

hear_attack_analysis

February 12, 2025

```
[ ]: # import library
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import statsmodels.api as sm
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import accuracy_score, confusion_matrix, \
    classification_report
```

```
[2]: # import the data
heart_data = pd.read_csv("/home/student/Documents/AIMS/COOP_tasks/
    data_science_projects/Heart_Attack_Risk_Assessment/data/updated_version.csv")
```

0.1 Description of our dataset:

This dataset contains 1,000 patient records generated for health risk assessment. It includes biometric health indicators commonly used in cardiovascular and general health research. Each record captures age, cholesterol levels, blood pressure, smoking habits, diabetes status, and heart attack history—key factors influencing cardiovascular diseases.

On this dataset we are intempt to run: * Exploratory Data Analysis (EDA) * Statistical analysis
* Machine learning classification tasks

These are the columns: * age: Patient's age (years) * sex: Biological sex (0 = Female, 1 = Male) * total_cholesterol: Total cholesterol level (mg/dL) * ldl: Low-Density Lipoprotein (LDL) cholesterol (mg/dL) * hdl: High-Density Lipoprotein (HDL) cholesterol (mg/dL) * systolic_bp: Systolic blood pressure (mmHg) * diastolic_bp: Diastolic blood pressure (mmHg) * smoking: Smoking status (0 = Non-Smoker, 1 = Smoker) * diabetes: Diabetes status (0 = No, 1 = Yes) * heart_attack: History of heart attack (0 = No, 1 = Yes)

Note: This dataset is synthetically generated and does not represent real patients. It is meant for research and educational purposes only.

```
[3]: # look at the shape of our dataset
print(heart_data.shape)
```

(1000, 10)

```
[4]: # take a look of some rows in our dataset
heart_data.head(10)
```

```
[4]:   age  sex  total_cholesterol      ldl      hdl  systolic_bp  \
0   57   1         229.463642  175.879129  39.225687   124.070127
1   58   1         186.464120  128.984916  34.950968    95.492552
2   37   1         251.300719  152.347592  45.913288    99.519335
3   55   1         192.058908  116.803684  67.208925   122.460002
4   53   1         151.203448  107.017396  60.693838   123.022257
5   39   1         236.033455  153.880809  31.208614   121.857396
6   65   0         174.615665  114.029407  55.692586   135.605050
7   33   0         242.919402  147.951375  54.439475   123.511557
8   49   0          95.804359   83.304875  60.758929   111.697488
9   55   0         181.360943  106.011783  50.576747   129.576418

      diastolic_bp  smoking  diabetes  heart_attack
0      91.378780         0         0             0
1      64.355040         1         0             0
2      64.953147         0         1             0
3      73.821382         0         0             0
4      81.121946         0         1             0
5      79.589069         0         0             0
6      85.529955         0         0             0
7      77.331714         0         0             0
8      77.630529         1         0             0
9      87.588781         0         0             0
```

0.2 Exploratory Data Analysis (EDA)

```
[5]: # We check for missng and duplicated observations
number_of_missing = heart_data.isnull().sum()
number_of_duplicated = heart_data.duplicated().sum()

print(number_of_missing, "\n ")
print(number_of_duplicated, "duplicated observations")
```

```
age          0
sex          0
total_cholesterol  0
ldl          0
hdl          0
systolic_bp  0
diastolic_bp  0
smoking      0
diabetes     0
heart_attack  0
dtype: int64
```

0 duplicated observations

```
[6]: # summary of continuous variable in our dataset
continuous_variable =
    heart_data[["age", "total_cholesterol", "ldl", "hdl", "systolic_bp", "diastolic_bp"]]
continuous_variable.describe()
```

```
[6]:
```

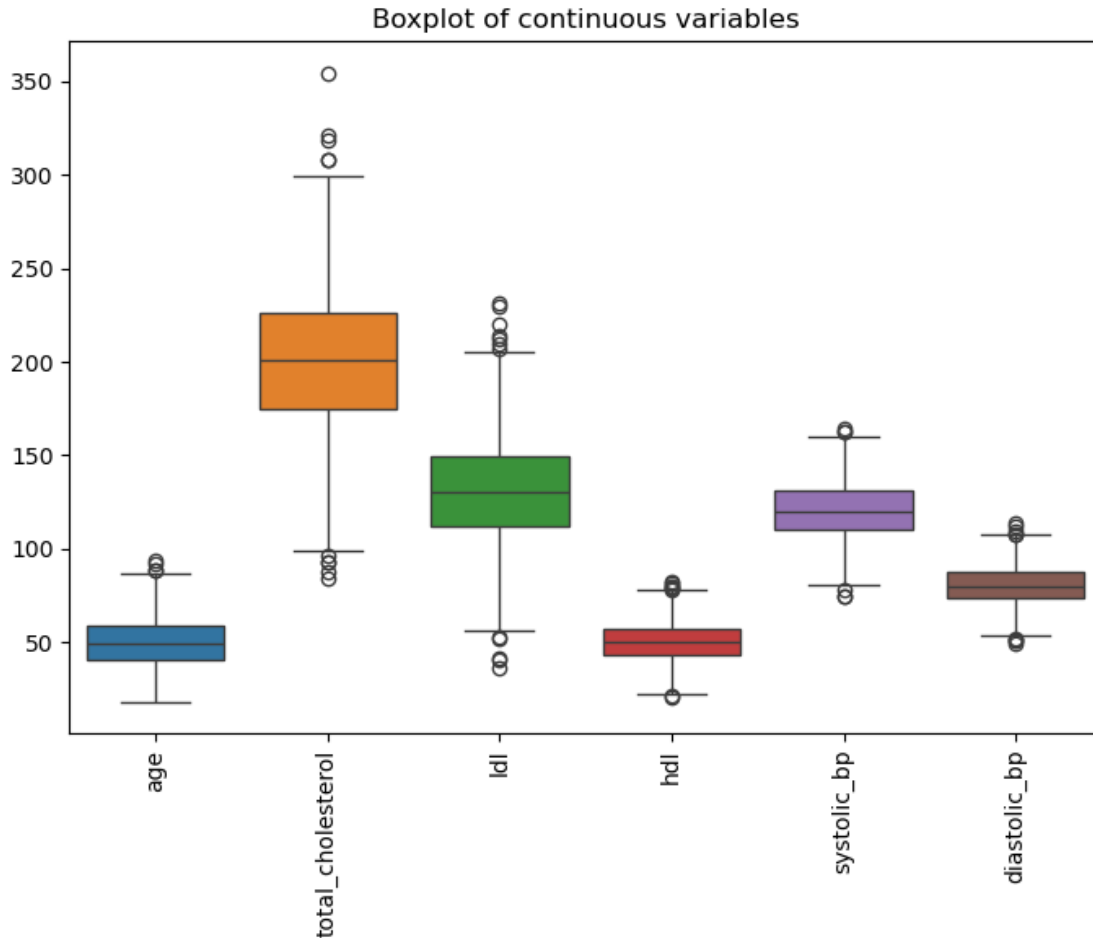
	age	total_cholesterol	ldl	hdl	systolic_bp \
count	1000.000000	1000.000000	1000.000000	1000.000000	1000.000000
mean	49.886000	201.087486	130.047807	49.811244	120.312687
std	14.209466	40.042655	30.041659	10.247178	15.507493
min	18.000000	84.165932	36.259745	20.600644	74.433950
25%	40.000000	174.707208	111.963197	42.622102	110.062952
50%	49.000000	201.191547	130.678540	49.682809	120.042175
75%	59.000000	226.251708	149.732446	56.703598	130.911804
max	94.000000	354.660015	231.376631	82.319810	164.080967

	diastolic_bp
count	1000.000000
mean	80.231248
std	10.235917
min	49.296305
25%	73.277119
50%	79.912592
75%	87.084443
max	113.848127

```
[7]: # boxplot of the continuous variables
plt.figure(figsize=(7, 6))
sns.boxplot(data=continuous_variable)

# Set labels and title
plt.title('Boxplot of continuous variables')
plt.xticks(rotation=90)
plt.tight_layout()

# Show the plot
plt.show()
```

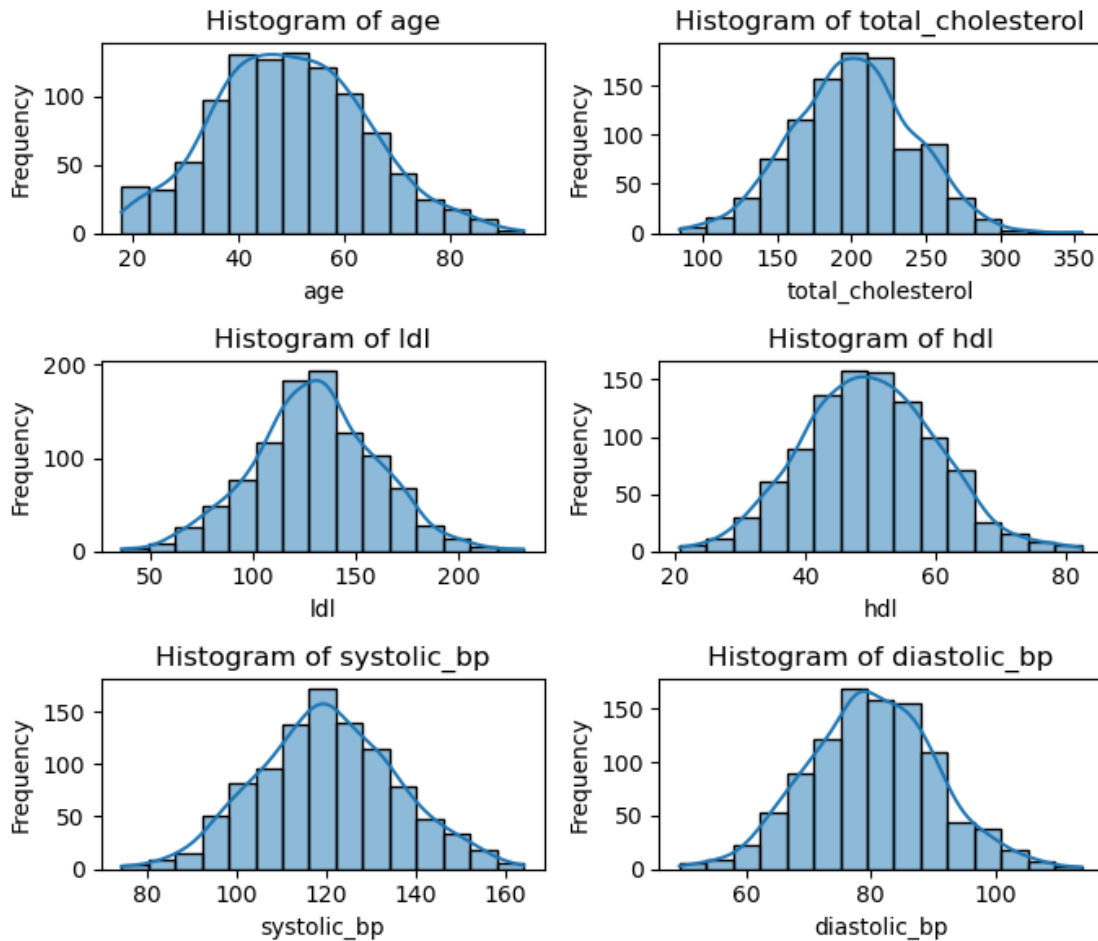


0.3 Insight and interpretation:

Here we can see all the boxes are centered around the mean, which suggests that they all follow a symmetric distribution. In each case, there are some extreme values which are far from the other values. We should pay attention to those observations as they have a direct effect on increasing the risk of getting a heart attack.

```
[8]: # What about the distribution of the variable in our dataset
# Histogramme of the continuous variables.
plt.figure(figsize=(7, 6))
# Loop through each continuous variable to plot their histogram
for i, col in enumerate(continuous_variable.columns, 1):
    plt.subplot(3, 2, i) # Adjust the number of rows and columns for subplots
    sns.histplot(continuous_variable[col], bins=15, kde=True) # Plot histogram
    # with KDE
    plt.title(f'Histogram of {col}')
    plt.xlabel(col)
```

```
plt.ylabel('Frequency')
# Adjust layout
plt.tight_layout()
# Show the plot
plt.show()
```



Here we can see that the suggestions give by the box plot are confirmed since all these variables are following normal distribution.

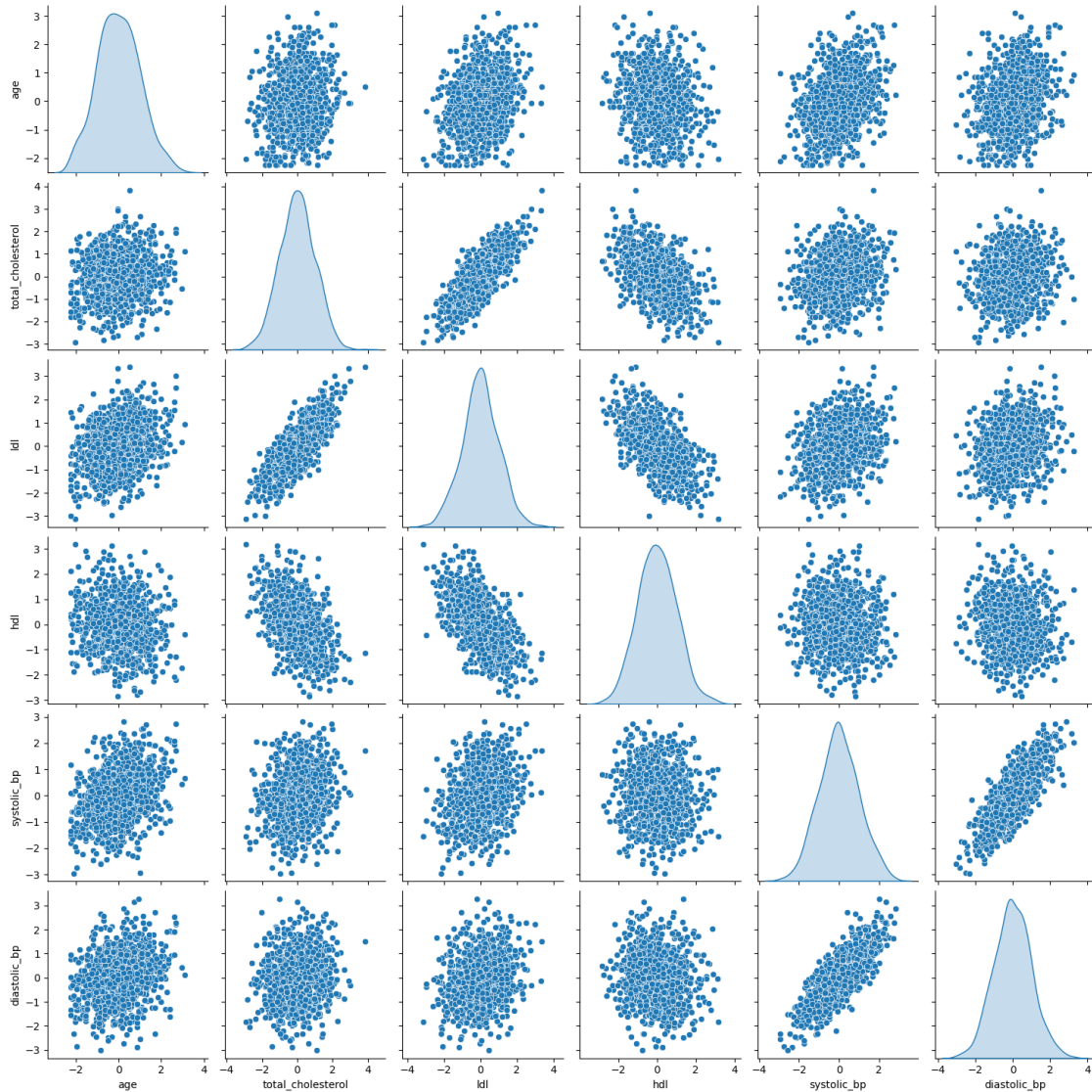
```
[9]: # Scaling the variable to be able to plot their scatter plot
scaler = StandardScaler()
scaled_data = scaler.fit_transform(continuous_variable)
# Convert back to a DataFrame with the same column names
scaled_df = pd.DataFrame(scaled_data, columns=continuous_variable.columns)
scaled_df.describe()
```

```
[9]:
```

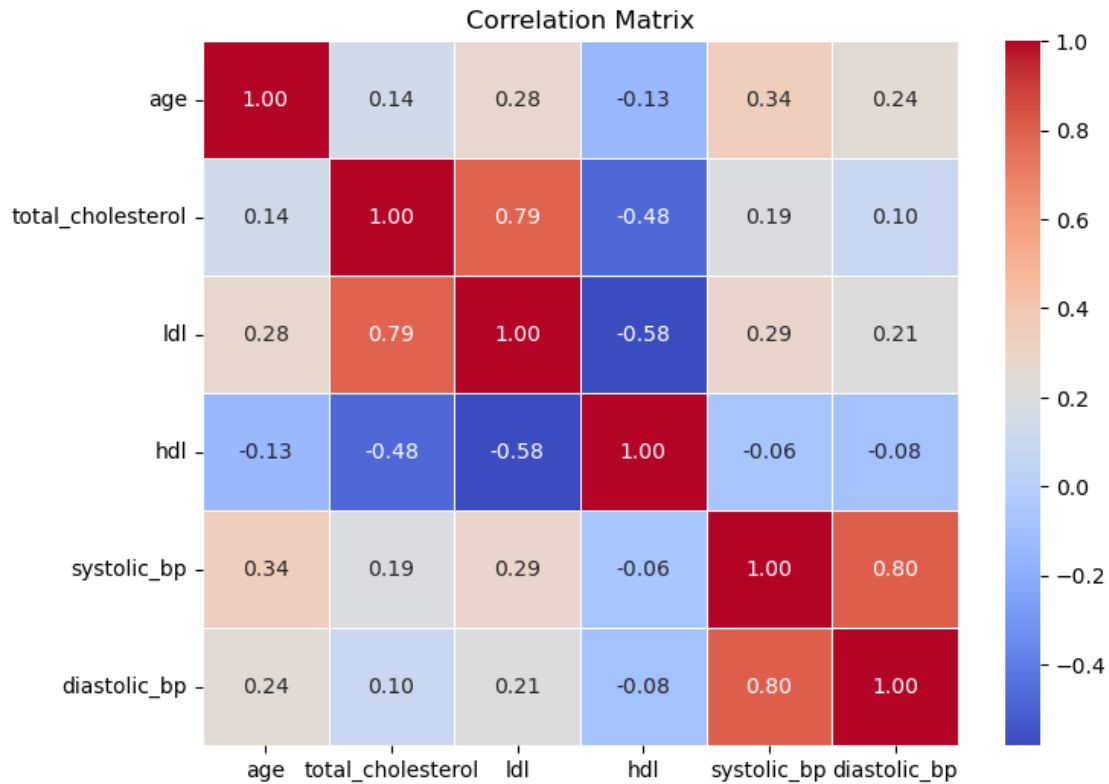
	age	total_cholesterol	ldl	hdl \
count	1.000000e+03	1.000000e+03	1.000000e+03	1.000000e+03
mean	-2.025047e-16	-2.842171e-16	-4.316547e-16	6.217249e-17
std	1.000500e+00	1.000500e+00	1.000500e+00	1.000500e+00
min	-2.245120e+00	-2.921386e+00	-3.123496e+00	-2.852026e+00
25%	-6.960815e-01	-6.591341e-01	-6.022856e-01	-7.019240e-01
50%	-6.238400e-02	2.600065e-03	2.100577e-02	-1.253994e-02
75%	6.417244e-01	6.287499e-01	6.555726e-01	6.729465e-01
max	3.106104e+00	3.837143e+00	3.374631e+00	3.174028e+00

	systolic_bp	diastolic_bp
count	1.000000e+03	1.000000e+03
mean	-7.460699e-17	5.710987e-16
std	1.000500e+00	1.000500e+00
min	-2.959969e+00	-3.023708e+00
25%	-6.612845e-01	-6.797250e-01
50%	-1.745272e-02	-3.114672e-02
75%	6.838256e-01	6.698593e-01
max	2.823808e+00	3.285851e+00

```
[ ]: # Scatter plot matrix of our variables
      # Create a scatter plot matrix
      sns.pairplot(scaled_df, diag_kind='kde', markers='o')
      # Show the plot
      plt.show()
```



```
[ ]: # Correlation matrix of the dataset
# Compute the correlation matrix
correlation_matrix = scaled_df.corr(numeric_only=True) # Exclude non-numeric
# Plot the heatmap
plt.figure(figsize=(8, 6))
sns.heatmap(correlation_matrix, annot=True, cmap="coolwarm", fmt=".2f",
            linewidths=0.5)
# Show the plot
plt.title("Correlation Matrix")
plt.show()
```



1 Statistical analysis

1.1 What is the effect of smoking on heart_attack

```
[35]: # Compare heart attack rates for smokers and non-smokers
summary = heart_data.groupby('smoking')['heart_attack'].mean()
print(summary)
```

```
smoking
0    0.071429
1    0.232673
Name: heart_attack, dtype: float64
```

1.1.1 Interpretation of our result

Smokers have a heart attack rate that is more than three times higher than non-smokers. this suggest us there is possibly significant impact between being a smoker and get a heart attack to confirm this suggestion we are going to run a chi square test for independance between the smoking and heart_attack variables.

```
[ ]: import scipy.stats as stats
# Create contingency table
```



```
contingency_table = pd.crosstab(heart_data['smoking'],
    ↪heart_data['heart_attack'])
# Perform Chi-Square Test
chi2, p, dof, expected = stats.chi2_contingency(contingency_table)
print(f"Chi-Square Statistic: {chi2}, p-value: {p}")
```

Chi-Square Statistic: 43.26255952585468, p-value: 4.786553985919395e-11

1.1.2 Interpretation of our result

Since the p-value (4.79e-11) is much smaller than 0.05, we reject the null hypothesis. This means that there is a statistically significant association between smoking and heart attacks. In other words, smoking has a meaningful effect on the likelihood of having a heart attack.

```
[ ]: import numpy as np
# Chi-Square statistic and sample size
n = len(heart_data) # Total sample size
k = min(len(heart_data['smoking'].unique()), len(heart_data['heart_attack'].
    ↪unique())) # Smallest category count
# Compute Cramér's V
cramers_v = np.sqrt(chi2 / (n * (k - 1)))
print(f"Cramér's V: {cramers_v:.4f}")
```

Cramér's V: 0.2080

1.1.3 Interpretation of our cramer's V value

0.2080 suggests a small-to-moderate association between smoking and heart attacks. We can say at this time that, while statistically significant, the association between smoking and heart_attack is not very strong.

```
[42]: # Does the level of cholesterol and the age increase the risk of heart attack ?

X = heart_data[['age', 'total_cholesterol']] # Independent variables
X = sm.add_constant(X) # Adds an intercept (constant term)
y = heart_data['heart_attack'] # Dependent variable (1 = heart attack, 0 = no
    ↪heart attack)

# Fit logistic regression model
logit_model = sm.Logit(y, X)
result = logit_model.fit()

# Print summary
print(result.summary())
```

Optimization terminated successfully.

Current function value: 0.314030

Iterations 7

Logit Regression Results

```

=====
Dep. Variable:          heart_attack    No. Observations:          1000
Model:                  Logit           Df Residuals:              997
Method:                 MLE            Df Model:                  2
Date:                  Wed, 12 Feb 2025    Pseudo R-squ.:            0.05918
Time:                  12:33:44          Log-Likelihood:           -314.03
converged:              True            LL-Null:                  -333.78
Covariance Type:        nonrobust        LLR p-value:              2.636e-09
=====
=====
              coef      std err          z      P>|z|      [0.025
0.975]
-----
const          -6.2149      0.706     -8.804      0.000     -7.598
-4.831
age             0.0208      0.007      2.785      0.005      0.006
0.036
total_cholesterol  0.0142      0.003      5.138      0.000      0.009
0.020
=====
=====

```

1.1.4 Interpretation of our logistic regression results

1. For the P-Values both age and cholesterol P-values respectively ($p = 0.005$ and $p < 0.001$) are statistically significant which suggest that an increase either in age or in cholesterol levels increase heart attack risk. Intercept ($p < 0.001$) → Indicates the baseline log-odds when all predictors are zero.
2. After computing the odds ratios: age: $e^{0.0208} = 1.021e^{0.0208} = 1.021$ → each additional year of age increases the odds of a heart attack by ~ 2.1%. cholesterol: $e^{0.0142} = 1.014e^{0.0142} = 1.014$ → Each unit increase in cholesterol increases the odds of a heart attack by ~ 1.4%.
3. Model Fit & Strength Pseudo $R^2 = 0.059$ → which suggests a weak-to-moderate explanatory power. The LLR p-value (2.636e-09) indicates that the model as a whole is statistically significant.

1.1.5 Key insight:

Both age and cholesterol significantly increase the risk of heart attacks. Cholesterol has a slightly stronger effect than age (based on the z-scores and odds ratios). The model is statistically significant but doesn't explain all variation, meaning other factors (e.g., smoking,) should be included to improve prediction.

```
[ ]: # Does it have more heart attack in a given category of sexe ?
summary = heart_data.groupby('sex')['heart_attack'].mean()
print(summary)
```

```
sex
0    0.067653
1    0.136622
Name: heart_attack, dtype: float64
```

1.1.6 Interpretation:

Our output indicates the mean heart attack rates for each sex category in the dataset: For sex = 0 (likely representing females), the heart attack rate is approximately 0.0677 (or 6.77%). For sex = 1 (likely representing males), the heart attack rate is approximately 0.1366 (or 13.66%).

So in this sample, males have a higher heart attack rate than females.

```
[43]: # how change the risk of heart attack with the ldl: Low-Density Lipoprotein,
      ↪ hdl: High-Density Lipoprotein
correlation_ldl = heart_data['ldl'].corr(heart_data['heart_attack'])
correlation_hdl = heart_data['hdl'].corr(heart_data['heart_attack'])

print(f"Correlation between LDL and heart attack: {correlation_ldl}")
print(f"Correlation between HDL and heart attack: {correlation_hdl}")
```

```
Correlation between LDL and heart attack: 0.15896458816666045
Correlation between HDL and heart attack: -0.14533819654661279
```

The correlation results indicate the following: - LDL and heart attack: The correlation is 0.159, suggesting a weak positive relationship between LDL and the occurrence of heart attacks. This means that as LDL levels increase, the likelihood of a heart attack slightly increases, but the correlation is not very strong. - HDL and heart attack: The correlation is -0.145, suggesting a weak negative relationship between HDL and heart attacks. This implies that as HDL levels increase, the likelihood of a heart attack slightly decreases, but again, the relationship is weak.

These correlations are not very strong, which could suggest that other factors in the dataset might also play a role in determining heart attack risk.

```
[45]: # Assume 'ldl' and 'hdl' are columns in the dataset
X = heart_data[['ldl', 'hdl']] # Independent variables
X = sm.add_constant(X) # Adds a constant term to the model
y = heart_data['heart_attack'] # Dependent variable

model = sm.Logit(y, X) # Logistic regression for binary outcomes
result = model.fit()

print(result.summary())
```

Optimization terminated successfully.

Current function value: 0.318750

Iterations 7

Logit Regression Results

```
=====
Dep. Variable:          heart_attack    No. Observations:          1000
```

Model:	Logit	Df Residuals:	997			
Method:	MLE	Df Model:	2			
Date:	Wed, 12 Feb 2025	Pseudo R-squ.:	0.04504			
Time:	12:55:19	Log-Likelihood:	-318.75			
converged:	True	LL-Null:	-333.78			
Covariance Type:	nonrobust	LLR p-value:	2.956e-07			
=====						
	coef	std err	z	P> z	[0.025	0.975]

const	-2.5420	1.073	-2.369	0.018	-4.645	-0.439
ldl	0.0125	0.004	2.880	0.004	0.004	0.021
hdl	-0.0274	0.013	-2.121	0.034	-0.053	-0.002
=====						

1.1.7 interpretation of the coefficient of our model :

: cholesterol (mg/dL) * ldl (mg/dL):

The coefficient for Low-Density Lipoprotein (LDL) cholesterol is 0.0125 with a p-value of 0.004, which is highly significant ($p < 0.05$). Suggesting a positive relationship between LDL and heart attack risk LDL, with an increase of the odds of having a heart attack by a factor of $e^{0.0125} \approx 1.0126$ (i.e., a 1.26%) when the increasing the LDL per mg/dl,

- hdl (mg/dL): The coefficient for High-Density Lipoprotein (HDL) cholesterol is -0.0274 with a p-value of 0.034, which is also statistically significant ($p < 0.05$). This indicates a negative relationship between HDL and heart attack risk with a decrease in the odds of having a heart attack decrease by a factor of $e^{-0.0274} \approx 0.973$ (i.e., a 2.7%) for every mg/dl increase in HDL.

1.1.8 Overall model fit:

- Pseudo R-squared: 0.04504, which is relatively low, indicating that the model explains only a small portion of the variance in heart attack risk. This suggests that there are other important factors not included in the model that affect heart attack risk.
- Log-Likelihood: -318.75, and the LLR p-value is 2.956e-07, which indicates that the model is statistically significant and provides a better fit than the null model (no predictors).

Insights :

- LDL: Higher levels of LDL increase the risk of a heart attack, supporting the well-known link between high LDL cholesterol and heart disease.
- HDL: Higher levels of HDL appear to lower the risk of a heart attack, which aligns with the protective role of HDL in cardiovascular health.

1.2 Build classification machine learning models

```
[ ]: # we train first a full linear classification model.
# Define the features (independent variables) and target (dependent variable)
X = heart_data[['age', 'sex', 'total_cholesterol', 'ldl', 'hdl', 'systolic_bp',
↪ 'diastolic_bp', 'smoking', 'diabetes']]
y = heart_data['heart_attack']
```

```

# Standardize continuous variables for better model performance
scaler = StandardScaler()
X[['age', 'total_cholesterol', 'ldl', 'hdl', 'systolic_bp', 'diastolic_bp']] = scaler.fit_transform(
    X[['age', 'total_cholesterol', 'ldl', 'hdl', 'systolic_bp', 'diastolic_bp']]
)

# Add a constant term for statsmodels
X = sm.add_constant(X)

```

```
[46]: # interpretation of it's coefficient and variance analysis.
```

```
[ ]: # we perform feature importance to see if we can build small model which perform
      ↪ better
```

```
[ ]: # build a model with the important features
```

```
[ ]: # interpretation of it's coefficient and variance analysis
```

```
[ ]: # select feature by step-wise selection and see if the selected features will
      ↪ be the same selected when we were running the feature importance
```

```
[ ]: # we build different models using different algorithms
```

```
[ ]: # we compare the different obtained models
```