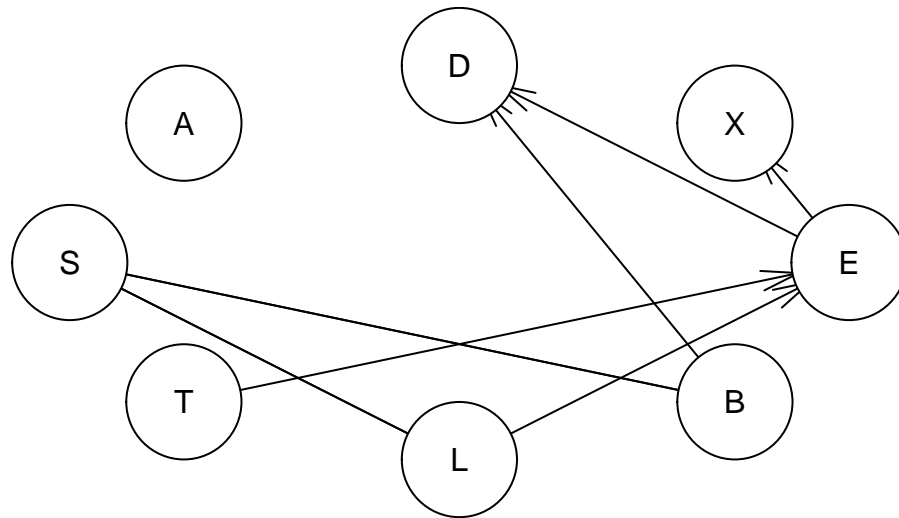


Lab_1

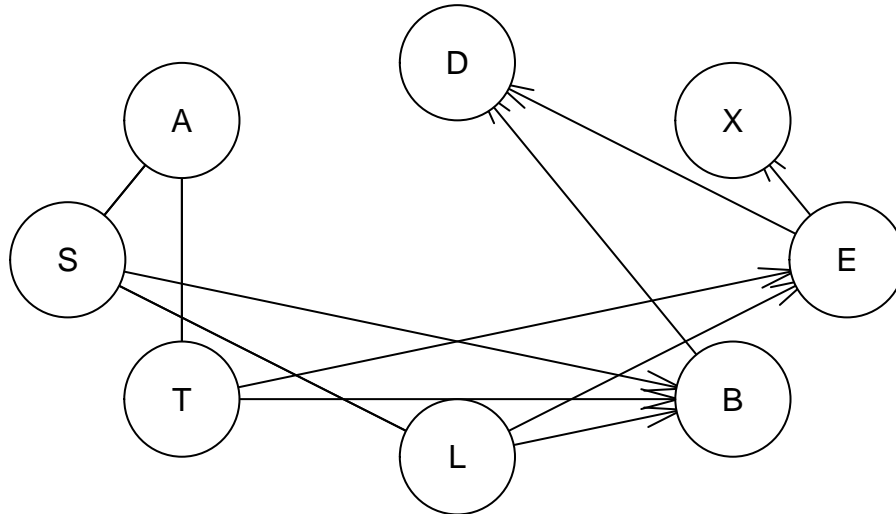
Axel Holmberg (axeho681)

Task 1

Show that multiple runs of the hill-climbing algorithm can return non-equivalent Bayesian network (BN) structures. Explain why this happens. Use the Asia dataset which is included in the bnlearn package. To load the data, run `data("asia")`.



Above is the junction tree for the data with 5 restarts.



Above is a junction tree for the data where the score is set by Akaike Information Criterion score(aic) instead of the default Bayesian Information Criterion.

```
## [1] "Different number of directed/undirected arcs"
```

As one can see in the print statement above, as well as in the plots, the hill-climbing algorithm can return non-equivalent BN-structures. The reason for this is mainly the score evaluation method as that changes how the hill-climbing method is made. In addition, the restarts indicates how many times the algorithm restarts and thereby can move away from potential “bad” local optimums, which also contributes to how the BN ends up at the end of the algorithm.

Task 2

Learn a BN from 80 % of the Asia dataset. The dataset is included in the bnlearn package. To load the data, run `data("asia")`. Learn both the structure and the parameters. Use any learning algorithm and settings that you consider appropriate. Use the BN learned to classify the remaining 20 % of the Asia dataset in two classes: $S = \text{yes}$ and $S = \text{no}$. In other words, compute the posterior probability distribution of S for each case and classify it in the most likely class. To do so, you have to use exact or approximate inference with the help of the bnlearn and gRain packages, i.e. you are not allowed to use functions such as `predict`. Report the confusion matrix, i.e. true/false positives/negatives. Compare your results with those of the true Asia BN.

```
##
## bn_pred  no yes
##      no  337 121
##      yes 176 366
```

Above is the confusion matrix based on BN learned by the training data.

```
##
```

```
## bn_pred_true  no yes
##              no  337 121
##              yes  176 366
```

Above is the confusion matrix based on BN from the true Asia BN.

Both of the confusion matrices are the same, which can be attributed to that both of the networks has the same connections to node “S”, so “S” is evaluated the same in both of the networks.

Task 3

In the previous exercise, you classified the variable S given observations for all the rest of the variables. Now, you are asked to classify S given observations only for the so-called Markov blanket of S , i.e. its parents plus its children plus the parents of its children minus S itself. Report again the confusion matrix.

```
##
## bn_pred_mb  no yes
##              no  337 121
##              yes  176 366
```

Above is the confusion matrix when S is classified only for the Markov Blanket of S .

```
##
## bn_pred_mb_true  no yes
##                  no  337 121
##                  yes  176 366
```

Above is the confusion matrix when S is classified only for the Markov Blanket of S for the true Asia BN.

Task 4

Repeat the exercise (2) using a naive Bayes classifier, i.e. the predictive variables are independent given the class variable. See p. 380 in Bishop’s book or Wikipedia for more information on the naive Bayes classifier. Model the naive Bayes classifier as a BN. You have to create the BN by hand, i.e. you are not allowed to use the function `naive.bayes` from the `bnlearn` package.

```
##
## bn_naive_pred  no yes
##                no  359 180
##                yes 154 307
```

Task 5

Explain why you obtain the same or different results in the exercises (2-4).

The reason for the confusion matrices being the same in task 2 and task 3 is because the dependence of “S” is the same in both cases, which is “L” and “B” (which is the Markow blanket for S). With knowledge about the state of L and B we have full information and therefore the result is the same when doing prediction for the whole network and for the Markow blanket. This applies as we complete data.

The confusion matrix is different in task 4 because of completely different dependency where all the nodes are dependent on “S” instead. Considering Bayes Theorem this affects the construction of the posterior distribution of S. We know the true structure of the BN and that the posterior should only consist of the distributions of L, B and the prior of S. We can then understand that the predictions are different.

Appendix for code

```
library(bnlearn)
library(Rgraphviz)

if (!requireNamespace("BiocManager", quietly = TRUE))
  install.packages("BiocManager")

BiocManager::install("Rgraphviz")
BiocManager::install("RBGL")
library(gRain)

data("asia")

data <- asia

##### TASK 1 #####
# Show that multiple runs of the hill-climbing algorithm can
# return non-equivalent Bayesian
# network (BN) structures. Explain why this happens

# Runs hill climbing algorithm with 5 restarts.
set.seed(12345)
bn_hc <- hc(x = data, restart = 5)
bn_hc <- cpdag(bn_hc)
plot(bn_hc)

set.seed(12345)
bn_hc_2 <- hc(x = data,
              score = "aic")
bn_hc_2 <- cpdag(bn_hc_2)
plot(bn_hc_2)

# Check equality
print(all.equal(bn_hc, bn_hc_2))

##### TASK 2 #####
# Learn a BN from 80% of the data set. Use BN to classify the rest of the data.
# Use exact or approximate inference with the help of the bnlearn and gRain
# packages. You are not allowed to use functions such as predict. Report the
# confusion matrix, i.e. true/false positives/negatives. Compare your results
# with those of the true Asia BN.

set.seed(12345)
n <- dim(data)[1]
id <- sample(1:n, floor(n * 0.8))
train <- data[id,]
test <- data[-id,]

bn_structure <- hc(train, restart = 5)
```

```

plot(bn_structure)

bn_fit <- bn.fit(bn_structure, train, method = "bayes")

# True Asia BN
bn_true <- model2network("[A] [S] [T|A] [L|S] [B|S] [D|B:E] [E|T:L] [X|E]")
bn_true_fit <- bn.fit(bn_structure, train, method = "bayes")
plot(bn_true)

#Converts the fit from above into a gRain-object
bn_grain <- as.grain(bn_fit)
bn_true_grain <- as.grain(bn_true_fit)

#Compiles the junction
j_tree <- compile(bn_grain)
j_tree_true <- compile(bn_true_grain)

prediction <- function(tree, data, obs, pred) {
  predictions <- c()

  for (i in 1:dim(data)[1]) {
    x_vals <- c()

    for (j in obs) {
      x_vals[j] <- if (data[i, j] == "yes")
        "yes"
      else
        "no"
    }

    evidence <- setEvidence(tree, obs, x_vals)
    probability <- querygrain(evidence, pred)$S["yes"]
    predictions[i] <- if (probability >= 0.5)
      "yes"
    else
      "no"
  }
  return(predictions)
}

bn_pred <-
  prediction(j_tree, test, c("A", "T", "L", "B", "E", "X", "D"), c("S"))
bn_pred_true <-
  prediction(j_tree_true, test, c("A", "T", "L", "B", "E", "X", "D"), c("S"))

# Confusion matrices by comparing with original data
confusion_matrix_fit <- table(bn_pred, test$S)

```

```

print(confusion_matrix_fit)
confusion_matrix_true <- table(bn_pred_true, test$S)
print(confusion_matrix_true)

# As one can see the confusion matrix is the exact same as the ones above.
# The reason for this is that the node is independent from

##### TASK 3 #####
# In the previous exercise, you classified the variable S given observations for all the
# rest of the variables. Now, you are asked to classify S given observations only for the
# so-called Markov blanket of S, i.e. its parents plus its children plus the parents of its
# children minus S itself. Report again the confusion matrix

bn_pred_mb <-
  prediction(j_tree, test, mb(bn_fit, "S"), c("S"))
bn_pred_mb_true <-
  prediction(j_tree_true, test, mb(bn_true_fit, "S"), c("S"))

confusion_matrix_fit_mb <- table(bn_pred_mb, test$S)
print(confusion_matrix_fit_mb)
confusion_matrix_true_mb <- table(bn_pred_mb_true, test$S)
print(confusion_matrix_true_mb)

##### TASK 4 #####
# Repeat the exercise (2) using a naive Bayes classifier, i.e. the predictive variables are
# independent given the class variable. See p. 380 in Bishop's book or Wikipedia for
# more information on the naive Bayes classifier. Model the naive Bayes classifier as a
# BN. You have to create the BN by hand, i.e. you are not allowed to use the function
# naive.bayes from the bnlearn package.

bn_naive <- empty.graph(c("A", "S", "T", "L", "B", "E", "X", "D"))
arc.set <- matrix(
  c("S", "A",
    "S", "T",
    "S", "L",
    "S", "B",
    "S", "E",
    "S", "X",
    "S", "D"),
  ncol = 2,
  byrow = TRUE,
  dimnames = list(NULL, c("from", "to")))
)

# Above is the BN set. In a naive bayes network all nodes are dependent on the independent classifier
# node. If S is observed, all other nodes are independent

arcs(bn_naive) <- arc.set

bn_naive_fit <- bn.fit(bn_naive, data=train)

bn_naive_grain <- as.grain(bn_naive_fit)

```

```
bn_naive_comp <- compile(bn_naive_grain)

bn_naive_pred <- prediction(bn_naive_comp, test, c("A", "T", "L", "B", "E", "X", "D"), c("S"))

confusion_matrix_fit_naive <- table(bn_naive_pred, test$S)
confusion_matrix_fit_naive
```