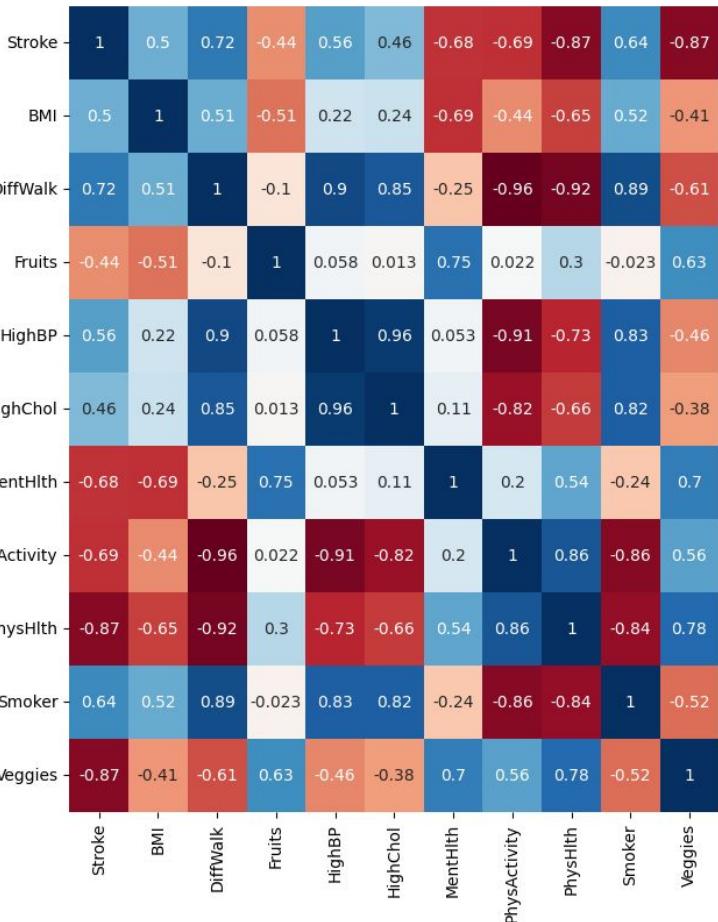


Correlation Analysis

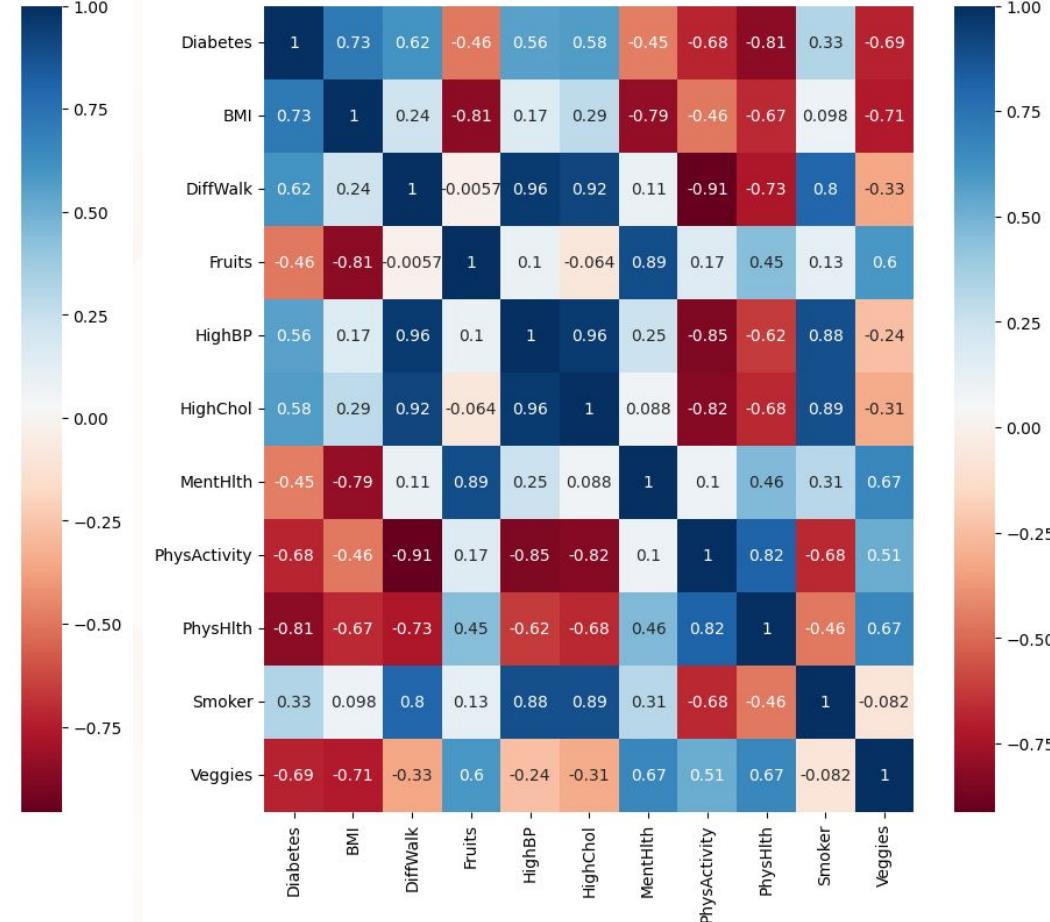




Stroke Correlation Matrix



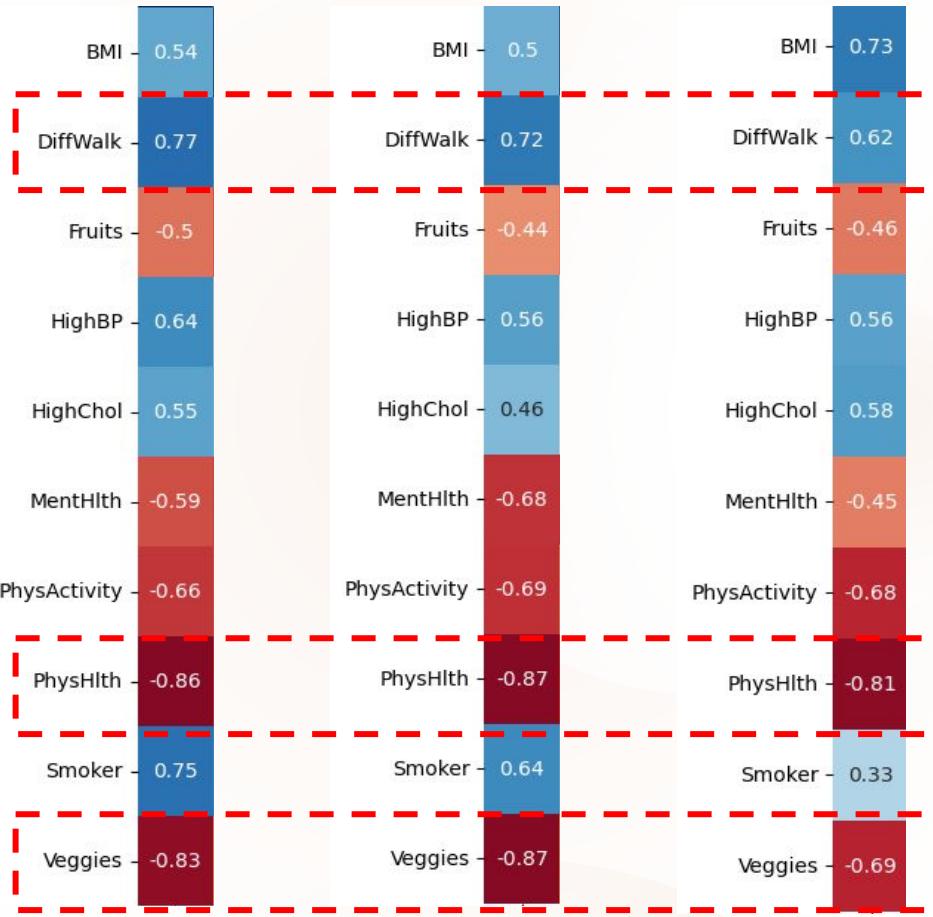
Diabetes Correlation Matrix



CHD

Stroke

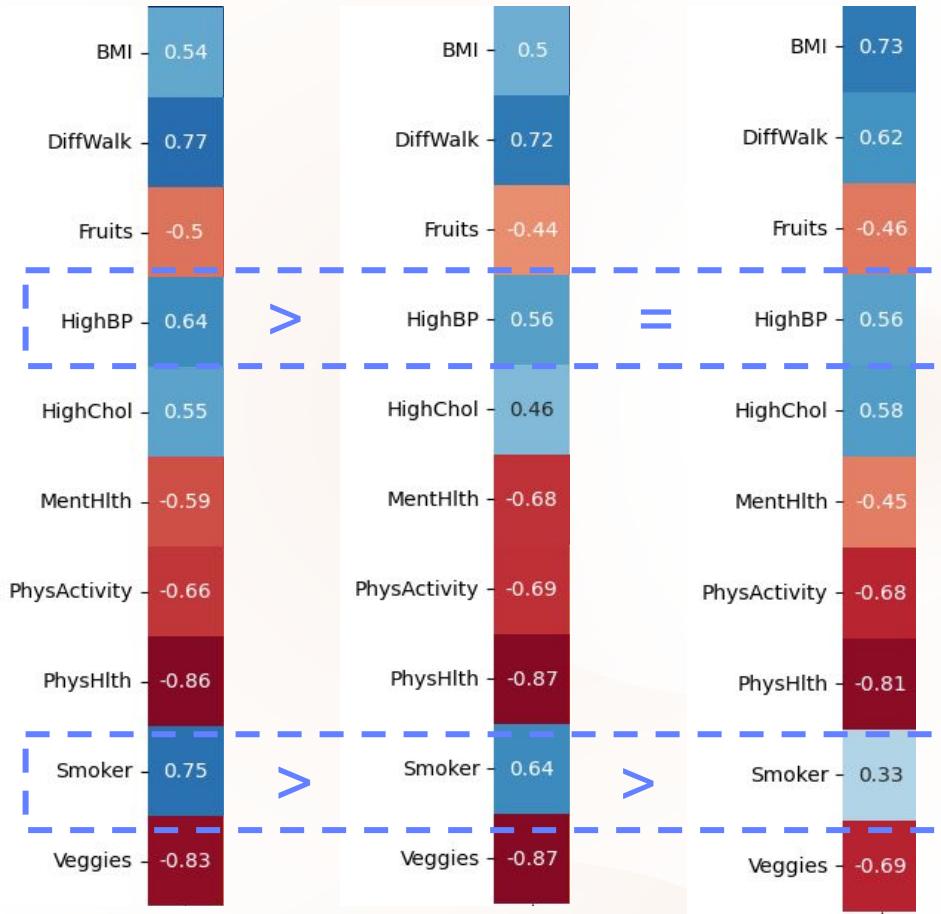
Diabetes



CHD

Stroke

Diabetes



Heart Disease Accuracy Prediction

```
log_cols=["Classifier", "Accuracy", "Log Loss"]
log = pd.DataFrame(columns=log_cols)

for cla in algorithm:
    cla.fit(X_train, Y_train)
    name = cla.__class__.__name__
    print("=" * 30)
    print(name)

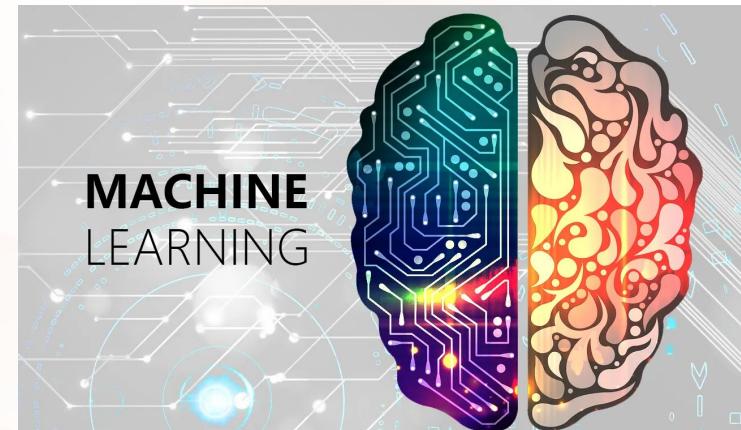
    train_predictions = cla.predict(X_test)
    acc = accuracy_score(Y_test, train_predictions)
    print("Accuracy: {:.4%}".format(acc))

    train_predictions = cla.predict(X_test)
    ll = log_loss(Y_test, train_predictions)
    print("Log Loss: {}".format(ll))

    log_entry = pd.DataFrame([[name, acc * 100, ll]], columns = log_cols)
    log = log.append(log_entry)
```

```
=====
LogisticRegression
****Results****
Accuracy: 90.6398%
Log Loss: 3.373764387119941
=====
KNeighborsClassifier
****Results****
Accuracy: 89.5676%
Log Loss: 3.760230554016814
=====
RandomForestClassifier
****Results****
Accuracy: 90.2318%
Log Loss: 3.5208204469501854
=====
XGBClassifier
****Results****
Accuracy: 90.6871%
Log Loss: 3.356714409168609
=====
AdaBoostClassifier
****Results****
Accuracy: 90.6851%
Log Loss: 3.3574248249165812
=====
GradientBoostingClassifier
****Results****
Accuracy: 90.7640%
Log Loss: 3.329008194997694
```

Conclusion





Reference List



Teboul, A. (2022). Heart Disease Health Indicators Dataset. Kaggle. Retrieved 30 April, 2023, from
<https://www.kaggle.com/datasets/alexteboul/heart-disease-health-indicators-dataset>



Pixabey - Wittyspark (n.d.). Retrieved 30 April, 2023, from
<https://towardsdatascience.com/introduction-to-machine-learning-for-beginners-eed6024fdb08>



Member contributions

- Au Yat Sin Candice: Desktop research, data cleaning and processing
- Chan Chung Yin: Desktop research, data visualization (Correlation Matrix) and Machine Learning Prediction
- Chui Yuen Tsun: Desktop research, data visualization (Scatter plot & Pairplot)

Thank you~