# Assignment 3

**Problem Statement:**

Write a program to construct a Bayesian network considering medical data. Use this model to demonstrate the diagnosis of heart patients using the standard Heart Disease Data Set (You can use Java/Python ML library classes/API.

**Objective:**
1. Evaluate and analyse retrieved information.
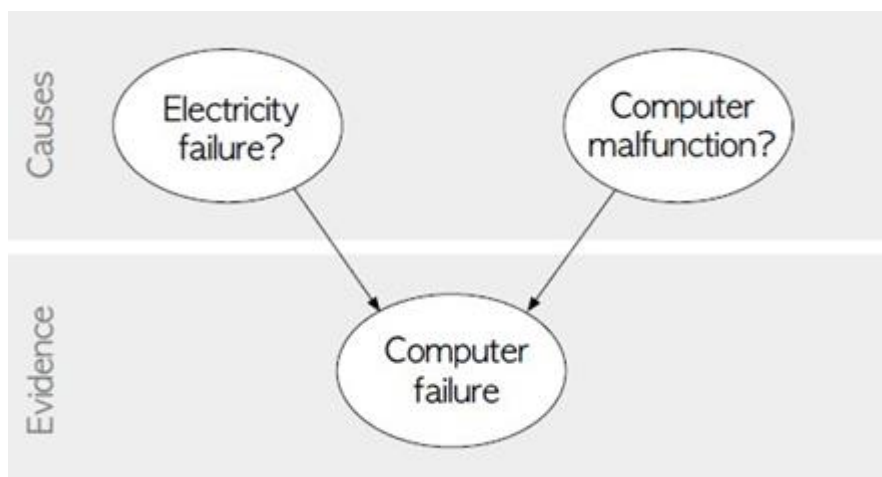2. To study Bayesian network model.

**Theory:**

A Bayesian network is a directed acyclic graph in which each edge corresponds to a conditional dependency, and each node corresponds to a unique random variable.

Bayesian network consists of two major parts: a directed acyclic graph and a set of conditional probability distributions

- The directed acyclic graph is a set of random variables represented by nodes.
- The conditional probability distribution of a node (random variable) is defined for every possible outcome of the preceding causal node(s).

For illustration, consider the following example. Suppose we attempt to turn on our computer, but the computer does not start (observation/evidence). We would like to know which of the possible causes of computer failure is more likely. In this simplified illustration, we assume only two possible causes of this misfortune: electricity failure and computer malfunction.

The corresponding directed acyclic graph is depicted in below figure.



The goal is to calculate the posterior conditional probability distribution of each of the possible unobserved causes given the observed evidence, i.e., P [Cause | Evidence].

**Data Set:**

**Title:** Heart Disease Databases
The Cleveland database contains 76 attributes, but all published experiments refer to using a subset of 14 of them. In particular, the Cleveland database is the only one that has been used by ML researchers to this date. The "Heartdisease" field refers to the presence of heart disease in the patient. It is integer valued from 0 (no presence) to 4.

Database:    0     1     2     3     4     Total
Cleveland:   164   55    36    35    13    303

**Attribute Information:**
1. age: age in years
2. sex: sex (1 = male; 0 = female)
3. cp: chest pain type
   1. Value 1: typical angina
   2. Value 2: atypical angina
   3. Value 3: non-anginal pain
   4. Value 4: asymptomatic
4. trestbps: resting blood pressure (in mm Hg on admission to the hospital)
5. chol: serum cholestoral in mg/dl
6. fbs: (fasting blood sugar > 120 mg/dl) (1 = true; 0 = false)
7. restecg: resting electrocardiographic results
   1. Value 0: normal
   2. Value 1: having ST-T wave abnormality (T wave inversions and/or ST elevation or depression of > 0.05 mV)
   3. Value 2: showing probable or definite left ventricular hypertrophy by Estes' criteria
8. thalach: maximum heart rate achieved
9. exang: exercise induced angina (1 = yes; 0 = no)
10. oldpeak = ST depression induced by exercise relative to rest
11. slope: the slope of the peak exercise ST segment
    1. Value 1: upsloping
    2. Value 2: flat
    3. Value 3: downsloping
12. thal: 3 = normal; 6 = fixed defect; 7 = reversable defect
13. Heartdisease: It is integer valued from 0 (no presence) to 4.

**Some instance from the dataset:**

| age | sex | cp | trestbps | chol | fbs | restecg | thalach | exang | oldpeak | slope | ca | thal | Heart disease |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 63 | 1 | 1 | 145 | 233 | 1 | 2 | 150 | o | 2.3 | 3 | o | 6 | o |
| 67 | 1 | 4 | 160 | 286 | o | 2 | 108 | 1 | 1.5 | 2 | 3 | 3 | 2 |

| 67 | 1 | 4 | 120 | 229 | o | 2 | 129 | 1 | 2.6 | 2 | 2 | 7 | 1 |
|----|---|---|-----|-----|---|---|-----|---|-----|---|---|---|---|
| 41 | o | 2 | 130 | 204 | o | 2 | 172 | o | 1.4 | 1 | o | 3 | o |
| 62 | o | 4 | 140 | 268 | o | 2 | 160 | o | 3.6 | 3 | 2 | 3 | 3 |
| 60 | 1 | 4 | 130 | 206 | o | 2 | 132 | 1 | 2.4 | 2 | 2 | 7 | 4 |

**Python Program to Implement and Demonstrate Bayesian network using pgmpy Machine Learning**

```
import numpy as np
import pandas as pd
import csv
from pgmpy.estimators import MaximumLikelihoodEstimator
from pgmpy.models import BayesianModel
from pgmpy.inference import VariableElimination
#read Cleveland Heart Disease data
heartDisease = pd.read_csv('heart.csv')
heartDisease = heartDisease.replace('?',np.nan)
#display the data
print('Sample instances from the dataset are given below')
print(heartDisease.head())
print('\n Attributes and datatypes')
print(heartDisease.dtypes)
#Model Bayesian Network
model=BayesianModel([('age','heartdisease'),('sex','heartdisease'),('exang','heartdisease'),('cp','heartdisea
se'),('heartdisease','restecg'),('heartdisease','chol')])
#Learning CPDs using Maximum Likelihood Estimators
print('\nLearning CPD using Maximum likelihood estimators')
model.fit(heartDisease,estimator=MaximumLikelihoodEstimator)
# Inferencing with Bayesian Network
print('\n Inferencing with Bayesian Network:')
HeartDiseasetest_infer = VariableElimination(model)
#computing the Probability of HeartDisease given Age
print('\n 1. Probability of HeartDisease given evidence= restecg')
q1=HeartDiseasetest_infer.query(variables=['heartdisease'],evidence={'restecg':1})
print(q1)
#computing the Probability of HeartDisease given cholesterol
print('\n 2. Probability of HeartDisease given evidence= cp ')
q2=HeartDiseasetest_infer.query(variables=['heartdisease'],evidence={'cp':2})
print(q2)
```

**Output**

```
Learning CPD using Maximum likelihood estimators

 Inferencing with Bayesian Network:

1. Probability of HeartDisease given evidence= restecg

    +----------------+---------------------+
    | heartdisease   |    phi(heartdisease) |
    +================+=====================+
    | heartdisease(0) |             0.1012 |
    +----------------+---------------------+
    | heartdisease(1) |             0.0000 |
    +----------------+---------------------+
    | heartdisease(2) |             0.2392 |
    +----------------+---------------------+
    | heartdisease(3) |             0.2015 |
    +----------------+---------------------+
    | heartdisease(4) |             0.4581 |
    +----------------+---------------------+
2. Probability of HeartDisease given evidence= cp

    +----------------+---------------------+
    | heartdisease   |    phi(heartdisease) |
    +================+=====================+
    | heartdisease(0) |             0.3610 |
    +----------------+---------------------+
    | heartdisease(1) |             0.2159 |
    +----------------+---------------------+
    | heartdisease(2) |             0.1373 |
    +----------------+---------------------+
    | heartdisease(3) |             0.1537 |
    +----------------+---------------------+
    | heartdisease(4) |             0.1321 |
    +----------------+---------------------+
```