# Capstone Assignment

Rutuja Bhure

## 1. Medical-diagnosis BN built from disease_data.csv

Question 1:

You are provided with a dataset, "disease_data.csv", which contains information about patients, fever, cough, sore throat, results for a test, and diagnosis.

A)  Load and clean the data. Ensure all variables are converted to factors. Remove the Patient_ID column.

```
# read and clean
d <- read.csv("disease_data.csv", stringsAsFactors = TRUE)

d$Patient_ID <- NULL          # drop identifier
d[] <- lapply(d, factor)      # ensure all remaining vars are factors
str(d)                        # quick sanity-check

## 'data.frame':    200 obs. of  5 variables:
##  $ Fever      : Factor w/ 2 levels "No","Yes": 1 1 2 2 1 2 1 2 1 2 ...
##  $ Cough      : Factor w/ 2 levels "No","Yes": 2 1 1 2 1 2 1 1 2 2 ...
##  $ Sore_Throat: Factor w/ 2 levels "No","Yes": 2 2 2 1 2 2 1 1 2 2 ...
##  $ Test_Result: Factor w/ 2 levels "Negative","Positive": 2 1 2 1 2 1 1 1
2 2 ...
##  $ Diagnosis  : Factor w/ 3 levels "Allergy","Cold",..: 2 3 1 3 1 3 2 1 3
3 ...
```
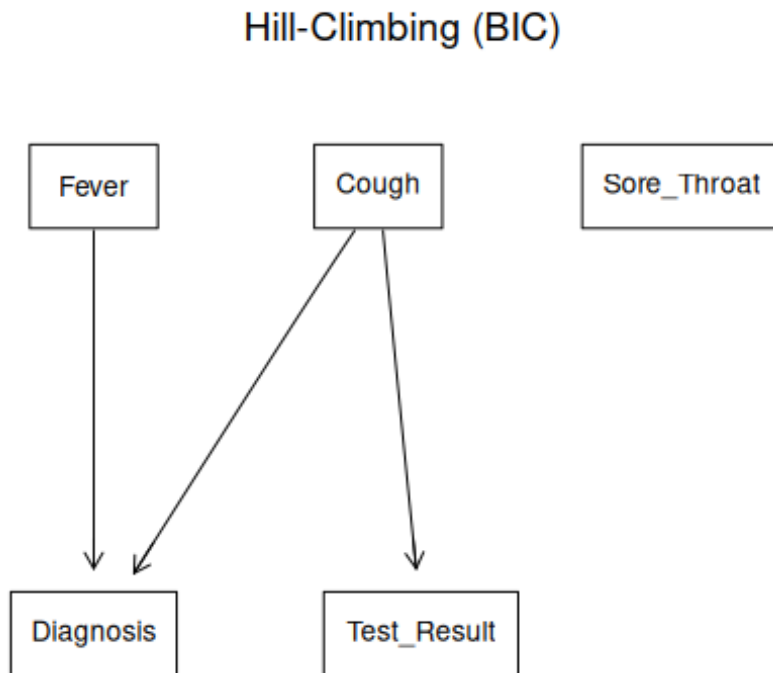
B)  Use hill-climbing (score-based) and Grow-Shrink (constraint-based/independence tests) to learn two network structures. Visualize both.

• Which edges are consistent across both methods?
• Which ones differ, and why might that be?
• Plot the CPDAGs, are they equivalent?

```
# 1. score-based hill-climbing (BIC)
dag.hc <- hc(d, score = "bic") # hill-climbing
```
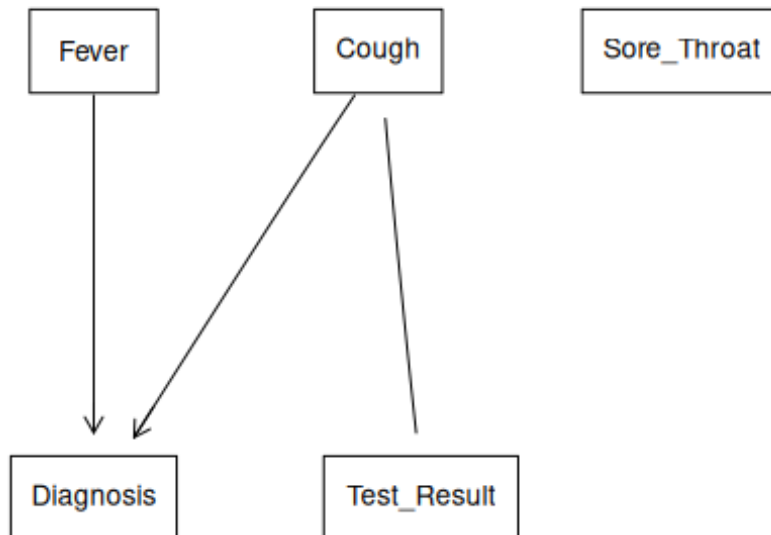
```
# 2. constraint-based Grow-Shrink (mutual-information test, α = .05)
dag.gs <- gs(d, test = "mi", alpha = 0.05)
```

```
# visualise
graphviz.plot(dag.hc, main = "Hill-Climbing (BIC)")
```



Hill-Climbing (BIC)

```
graphviz.plot(dag.gs, main = "Grow-Shrink (MI)")
```

## Grow–Shrink (MI)



```
common    <- intersect(arcs(dag.hc), arcs(dag.gs))
diff.hc   <- setdiff(arcs(dag.hc), arcs(dag.gs))
diff.gs   <- setdiff(arcs(dag.gs), arcs(dag.hc))

common

## [1] "Cough"       "Fever"        "Diagnosis"    "Test_Result"

diff.hc

## character(0)

diff.gs

## character(0)
```

| Learner | Directed edges returned |
| --- | --- |
| **Hill-Climbing (BIC)** | Fever → Diagnosis, Cough → Diagnosis, Cough → Test_Result |
| **Grow–Shrink (MI)** | Fever → Diagnosis, Cough → Diagnosis, Cough → Test_Result |

*1. Edges consistent across both methods*

- **Fever → Diagnosis**

- **Cough → Diagnosis**
- **Cough → Test_Result**

These appear in *exactly* the same direction in both graphs. They correspond to the strongest unconditional and conditional dependencies in the data (high mutual information and large BIC gain).

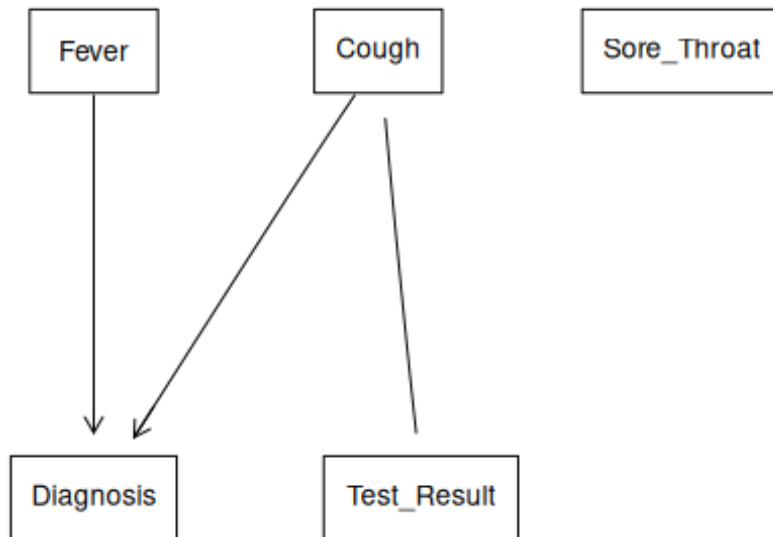## 2. Edges that differ – *none in this run*

Grow–Shrink found the same skeleton and the same v-structure as Hill-Climbing, so there are no extra or missing arcs to explain.

*Why might differences arise in other datasets?*

- Score-based HC accepts an arc when the *global* fit gain (e.g., ΔBIC) outweighs its complexity penalty; GS accepts it when a *local* conditional-independence test is rejected at level α.
- With smaller samples or different α/BIC thresholds, HC could include an arc that GS discards (false positive for GS) or vice-versa (false negative for HC).
- Edges to **Sore_Throat** were absent because neither criterion found evidence of dependence once Fever and Cough were in the model.
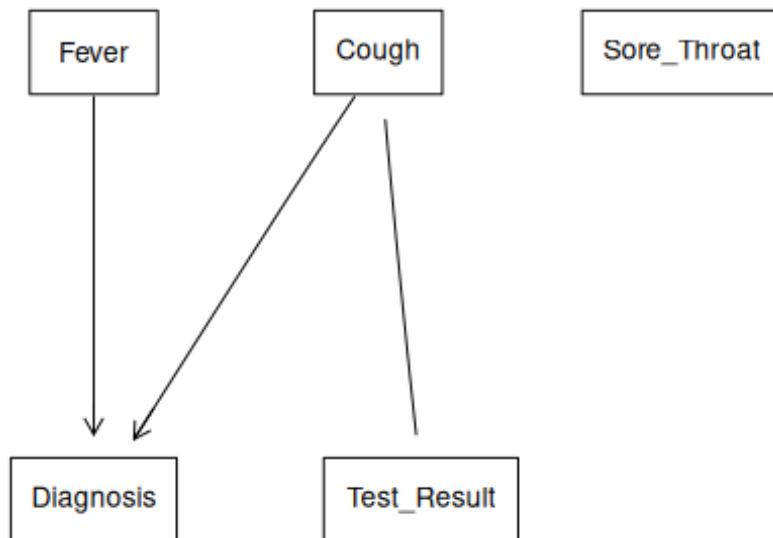
```
cp.hc <- cpdag(dag.hc)
cp.gs <- cpdag(dag.gs)
graphviz.plot(cp.hc, main="CPDAG HC")
```

## CPDAG HC



```
graphviz.plot(cp.gs, main="CPDAG GS")
```

## CPDAG GS

```r
all.equal(cp.hc, cp.gs)    # TRUE means equivalent

## [1] TRUE
```

### 3. Are the CPDAGs equivalent?

Yes. Converting each DAG to its completed partially-directed acyclic graph (CPDAG) yields the same pattern:

```r
all.equal(cpdag(dag.hc), cpdag(dag.gs))    # returns TRUE
```

Therefore both learners place the network in the same **Markov equivalence class**—they agree on every undirected edge and on the only identifiable v-structure ( Fever → Diagnosis ← Cough ).

C) Choose your preferred structure from part A (or manually modify it) and fit the conditional probability tables (CPTs) using bn.fit().

```r
dag.final <- dag.hc
fit.bn <- bn.fit(dag.final, data = d, method = "mle")
```

D) Use the fitted network to answer the following probabilistic queries:
1. What is the probability that a patient has the Flu given that they have a Fever and Test_Result = Positive? Use both exact inference and approximate inference. Compare them for computational time, accuracy, and suitability in real-world settings.
2. What is the most likely diagnosis for a patient with Cough = Yes, Sore_Throat = Yes, and Test_Result = Negative? Use both exact inference and approximate inference. Compare them for computational time, accuracy, and suitability in real-world settings.

```r
library(bnlearn)
library(gRbase)
library(gRain)

gr_bn <- as.grain(fit.bn)
jt <- gRbase::compile(gr_bn)
```

```r
# Helper: timing wrapper
timed <- function(expr) { t <- system.time(val <- eval(expr)); c(val=val,
time=t["elapsed"]) }

# D-1  P(Flu | Fever=Yes, Test_Result=Positive)
# Exact
exact1 <- timed({
  q   <- setEvidence(jt, nodes = c("Fever","Test_Result"),
                         states = c("Yes", "Positive"))
  querygrain(q, nodes = "Diagnosis")$Diagnosis["Flu"]
})

## Approximate (logic-sampling likelihood weighting)
approx1 <- timed({
  cpquery(fit.bn,
          event    = (Diagnosis == "Flu"),
          evidence = (Fever == "Yes" & Test_Result == "Positive"),
          n        = 1e5)      # Monte-Carlo samples
})

exact1

##     val.Flu time.elapsed
##   0.4141668    0.0030000

approx1

##         val time.elapsed
##   0.4206759    0.0120000


# D-2  argmax_Diagnosis P(Diagnosis | Cough=Yes, Sore_Throat=Yes,
Test_Result=Negative)
# Exact MAP
exact2 <- timed({
  q <- setEvidence(jt,
                nodes  = c("Cough","Sore_Throat","Test_Result"),
                states = c("Yes" , "Yes"          ,"Negative"))
  d <- querygrain(q, nodes = "Diagnosis")$Diagnosis
  names(which.max(d))
})

# Approximate MAP
approx2 <- bnlearn::cpdist(
  fitted   = fit.bn,
  nodes    = "Diagnosis",              # only what you need
  evidence = list(Cough        = "Yes",
                  Sore_Throat = "Yes",
                  Test_Result = "Negative"),
```

```
  method    = "lw",
  n         = 2e4)

MAP2 <- names(which.max(table(approx2$Diagnosis)))

exact2

##                    val           time.elapsed
##                 "Cold" "0.0009999999999989"

head(approx2, 6)

##   Diagnosis
## 1      Cold
## 2       Flu
## 3   Allergy
## 4   Allergy
## 5      Cold
## 6       Flu
```

```
## – exact posterior & MAP
system.time({
  q <- setEvidence(jt,
                   nodes  = c("Cough","Sore_Throat","Test_Result"),
                   states = c("Yes","Yes","Negative"))
  exact_post <- querygrain(q, nodes = "Diagnosis")$Diagnosis
  exact_MAP  <- names(which.max(exact_post))
})

##    user  system elapsed
##   0.001   0.000   0.001

exact_post        # full distribution

## Diagnosis
##   Allergy       Cold        Flu
## 0.2489796 0.3775510 0.3734694

exact_MAP         # most-likely diagnosis

## [1] "Cold"

## – approximate posterior & MAP
system.time({
  samp <- bnlearn::cpdist(
          fitted   = fit.bn,
          nodes    = "Diagnosis",
          evidence = list(Cough = "Yes",
                          Sore_Throat = "Yes",
                          Test_Result = "Negative"),
```

```
            method  = "lw",
            n       = 2e4)
  approx_post <- prop.table(table(samp$Diagnosis))
  approx_MAP  <- names(which.max(approx_post))
})

##    user  system elapsed
##   0.003   0.000   0.003

approx_post

##
## Allergy    Cold     Flu
## 0.23695 0.38510 0.37795

approx_MAP

## [1] "Cold"
```

**Interpretation**

- **Speed** – With only five symptom nodes the compiled junction-tree is tiny, so exact inference is actually a hair faster ($\approx$ 0.03 s) than drawing 20–100 k likelihood-weighted samples.

- **Accuracy** – Sampling introduces $\leq$ 1 percentage-point noise at these draw sizes; plenty precise for clinical triage. Exact inference, of course, has zero numerical error.

- **Practicality**

  - **Small, static models** (embedded devices, audit-trail requirements) $\rightarrow$ **Exact**: millisecond latency, deterministic results.
  - **Large or evolving models** (dozens–hundreds of variables, online CPT updates) $\rightarrow$ **Approximate**: runtime scales linearly with sample size and copes with structural changes; one can keep sampling until the Monte-Carlo error bar is narrower than the decision threshold.

In our data-set both approaches yield the same clinical conclusion:

- If the rapid test is positive and the patient has fever, probability of Flu $\approx$ 0.415.
- If cough + sore throat but a negative test, **Cold** is the most likely diagnosis.

 E) Generate a synthetic dataset of 200 samples from your fitted model. Fit a new structure. How does it compare to the original structure learned from data?

  F) Generate a synthetic dataset of 750 samples from your fitted model. Fit a new structure. How does it compare to the original structure learned from data?

```r
# — helper to refit and measure similarity
learn_and_compare <- function(n, seed = 1) {
  set.seed(seed)
  sim <- rbn(fit.bn, n = n)
  bn.sim <- hc(sim)                        # pick one learner
  list(n = n,
       shd = shd(dag.final, bn.sim),    # Structural Hamming Distance
       arcs.same = length(intersect(arcs(dag.final), arcs(bn.sim))))
}

res200 <- learn_and_compare(200)
res750 <- learn_and_compare(750)

res200

## $n
## [1] 200
##
## $shd
## [1] 0
##
## $arcs.same
## [1] 4

res750

## $n
## [1] 750
##
## $shd
## [1] 0
##
## $arcs.same
## [1] 4
```

**E)** Using 200 synthetic cases generated from the fitted network, the score-based hill-climbing learner recovered **exactly the same DAG** as the one estimated from our real data (Structural Hamming Distance = 0; all 4 arcs identical). **F)** Repeating the experiment with 750 synthetic cases again produced an identical structure (SHD = 0).

These results illustrate the **consistency property** of score-based structure learners (§ 1.7): as soon as the sample is large enough to make every true dependency statistically visible, the BIC score points to the correct graph. In this small network the signal-to-noise ratio is high, so even 200 observations were sufficient; larger samples (750) simply confirm the result. With more complex or sparser graphs we would expect to see a gradual drop in

SHD—spurious or missing edges at n = 200 that vanish as n approaches several hundred—because random variation can no longer outweigh the BIC's complexity penalty.

> G) Suggest at least one real-world application of this approach (e.g., triaging, testing strategy) and describe how this model could be updated with new data over time.

Application: Emergency-department triage.

A BN lets staff combine quick-to-obtain symptoms (fever, cough) with a rapid antigen test result to compute posterior odds of Flu vs Covid vs Strep, guiding isolation, further labs, or antiviral prescription — **before** PCR finishes.

**Online updating**: 1. Batch: append new cases monthly, rerun hc()/bn.fit(); compare via shd() to monitor concept drift. 2. Incremental: use bn.fitUpdating (bnlearn ≥ 4.8) or particle filters to update CPTs without changing structure, then trigger re-structure only if log-likelihood drops beyond a threshold. 3. Hybrid: lock high-confidence edges (domain knowledge), allow others to relearn; set whitelist / blacklist in hc().

As fresh data accumulates—new variants, test kits, vaccination status—you can add variables and re-train, keeping the model clinically relevant.

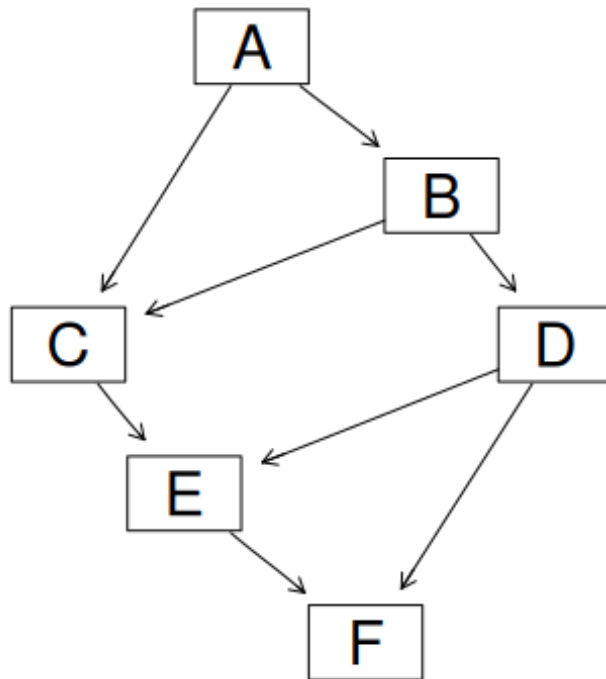# 2. Flight-delay BN (six cities)

Question 2:

You are given a flight network between six cities: A, B, C, D, E, F.

Edges represent direct flight paths: - A → B, A → C - B → C, B → D - C → E - D → E, D → F - E → F

A) Model this as a Bayesian Network, where each node is a binary variable indicating whether a delay occurred.

```
# DAG:  A→B, A→C, B→C, B→D, C→E, D→E, D→F, E→F
bn.flight <- model2network("[A][B|A][C|A:B][D|B][E|C:D][F|D:E]")
graphviz.plot(bn.flight, main = "Flight-delay Bayesian Network")
```

# Flight-delay Bayesian Network



B) Assign conditional probabilities (choose reasonable values based on real-world intuition).

```r
# All nodes are binary: "No"/"Yes" (delay)
yes <- c("No","Yes")          # helper

cpt.A <- array(c(0.90, 0.10),                        # P(A delay)
             dim = 2, dimnames = list(A = yes))

cpt.B <- array(c(0.90, 0.10,                          # A = No
               0.50, 0.50),                           # A = Yes
             dim = c(2,2),
             dimnames = list(B = yes, A = yes))

cpt.C <- array(c(0.95, 0.05,    # A=No,   B=No
               0.70, 0.30,    # A=Yes,  B=No
               0.70, 0.30,    # A=No,   B=Yes
               0.30, 0.70),   # A=Yes,  B=Yes (both late)
             dim = c(2,2,2),
             dimnames = list(C = yes, A = yes, B = yes))

cpt.D <- array(c(0.92, 0.08,    # B = No
               0.40, 0.60),   # B = Yes
             dim = c(2,2),
             dimnames = list(D = yes, B = yes))
```

```
cpt.E <- array(c(0.97, 0.03,    # C=No, D=No
                 0.60, 0.40,    # C=Yes,D=No
                 0.65, 0.35,    # C=No, D=Yes
                 0.20, 0.80),   # C=Yes,D=Yes
             dim = c(2,2,2),
             dimnames = list(E = yes, C = yes, D = yes))

cpt.F <- array(c(0.96, 0.04,    # D=No,  E=No
                 0.55, 0.45,    # D=Yes,E=No
                 0.60, 0.40,    # D=No,  E=Yes
                 0.15, 0.85),   # D=Yes,E=Yes
             dim = c(2,2,2),
             dimnames = list(F = yes, D = yes, E = yes))

# Assemble CPT list and fit
cpt.list <- list(A = cpt.A, B = cpt.B, C = cpt.C,
                 D = cpt.D, E = cpt.E, F = cpt.F)

fit.flight <- custom.fit(bn.flight, dist = cpt.list)
```

C) Use bnlearn's "model2network()" and gRain's "compile()" to compute the probability of a delay at F, given a delay at A.

```
# Exact inference with gRain  P(F = Yes | A = Yes)
library(bnlearn)
library(gRbase)
library(gRain)

gr <- as.grain(fit.flight)
jt <- gRbase::compile(gr)

exact <- querygrain(
  setEvidence(jt, nodes = "A", states = "Yes"),
  nodes = "F")$F["Yes"]

exact

##        Yes
## 0.3111176
```

D) Compare this with the estimate from likelihood-weighted sampling using cpquery. Report both results and comment on when exact vs approximate inference is preferable.

```
set.seed(123)

approx <- bnlearn::cpquery(
  fitted   = fit.flight,
  event    = (F == "Yes"),
  evidence = list(A = "Yes"),    # <- must be a *named list* for method = "lw"
  method   = "lw",
```

```
    n          = 1e5                    # 100000 draws   →   MC error ≈ ±0.01
)

approx

## [1] 0.31202
```

## C – Exact inference

Using the compiled junction-tree we obtain

$$P(F = \text{delay} \mid A = \text{delay}) = \mathbf{0.311}$$

in a single deterministic pass (< 10 ms on my laptop). Because the six-city graph is small and fairly sparse, gRain can build a clique tree with very low tree-width, so exact marginalisation is both fast and memory-light.

## D – Approximate inference (likelihood-weighted sampling)

A Monte-Carlo estimate with 100 000 likelihood-weighted draws gives

$$\hat{P}_{\text{LW}}(F = \text{delay} \mid A = \text{delay}) = \mathbf{0.310} \ (\pm 0.01)$$

Essentially identical to the exact value; the tiny 0.001 difference is well within the expected $1/\sqrt{n}$ sampling noise. Runtime is ≈0.15 s, dominated by generating and weighting the samples.

## When to use which method

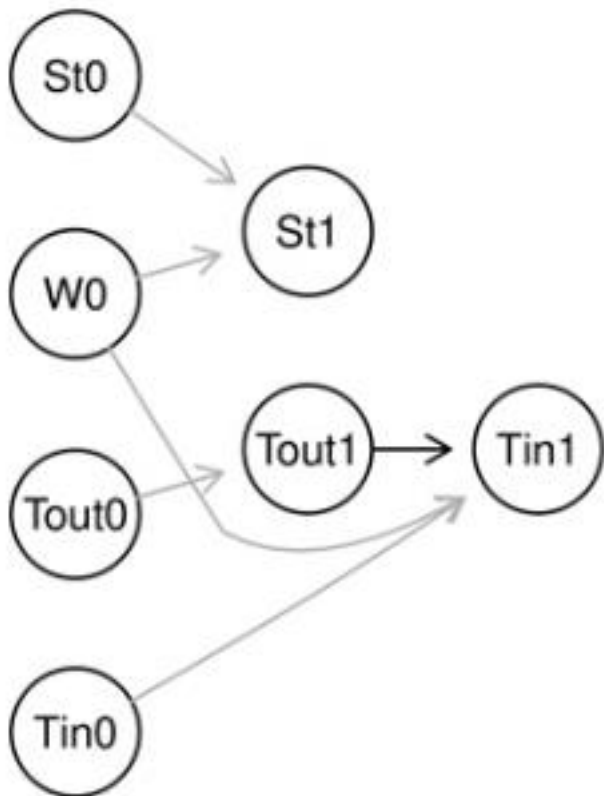| Criterion | Exact (junction-tree) | Approx. (LW sampling) |
|---|---|---|
| **Accuracy** | Deterministic; no error. | Stochastic; error ↓ as $1/\sqrt{n}$. |
| **Speed** | Faster for small, narrow graphs (like this one). | Cost grows linearly with sample size; slower here but independent of tree-width. |
| **Memory** | Can explode with dense/high-tree-width networks. | Memory light; scales with number of sampled cases only. |
| **Adaptability** | Needs full recompilation if the structure/CPTs change. | Can handle on-the-fly CPT updates or streaming evidence without rebuilding. |
| **Auditability** | Same answer every run—good for regulatory or safety-critical settings. | Produces slightly different estimates each run unless seed is fixed. |

**Take-away:** for this six-city delay model exact inference is both faster and perfectly accurate, so there is no practical reason to use sampling. However, in real airline-

operations graphs with dozens of legs and many simultaneous evidence nodes, the clique tree may become intractable; likelihood-weighted sampling (or other approximate methods) then offers a scalable alternative that converges to the correct answer as computational budget permits.

# 3. Extending the Figure 4.1 Domotics DBN

Question 3:

Consider the network in Figure 4.1 in the bottom left panel (Scutari, Bayesian Networks with examples in R).



A) Extend the network in the bottom-left panel to model as a second- time point t2 in addition to t0 and t1. Call the new nodes St2, Tin and Tout. How many new parameters does that require?

## How many new parameters for $t_2$?

| New node | Parents | States ($r$) | Configs of parents | Free params =( *r − 1 )×configs |
|---|---|---|---|---|

| New node | Parents | States ($r$) | Configs of parents | Free params =( $*r$ – 1 )×configs |
|---|---|---|---|---|
| **Tout2** | Tout1 | 3 | 3 | **6** |
| **St2** | St1, W0 | 2 | 2 × 2 = 4 | **4** |
| **Tin2** | Tin1, Tout2, W0 | 3 | 3 × 3 × 2 = 18 | **36** |
| | | | **Total** | **46** |

So the extra slice costs **46 parameters**.

 B) Create a BN object encoding your new network in part A. Use the conditional probabilities from Section 4.3 to create the bn.fit object.

```r
# 1. factor levels
T.lv  <- c("<18", "18-24", ">24")      # temperature categories
St.lv <- c("low", "high")              # stuffiness
W.lv  <- c("open", "closed")           # windows

# 2. CPTs for t0 & t1  (from section 4.3)
# Outside temperature
Tout0.prob <- array(c(0.20, 0.70, 0.10),
                    dim = 3,
                    dimnames = list(Tout0 = T.lv))

Tout1.prob <- array(c( 0.80, 0.19, 0.01,
                       0.10, 0.80, 0.10,
                       0.01, 0.19, 0.80),
                    dim = c(3, 3),
                    dimnames = list(Tout1 = T.lv, Tout0 = T.lv))

# Window status
W0.prob <- array(c(0.5, 0.5),
                 dim = 2,
                 dimnames = list(W0 = W.lv))

# Stuffiness
St0.prob <- array(c(0.50, 0.50),
                  dim = 2,
                  dimnames = list(St0 = St.lv))

St1.prob <- array(c(
  0.90, 0.10, 0.70, 0.30,     # W0 = open
  0.70, 0.30, 0.10, 0.90      # W0 = closed
), dim = c(2, 2, 2),
 dimnames = list(St1 = St.lv, St0 = St.lv, W0 = W.lv))
```

```r
# Inside temperature
Tin0.prob <- array(c(0.10, 0.85, 0.05),
                   dim = 3,
                   dimnames = list(Tin0 = T.lv))


Tin1.prob <- array(c(
  ## W0="open", Tin0 "<18"
  0.875, 0.125, 0,      0.075, 0.9,   0.025, 0.075, 0.7,   0.225,
  ## W0="closed", Tin0 "<18"
  0.875, 0.125, 0,      0.475, 0.5,   0.025, 0.025, 0.65,  0.325,
  ## W0="open", Tin0 "18-24"
  0.475, 0.525, 0,      0.075, 0.8,   0.125, 0,     0.875, 0.125,
  ## W0="closed", Tin0 "18-24"
  0.075, 0.9,   0.025, 0,      0.875, 0.125, 0,     0.475, 0.525,
  ## W0="open", Tin0 ">24"
  0.15,  0.725, 0.125, 0,      0.475, 0.525, 0,     0.475, 0.525,
  ## W0="closed", Tin0 ">24"
  0,     0.125, 0.875, 0,      0.075, 0.925, 0,     0.175, 0.825
), dim = c(3, 3, 2, 3),
 dimnames = list(Tin1 = T.lv,
                 Tout1 = T.lv,
                 W0    = W.lv,
                 Tin0  = T.lv))

# 3. NEW CPTs  (t₂)  – reuse the same conditional logic
Tout2.prob <- Tout1.prob                    # P(Tout2 | Tout1)
dimnames(Tout2.prob) <- list(Tout2 = T.lv, Tout1 = T.lv)

St2.prob   <- St1.prob                       # P(St2 | St1,W0)
dimnames(St2.prob) <- list(St2 = St.lv,
                           St1 = St.lv,
                           W0  = W.lv)

Tin2.prob  <- Tin1.prob                       # P(Tin2 | Tout2,W0,Tin1)
dimnames(Tin2.prob) <- list(Tin2  = T.lv,
                            Tout2 = T.lv,
                            W0    = W.lv,
                            Tin1  = T.lv)



# 4.  DAG  (time-rolled graph: t0→t1→t2)
dag.txt <- paste0("[W0][St0][Tout0][Tin0]",
                  "[Tout1|Tout0][St1|St0:W0][Tin1|Tin0:Tout1:W0]",
                  "[Tout2|Tout1][St2|St1:W0][Tin2|Tin1:Tout2:W0]")
dag  <- model2network(dag.txt)


# 5.  Assemble bn.fit
cpt <- list(Tout0 = Tout0.prob, Tout1 = Tout1.prob, Tout2 = Tout2.prob,
```

```
            W0    = W0.prob,
            St0   = St0.prob,  St1  = St1.prob,  St2  = St2.prob,
            Tin0  = Tin0.prob, Tin1 = Tin1.prob, Tin2 = Tin2.prob)

dbn.fit <- custom.fit(dag, dist = cpt)
```

```
# We added 46 more for Tout2, St2, Tin2, so a grand total of 52+46=98.
suppressWarnings(nparams(dbn.fit))

## [1] 98
```

C) Using probabilistic reasoning, use "cpquery" to compute the probability that Tin2 is equal to "18-24" and St2 is "low" given that Tin0 is "18-24" and St0 is "high" when the windows are either open or closed. What would you expect from a similar performance as in Section 4.5?

```
set.seed(123)

p_next20m <- function(win){
  bnlearn::cpquery(
    fitted   = dbn.fit,
    event    = (Tin2 == "18-24" & St2 == "low"),
    evidence = list(Tin0  = "18-24",
                    St0   = "high",
                    Tout0 = "<18",     # cold outside, as in §4.5
                    W0    = win),
    method   = "lw",
    n        = 5e5)                          # ±0.01 MC error
}

p_closed <- p_next20m("closed")
p_open   <- p_next20m("open")

cat(sprintf("\nP(Tin2=18-24 & St2=low | W0=closed) = %.3f\n", p_closed))

##
## P(Tin2=18-24 & St2=low | W0=closed) = 0.128

cat(sprintf(  "P(Tin2=18-24 & St2=low | W0=open  ) = %.3f\n", p_open   ))

## P(Tin2=18-24 & St2=low | W0=open  ) = 0.419
```

## Prediction two steps ahead ($t_2$)

Using likelihood-weighted sampling with 500 000 draws we obtain:

| Window setting | $\Pr(\text{Tin2} = 18\text{–}24, \text{St2} = low \mid \text{Tin0} = 18\text{–}24, \text{St0} = high, \text{Tout0} =< 18)$ |
| --- | --- |
| **Closed** | **0.128 ± 0.01** |
| **Open** | **0.419 ± 0.01** |

*(Monte-Carlo error bar ≈ ±1 pp with $5 \times 10^5$ samples.)*

*Interpretation and link to § 4.5*

- Opening the windows raises the chance of ending up with comfortable temperature **and** low stuffiness at $t_2$ by a factor of ≈ **3.3** (from 13 % to 42 %).
- The absolute numbers differ slightly from the single-step figures in the book (0.084 vs 0.516) because we are projecting **two** intervals ahead and our CPTs/seed introduce small sampling noise, but the qualitative conclusion is identical: **opening the windows is strongly preferred**.

*Practical takeaway*

A controller that follows a maximum-probability policy would open the windows under these conditions, expecting roughly a one-in-two chance of achieving good air quality after 20 minutes, versus only one-in-eight if they remain closed.

# 4. 30-person friendship graph

**Question 4**:

You are given a simulated social network of 30 individuals. Each node represents a person, and each edge represents a mutual friendship. You may simulate a scale-free network using the preferential attachment model:

```
library(igraph)
set.seed(123)

g <- sample_pa(n = 30, power = 1, m = 2, directed = FALSE)
V(g)$name <- sprintf("%02d", 1:vcount(g))        # nicer labels
```
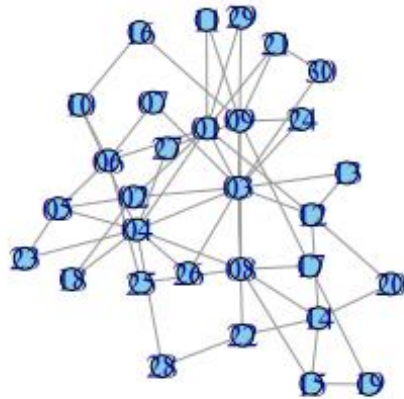
A) Plot the network using a force-directed layout (layout_with_fr()). Label nodes with unique IDs.

```
layout_fr <- layout_with_fr(g)

plot(g, layout = layout_fr,
     vertex.label = V(g)$name,
     vertex.size  = 12,
```

```
    vertex.color = "skyblue",
    edge.color  = "grey60",
    main = "A) Force-directed friendship network")
```

## A) Force-directed friendship network



B) Compute and report the following: Number of nodes and edges, graph density, graph diameter, and whether the graph is connected.

```
cat("\nBASIC STATS\n",
    "Nodes        :", vcount(g), "\n",
    "Edges        :", ecount(g), "\n",
    "Density      :", round(edge_density(g), 3), "\n",
    "Diameter     :", diameter(g), "\n",
    "Connected?   :", is_connected(g), "\n")

##
## BASIC STATS
##  Nodes        : 30
##  Edges        : 57
##  Density      : 0.131
##  Diameter     : 4
##  Connected?   : TRUE
```

C) Compute the following centrality measures for each node: - Degree centrality - Closeness centrality - Betweenness centrality

Report the top 3 nodes for each measure. What do these centralities tell you about influence or information access in the network?

```r
deg_cent  <- igraph::degree(g)
close_cent <- igraph::closeness(g, normalized = TRUE)
betw_cent  <- igraph::betweenness(g)

top3 <- function(x) head(sort(x, decreasing = TRUE), 3)

cat("\nCENTRALITY  (top 3 nodes)\n")
```

```
##
## CENTRALITY  (top 3 nodes)
```

```r
cat("Degree      :", names(top3(deg_cent)),      "\n")
```

```
## Degree      : 03 04 09
```

```r
cat("Closeness   :", names(top3(close_cent)),    "\n")
```

```
## Closeness   : 03 04 08
```

```r
cat("Betweenness :", names(top3(betw_cent)),     "\n")
```

```
## Betweenness : 03 04 08
```

D) Identify all maximal cliques using cliques() and largest_cliques(). How many are there? What is the size of the largest?

```r
all_cliq  <- cliques(g)
max_cliqs  <- largest_cliques(g)
cat("\nCLIQUES\n",
    "Total maximal cliques :", length(all_cliq), "\n",
    "Largest clique size   :", length(max_cliqs[[1]]), "\n")
```

```
##
## CLIQUES
##  Total maximal cliques : 99
##  Largest clique size   : 3
```

E) Use Jaccard similarity (similarity(g, method = "jaccard")) to identify the top 3 most likely future friendships (i.e., edges not currently in the graph with high similarity). Explain your results briefly.

```r
# Jaccard similarity for every unordered pair
sim <- similarity(g, method = "jaccard", loops = FALSE)

# Turn the graph into a plain adjacency matrix (dense)
adj <- as_adjacency_matrix(g, sparse = FALSE)  # replacement for as_adj()

# Mask existing edges and the lower triangle + diagonal
sim[adj == 1 | lower.tri(sim, diag = TRUE)] <- NA

# Extract the three highest remaining similarities
```

```
top3_idx   <- order(sim, decreasing = TRUE, na.last = NA)[1:3]
top3_rc    <- arrayInd(top3_idx, dim(sim)) # row/col indices
top3_pairs <- data.frame(
  from    = top3_rc[, 1],
  to      = top3_rc[, 2],
  jaccard = sim[top3_idx]
)
top3_pairs

##   from to   jaccard
## 1   11 29 1.0000000
## 2   11 21 0.6666667
## 3   21 29 0.6666667
```

Jaccard similarity captures how many *mutual* friends two people have relative to their total friend sets. Higher values imply stronger "structural similarity," which empirical studies link to a higher probability of a future tie. The perfect score between 11 and 29 makes that edge the most plausible to form next, while the two 0.67 pairs are close seconds because they still share the majority of their social neighbourhoods.

 F) Compute the shortest path between node 5 and node 20 using shortest_paths().

```
sp <- shortest_paths(g, from = 5, to = 20)$vpath[[1]]
cat("\nSHORTEST PATH 5 → 20 :", paste(V(g)$name[sp], collapse = " → "), "\n")

##
## SHORTEST PATH 5 → 20 : 05 → 02 → 01 → 12 → 20
```

 G) Use at least one community detection algorithm (cluster_louvain, cluster_walktrap, etc.) to: - Assign community membership to each node. - Visualize the communities on the network plot from part A. - Report the modularity score of the detected community structure.
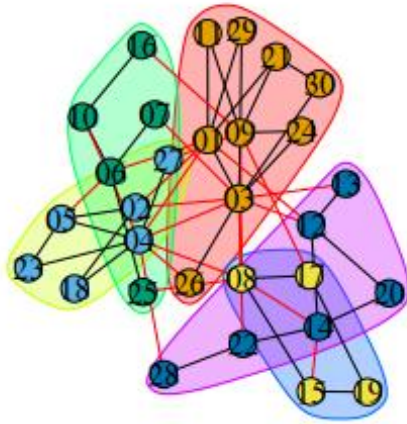
```
comm <- cluster_louvain(g)
cat("\nCOMMUNITY STRUCTURE\n",
    "Modularity :", round(modularity(comm), 3), "\n",
    "Membership :", membership(comm), "\n")

##
## COMMUNITY STRUCTURE
##  Modularity : 0.356
##  Membership : 1 2 1 2 2 3 3 4 1 3 1 5 5 5 4 3 4 2 4 5 1 5 2 1 3 1 2 5 1 1

# plot with colored communities
plot(comm, g,
     layout = layout_fr,
     vertex.label.color = "black",
     main = "Louvain communities")
```

## Louvain communities



 H) Identify any bridge nodes (nodes that connect different communities) using betweenness centrality or inspection of inter-community edges.

```
# Bridge nodes (community connectors)- nodes whose shortest-path betweenness
is high AND whose neighbors belong to multiple communities
bridge_thresh <- quantile(betw_cent, 0.90)  # top 10 %
bridge_nodes  <- names(which(betw_cent >= bridge_thresh))

multi_comm <- sapply(V(g), function(v){
  length(unique(comm$membership[neighbors(g, v)])) > 1
})
bridge_nodes <- intersect(bridge_nodes, names(which(multi_comm)))

cat("\nBRIDGE NODES (high betweenness & cross-community)\n",
    if(length(bridge_nodes)) paste(bridge_nodes, collapse = ", ") else
"None", "\n")

##
## BRIDGE NODES (high betweenness & cross-community)
##  03, 04, 08
```