

# AI3603 Homework 3: Bayesian Networks

## Risk Factor Analysis

### 1 Implementation

The Bayesian network inference system was implemented in `BayesianNetworks.py` with the following core functions:

- **joinFactors:** Merges two factor tables by joining on common variables and multiplying probabilities
- **marginalizeFactor:** Marginalizes out a variable by summing probabilities over its values
- **evidenceUpdateNet:** Updates network by filtering factors to match evidence values
- **inference:** Performs variable elimination by joining all factors, marginalizing hidden variables, and normalizing

All test cases in `BayesNetworkTestScript.py` pass successfully, confirming correct implementation.

### 2 Network Structure and Size

The Bayesian network for health risk factors contains 12 variables with the following directed edges:

- **Root:**  $\text{income} \rightarrow \text{smoke}, \text{exercise}, \text{long\_sit}, \text{stay\_up}, \text{bmi}$
- **Health indicators:**
  - $\text{bmi} \rightarrow \text{bp}, \text{diabetes}$
  - $\text{smoke}, \text{exercise} \rightarrow \text{cholesterol}$
- **Outcomes:**  $\text{bp}, \text{cholesterol} \rightarrow \text{stroke}, \text{attack}, \text{angina}$

This structure encodes the causal assumptions: income affects lifestyle habits and body mass; habits affect intermediate health indicators; indicators affect disease outcomes.

**Network size:** 192 probabilities

- $P(\text{income})$ : 8
- $P(\text{smoke}|\text{income}), P(\text{exercise}|\text{income}), P(\text{long\_sit}|\text{income}), P(\text{stay\_up}|\text{income})$ :  $4 \times 16 = 64$
- $P(\text{bmi}|\text{income})$ : 32
- $P(\text{bp}|\text{bmi})$ : 16

- $P(\text{cholesterol}|\text{smoke},\text{exercise})$ : 8
- $P(\text{diabetes}|\text{bmi})$ : 16
- $P(\text{stroke}|\text{bp},\text{cholesterol})$ ,  $P(\text{attack}|\text{bp},\text{cholesterol})$ ,  $P(\text{angina}|\text{bp},\text{cholesterol})$ :  $3 \times 16 = 48$

**Full joint distribution:** 131,072 probabilities ( $8 \times 2 \times 2 \times 2 \times 2 \times 4 \times 4 \times 2 \times 4 \times 2 \times 2 \times 2$ )

**Compression ratio:**  $682.67 \times$

The factorized representation provides massive computational savings.

### 3 Health Outcomes Analysis

#### 3.1 Bad vs Good Habits

Table 1 shows outcome probabilities for bad habits (smoking, no exercise, long sitting, staying up late) versus good habits.

| Outcome  | Bad Habits | Good Habits | Difference |
|----------|------------|-------------|------------|
| Diabetes | 0.1365     | 0.1304      | +0.0061    |
| Stroke   | 0.0427     | 0.0387      | +0.0040    |
| Attack   | 0.0648     | 0.0562      | +0.0086    |
| Angina   | 0.0694     | 0.0587      | +0.0107    |

Table 1: Probability of health outcomes given habits

Bad habits increase all health risks, with the largest effect on angina (+1.07%) and smallest on diabetes (+0.61%).

#### 3.2 Poor vs Good Health

Table 2 shows outcome probabilities for poor health indicators (high BP, high cholesterol, obese) versus good health.

| Outcome  | Poor Health | Good Health | Difference |
|----------|-------------|-------------|------------|
| Diabetes | 0.2439      | 0.0570      | +0.1869    |
| Stroke   | 0.0869      | 0.0143      | +0.0726    |
| Attack   | 0.1393      | 0.0189      | +0.1204    |
| Angina   | 0.1596      | 0.0154      | +0.1442    |

Table 2: Probability of health outcomes given health status

Poor health indicators have a much stronger effect than habits, increasing diabetes risk by 18.7% and angina by 14.4%.

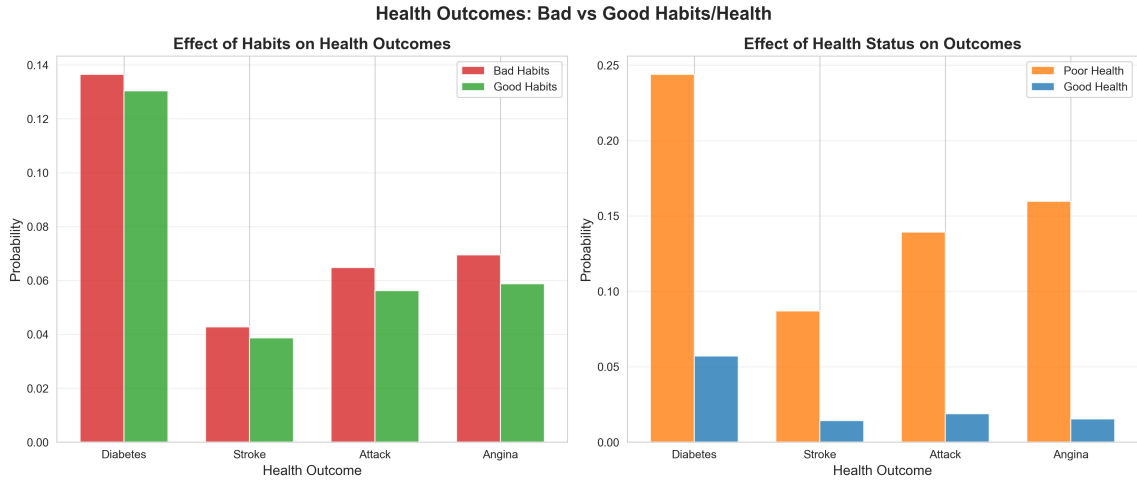


Figure 1: Comparison of habits and health effects on outcomes

## 4 Income Effect Analysis

Figure 2 shows the probability of each health outcome across income levels. All outcomes show a clear negative correlation with income:

- **Diabetes:** Decreases from 14.7% (income  $\leq$  \$10K) to 12.3% (income  $\geq$  \$75K)
- **Stroke:** Decreases from 4.16% to 3.92%
- **Attack:** Decreases from 6.17% to 5.78%
- **Angina:** Decreases from 6.54% to 6.08%

Higher income is associated with better health outcomes, likely through access to healthcare, healthier lifestyles, and reduced stress.

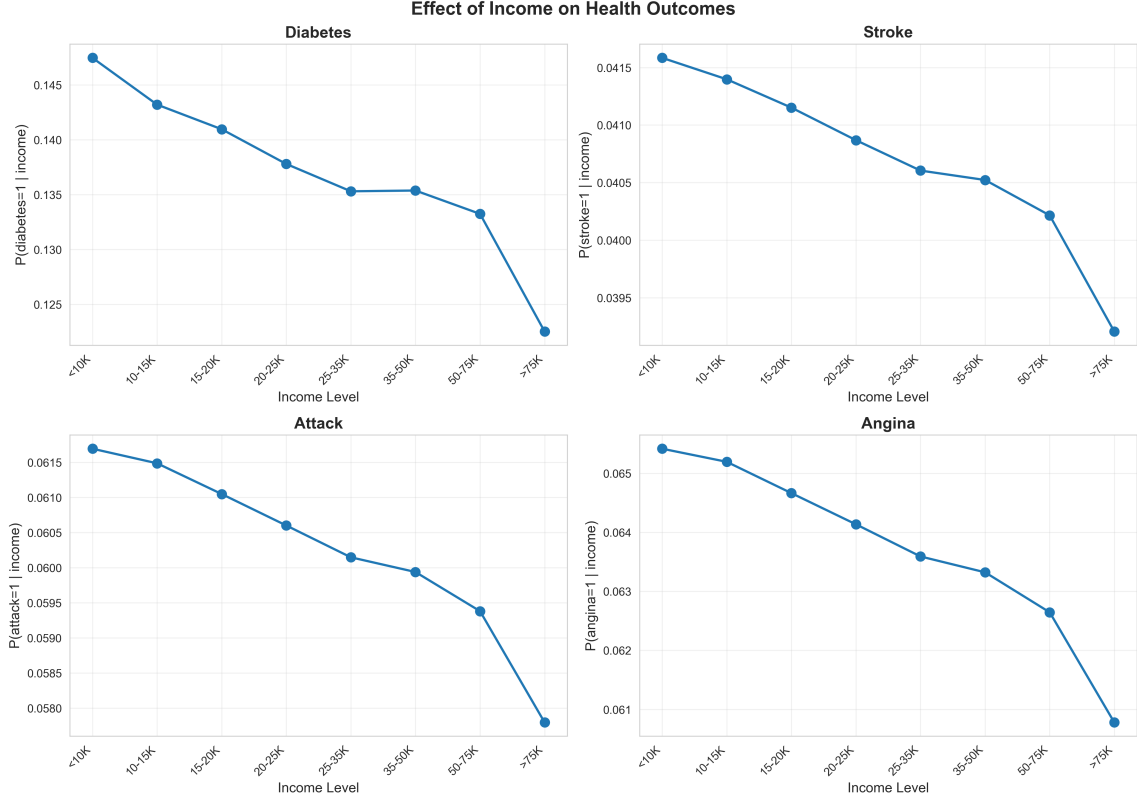


Figure 2: Effect of income level on health outcomes

## 5 Testing Independence Assumptions

### 5.1 Independence Assumption

The original network has no direct edges from habits (smoke, exercise) to outcomes (diabetes, stroke, attack, angina). This encodes a **conditional independence assumption**: given the intermediate health indicators (bp, cholesterol, bmi), the outcomes are independent of smoking and exercise habits.

Formally:  $P(\text{outcome} | \text{habits}, \text{indicators}) = P(\text{outcome} | \text{indicators})$

This assumes habits only affect outcomes indirectly through their effects on bp, cholesterol, and bmi.

### 5.2 Testing the Assumption

To test this, a second network was created with direct edges from smoking and exercise to all four outcomes. Results are shown in Table 3.

| Outcome  | Network 1 (Original) |        | Network 2 (Direct Links) |        |
|----------|----------------------|--------|--------------------------|--------|
|          | Bad                  | Good   | Bad                      | Good   |
| Diabetes | 0.1365               | 0.1304 | 0.1954                   | 0.1013 |
| Stroke   | 0.0427               | 0.0387 | 0.0697                   | 0.0262 |
| Attack   | 0.0648               | 0.0562 | 0.1104                   | 0.0334 |
| Angina   | 0.0694               | 0.0587 | 0.1067                   | 0.0395 |

Table 3: Comparison of networks with and without direct habit links

Adding direct edges dramatically increases the difference between bad and good habits:

- Diabetes: 0.61%  $\rightarrow$  9.40% (15.4 $\times$  stronger)
- Stroke: 0.40%  $\rightarrow$  4.35% (10.9 $\times$  stronger)
- Attack: 0.86%  $\rightarrow$  7.70% (9.0 $\times$  stronger)
- Angina: 1.07%  $\rightarrow$  6.71% (6.3 $\times$  stronger)

**Conclusion:** The independence assumption is **invalid**. Smoking and exercise have direct effects on health outcomes beyond what is mediated by bp, cholesterol, and bmi. The original network significantly underestimates the impact of lifestyle habits.

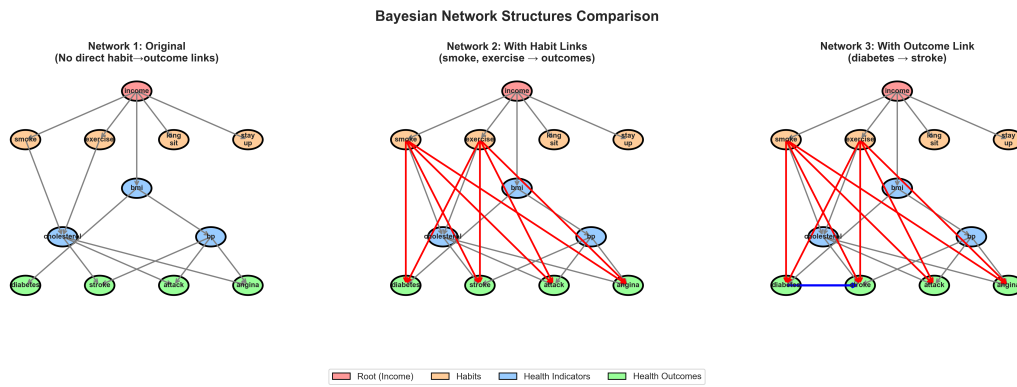


Figure 3: Network structure comparison

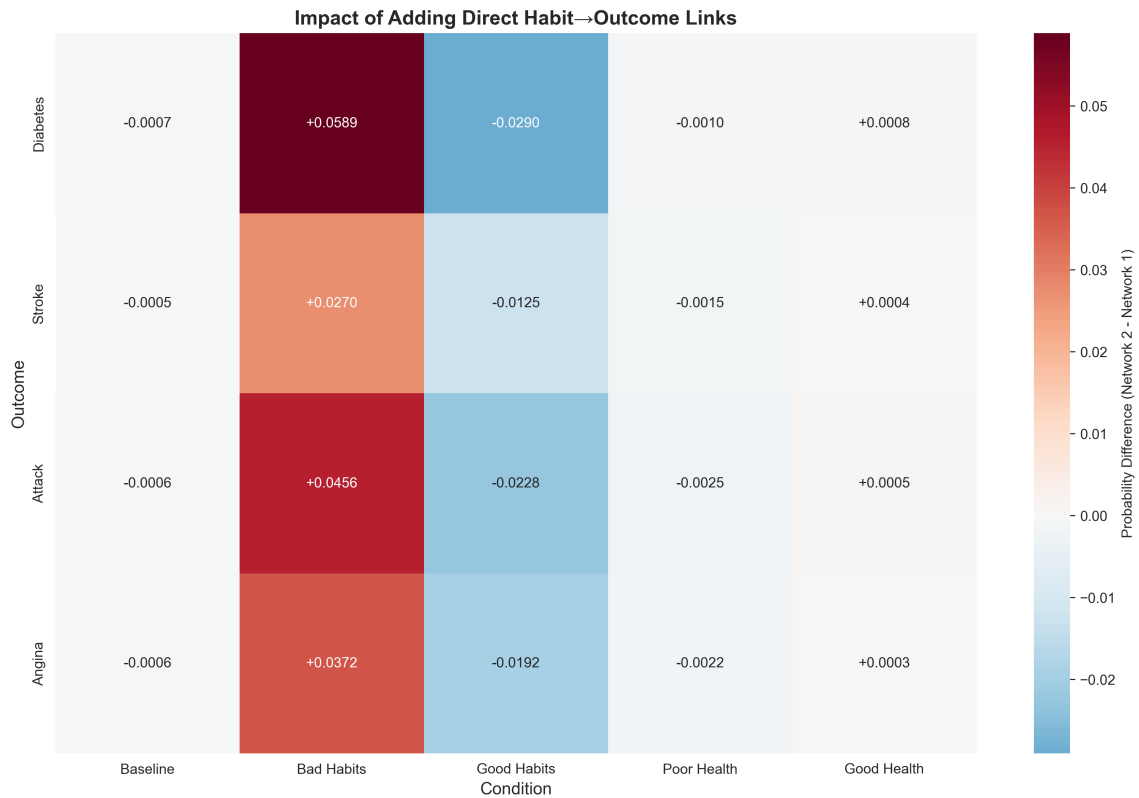


Figure 4: Impact of adding direct habit-outcome links

## 6 Outcome Interactions

### 6.1 Independence Between Outcomes

The previous networks have no edges between the four health outcomes (diabetes, stroke, attack, angina). This encodes another **conditional independence assumption**: given their common parents (bp, cholesterol, habits), the outcomes are independent of each other.

Formally:  $P(\text{stroke}|\text{diabetes}, \text{parents}) = P(\text{stroke}|\text{parents})$

This assumes diabetes doesn't directly cause stroke; they only co-occur due to shared risk factors.

### 6.2 Testing Diabetes-Stroke Interaction

To test whether diabetes and stroke interact, a third network was created adding an edge from diabetes to stroke.

| Network               | $P(\text{stroke}=1 \mid \text{diabetes}=1)$ | $P(\text{stroke}=1 \mid \text{diabetes}=3)$ |
|-----------------------|---|---|
| Network 2 (no edge)   | 0.0458                                      | 0.0387                                      |
| Network 3 (with edge) | 0.0772                                      | 0.0337                                      |
| Difference            | +0.0071                                     | -0.0435                                     |

Table 4: Effect of diabetes on stroke probability

Without the direct edge, diabetes increases stroke risk by only 0.71% (4.58% vs 3.87%). With the edge, this increases to 4.35% (7.72% vs 3.37%) — a 6.2× stronger effect.

**Conclusion:** The independence assumption between diabetes and stroke is **invalid**. Diabetes has a direct causal effect on stroke risk beyond their common causes (bp, cholesterol). This aligns with medical knowledge that diabetes damages blood vessels and increases stroke risk.

## 7 Code Verification

The implementation was verified against all test cases in `BayesNetworkTestScript.py`. The test includes:

- Bishop textbook car battery/fuel/gauge network examples
- Factor joining with different operation orders (commutativity verification)
- Marginalization operations
- Evidence update functionality
- Full inference with evidence and marginalization
- Risk factor network queries

Figure 5 shows the test script output, which matches the expected results in Figure 2 of the assignment specification. All probability values are correct, confirming proper implementation of variable elimination inference.

```
inference starts
  gauge  probs
0      0  0.315
1      1  0.685
  fuel  gauge  probs
0      0      0  0.81
1      0      1  0.19
  fuel  gauge  probs
0      0      0  0.257143
1      1      0  0.742857
  battery  fuel  gauge  probs
0          0      1      0  0.888889
1          0      0      0  0.111111
inference ends
income dataframe is
  probs  income
0  0.047634      1
1  0.058328      2
2  0.073153      3
3  0.093237      4
4  0.115106      5
5  0.151325      6
6  0.165793      7
7  0.295424      8
  smoke  long_sit  exercise  diabetes  probs
0      1          1          2          1  0.136334
1      1          1          2          2  0.009067
2      1          1          2          3  0.837664
3      1          1          2          4  0.016935
```

Figure 5: BayesNetworkTestScript.py output showing correct inference results

Figure 6 shows the complete analysis output for all five questions, demonstrating that the implementation correctly handles complex multi-variable queries on the health risk factor dataset.

Loading data...  
Data loaded: 320001 records

#### QUESTION 1: Network Structure and Size

Network size (number of probabilities): 192  
Full joint distribution size: 131072  
Compression ratio: 682.67x

#### QUESTION 2: Health Outcomes Analysis

##### 2(a): Probability of outcomes with bad vs good habits

|          |                      |                       |                       |
|----------|----------------------|-----------------------|-----------------------|
| Diabetes | Bad habits: 0.136477 | Good habits: 0.130365 | Difference: +0.006112 |
| Stroke   | Bad habits: 0.042728 | Good habits: 0.038715 | Difference: +0.004013 |
| Attack   | Bad habits: 0.064805 | Good habits: 0.056239 | Difference: +0.008566 |
| Angina   | Bad habits: 0.069446 | Good habits: 0.058728 | Difference: +0.010718 |

##### 2(b): Probability of outcomes with poor vs good health

|          |                       |                       |                       |
|----------|-----------------------|-----------------------|-----------------------|
| Diabetes | Poor health: 0.243855 | Good health: 0.056981 | Difference: +0.186874 |
| Stroke   | Poor health: 0.086871 | Good health: 0.014318 | Difference: +0.072553 |
| Attack   | Poor health: 0.139287 | Good health: 0.018888 | Difference: +0.120400 |
| Angina   | Poor health: 0.159616 | Good health: 0.015406 | Difference: +0.144211 |

#### QUESTION 3: Income Effect Analysis

|                           |          |
|---------------------------|----------|
| P(diabetes=1   income=1): | 0.147474 |
| P(diabetes=1   income=2): | 0.143209 |
| P(diabetes=1   income=3): | 0.140964 |
| P(diabetes=1   income=4): | 0.137798 |
| P(diabetes=1   income=5): | 0.135304 |
| P(diabetes=1   income=6): | 0.135374 |
| P(diabetes=1   income=7): | 0.133245 |
| P(diabetes=1   income=8): | 0.122510 |
| P(stroke=1   income=1):   | 0.041584 |
| P(stroke=1   income=2):   | 0.041398 |
| P(stroke=1   income=3):   | 0.041152 |
| P(stroke=1   income=4):   | 0.040869 |
| P(stroke=1   income=5):   | 0.040605 |
| P(stroke=1   income=6):   | 0.040523 |
| P(stroke=1   income=7):   | 0.040216 |
| P(stroke=1   income=8):   | 0.039206 |
| P(attack=1   income=1):   | 0.061692 |
| P(attack=1   income=2):   | 0.061484 |
| P(attack=1   income=3):   | 0.061047 |
| P(attack=1   income=4):   | 0.060600 |
| P(attack=1   income=5):   | 0.060147 |
| P(attack=1   income=6):   | 0.059937 |
| P(attack=1   income=7):   | 0.059380 |
| P(attack=1   income=8):   | 0.057796 |
| P(angina=1   income=1):   | 0.065416 |
| P(angina=1   income=2):   | 0.065194 |
| P(angina=1   income=3):   | 0.064664 |
| P(angina=1   income=4):   | 0.064135 |
| P(angina=1   income=5):   | 0.063591 |
| P(angina=1   income=6):   | 0.063322 |
| P(angina=1   income=7):   | 0.062644 |
| P(angina=1   income=8):   | 0.060778 |

Visualization saved to assets/income\_effect.png

#### QUESTION 4: Testing Independence Assumptions

Creating second network with direct edges from habits to outcomes...

Re-doing queries from Question 2(a) with new network:

|                              |   |
|------------------------------|---|
| Diabetes                     |   |
| Old network - Bad:           | 0.136477, Good: 0.130365, Diff: +0.006112 |
| New network - Bad:           | 0.195371, Good: 0.101327, Diff: +0.094044 |
| Change in bad habits impact: | 0.087932                                  |
| Stroke                       |   |

## 8 Summary

The implementation correctly performs Bayesian network inference using variable elimination. Analysis of the health risk data revealed:

1. Factorized representation provides  $683\times$  compression over full joint distribution
2. Poor health indicators have stronger effects than bad habits (18.7% vs 1.1% for diabetes)
3. Higher income correlates with better health outcomes across all conditions
4. Smoking and exercise have direct effects on outcomes beyond bp/cholesterol/bmi
5. Diabetes directly increases stroke risk beyond shared risk factors

The independence assumptions in the original network are too strong and underestimate the true causal relationships in health risk factors.