# 732A96: Lab 1
# Advanced Machine Learning

### Carles Sans Fuentes
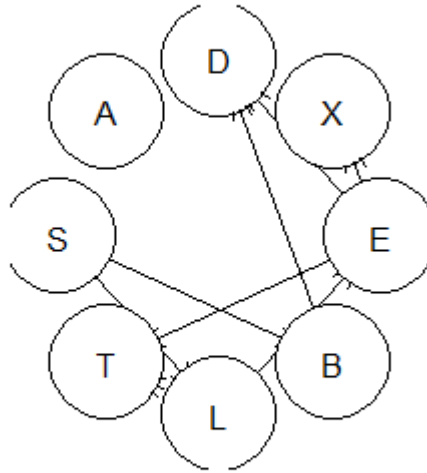
### October 9, 2017

---

## Assignment

### Question 1

In this question we are asked to show that multiple runs of the hill-climbing algorithm can return non-equivalent graphs. The hill-climbing algorithm is is a local greedy search optimization technique which iteratively starts with an arbitrary solution to a problem, and then it tries to find a better solution by incrementally changing a single element of the solution, such that if the change produces a better solution, then this new one is taken for the next incremental change until no further improvements are found. Intuitively, it can be understood that different ending outcomes can be found from different initial points if the function is complex enough. In order to prove this, the asia dataset from the bnlearn package in R has been taken. Two different starts has been tried, specifying in one an initial connection to exist and in the second one just totally random. It can be seen in the information and graphs below (table and 1 that the graphs are not equivalent.

```
 1  > mygraph
 2    Bayesian network learned via Score-based methods
 3    model:
 4      [partially directed graph]
 5    nodes:                                8
 6    arcs:                                 7
 7      undirected arcs:                    2
 8      directed arcs:                      5
 9    average markov blanket size:          2.25
10    average neighbourhood size:           1.75
11    average branching factor:             0.62
12    learning algorithm:                   Hill-Climbing
13    score:                                BIC (disc.)
14    penalization coefficient:             4.258597
15    tests used in the learning procedure: 77
16    optimized:                            TRUE
17  > cmygraph
18    Bayesian network learned via Score-based methods model:
19      [partially directed graph]
20    nodes:                                8
21    arcs:                                 8
22      undirected arcs:                    6
23      directed arcs:                      2
24    average markov blanket size:          2.25
25    average neighbourhood size:           2.00
26    average branching factor:             0.25
27    learning algorithm:                   Hill-Climbing
28    score:                                BIC (disc.)
29    penalization coefficient:             4.258597
30    tests used in the learning procedure: 82
31    optimized:                            TRUE
32  > all.equal(mygraph, cmygraph)
33  [1] "Different number of directed/undirected arcs"
```
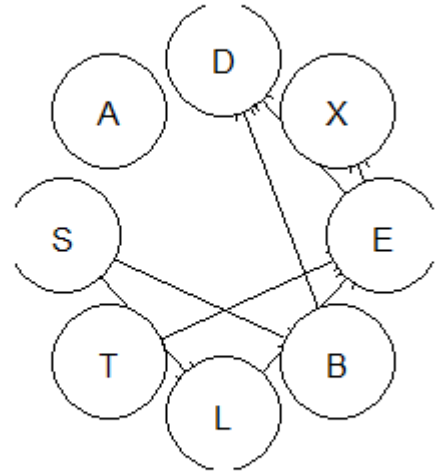
Figure 1: Graphical illustration of HC optimization for two different set ups

ALTERNATIVE SHOWING:

In order to show it in another way, I am going to use the data set alarm. I will be simulating from hill climbing as many times until it creates different graphs. And I will show this last ones graphically. This is what it can be seen in the following graph:
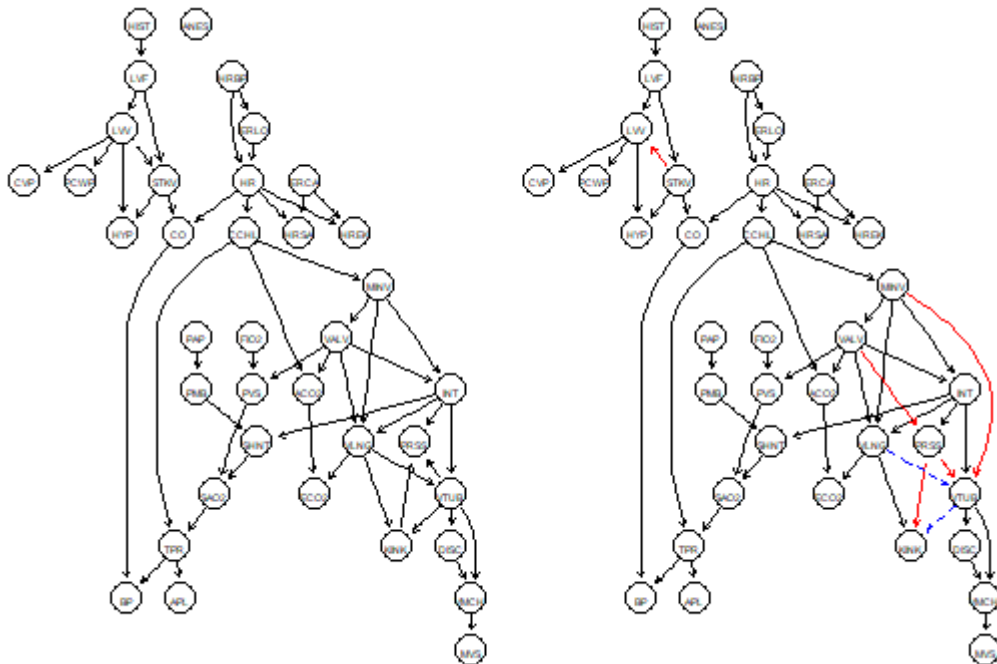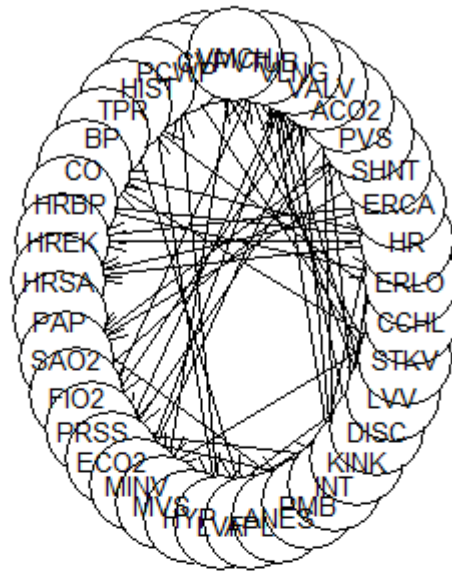


Figure 2: Graphical representation of the alarm data with different hc ending points
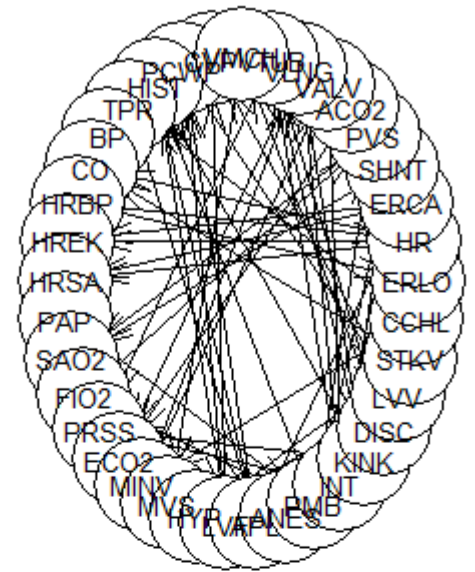
As explained previously, the hc() is a greedy algorithm which ends up in a local optima, so it tries to change arcs in order to get a better optima.
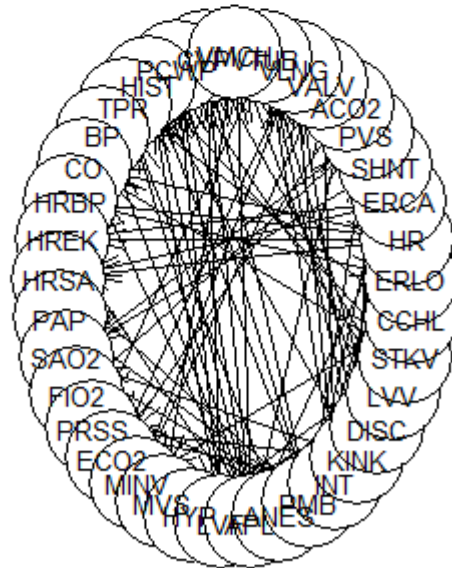
## Question 2

In figure 3 it can be found the graphical proof that increasing the sample size the BDeu decreases regularization.
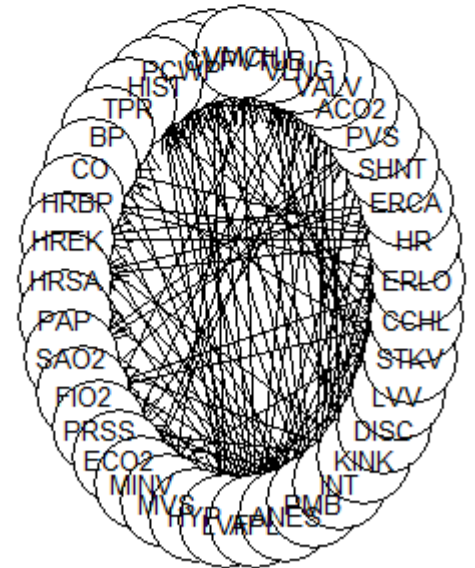
iss = 1

iss = 25

iss = 50

iss = 100

Figure 3: Graphical representation of the alarm data with different iss scores

The alpha.star results are the following ones:

```
1  [1]  6.779801 #for iss = 1
2  [1]  7.439901 #for iss = 25
3  [1]  10.34823 #for iss = 50
4  [1]  16.62416 #for iss = 100
```

The score got is lower because the more arc a graph has, the more likely to represent a graph from the data but the more complex it is so the more penalized the model is. Also, it is much probable that there will be a lot of parameters (e.g. arcs that do not exist in the real data) which will worsen the result.

## Question 3

2nd rewritten EXPLANATION

Now we are asked to compare the answers given some queries by exact and approximate inference algorithms.

The approximate algorithm uses samples from the distribution while it discards the samples that are not consistent with the conditioning set. Now we will iterate 4 times on both algorithm in order to compare how well they perform. For that, I have used the dataset called data.learning and I have compared the LS and the approximate algorithm in order to find the conditional distribution of A, E given $B = b$ and $F = b$. and then further below it can be found the conditional distribution of A, E given $B = b$ and $F = b$ , $D = b$ and $C = b$.

```
1  > res2
2  [[1]]
3             A A_Approx.Freq           E E_Approx.Freq
4  a 0.07394366    0.06010929 0.3167939    0.3333333
5  b 0.64964789    0.65391621 0.3664122    0.3679417
6  c 0.27640845    0.28597450 0.3167939    0.2987250
7
8  [[2]]
9             A A_Approx.Freq           E E_Approx.Freq
10 a 0.07394366    0.08791209 0.3167939    0.3150183
11 b 0.64964789    0.64285714 0.3664122    0.3791209
12 c 0.27640845    0.26923077 0.3167939    0.3058608
13
14 [[3]]
15            A A_Approx.Freq           E E_Approx.Freq
16 a 0.07394366        0.0752 0.3167939       0.2992
17 b 0.64964789        0.6352 0.3664122       0.3840
18 c 0.27640845        0.2896 0.3167939       0.3168
19
20 [[4]]
21            A A_Approx.Freq           E E_Approx.Freq
22 a 0.07394366    0.08448276 0.3167939    0.2948276
23 b 0.64964789    0.66724138 0.3664122    0.3862069
24 c 0.27640845    0.24827586 0.3167939    0.3189655
```

As previously mentioned, it can be found below the conditional distribution of A, E given $B = b$ and $F = b$ , $D = b$ and $C = b$

```
1  > res5
2  [[1]]
3             A A_Approx.Freq           E E_Approx.Freq
4  a 0.05115203    0.07142857 0.3167939    0.4285714
5  b 0.23753656    0.07142857 0.3664122    0.3214286
6  c 0.71131142    0.85714286 0.3167939    0.2500000
7
8  [[2]]
9             A A_Approx.Freq           E E_Approx.Freq
10 a 0.05115203    0.0000000 0.3167939    0.3529412
11 b 0.23753656    0.1764706 0.3664122    0.2352941
12 c 0.71131142    0.8235294 0.3167939    0.4117647
13
14 [[3]]
15            A A_Approx.Freq           E E_Approx.Freq
16 a 0.05115203    0.08333333 0.3167939    0.3333333
17 b 0.23753656    0.16666667 0.3664122    0.3333333
18 c 0.71131142    0.75000000 0.3167939    0.3333333
19
```

```
20 [[4]]
21           A A_Approx.Freq         E E_Approx.Freq
22 a 0.05115203        0.05 0.3167939          0.40
23 b 0.23753656        0.15 0.3664122          0.25
24 c 0.71131142        0.80 0.3167939          0.35
```

Answer: As a result, the LS algorithm gives the same results every time (being exact) whereas the approximate one gives every time different result.

Comparing both algorithms, the approximate one performs worse when the conditioning set is large because it is based on first obtaining a sample from the distribution and then discarding the samples that are not consistent with the conditioning set (the less samples the higher variance). So, when the conditioning set is large, most of the samples are discarded leading to a poor approximation of the exact result.

In other words, the variance of the output prediction is greater the more nodes are observed because there is less data to calculate it, leading to a worse prediction of the approximate model compared to the exact one. Nevertheless, the goodness of the model is higher since the prediction will be more accurate.

## Question 4

In this section We are asked to compute approximately the fraction of the 29281 DAGs that represent different independence models. In order to do so, I have used the random.graph function() sampling 10000 graphs. From these ones, I have check the unique ones and transformed them into DAGS. After that, I accounted for the the fraction of graphs resulted for different "burn.in" and "every" parameters getting the following results:

```
1 > result
2             burn in = 1 burn in = 100 burn in = 10000 burn in = 1e+06
3 every = 2     0.6178384     0.6073096       0.6229086        0.6252186
4 every = 20    0.5471494     0.5539439       0.5519067        0.5567743
5 every = 100   0.5531273     0.5551819       0.5510988        0.5536594
6 every = 200   0.5454000     0.5534984       0.5536512        0.5552473
```

Results seem to show that the true proportion of essential graphs is approximately 0.55-0.6 of the DAGS. In light of the results, it is preferable to perform structure learning in the space of essential graphs for computation simplicity given that the result must be quite well accurate. The "Burn in" in this case does not change the result. These might mean that the process is directly stationary from the first observations. On the other side, the "every" parameter affects more to the output. By doing each observation less correlated with the previous ones (e.g. increasing the "every" parameter) the existent correlation through data tends to a certain level (e.g. in our case with 20 it already converges). leading to a more accurate fraction of the number of unique elements graphs.

# Contributions

All results and comments presented have been developed and discussed together by the members of the group.

# Appendix

## Poisson regression-the MCMC way

```
1
2  ##Advanced Machine Learning lab 1
3
4
5  # source("http://bioconductor.org/biocLite.R")
6  # biocLite(c("graph", "RBGL", "Rgraphviz"))
7  #install.packages("gRain", dependencies=TRUE)
8  #install.packages("bnlearn")
9  #install.packages("gRbase")
10 #install.packages("gRain")
11 library(bnlearn)
12 library(gRain)
13 library(gRbase)
14
15 ###1
16 data(asia)
17
18 plot(asia)
19
20 par(mfrow=c(1,2))
21
22 wl = matrix(c("E", "T"), ncol = 2, byrow = TRUE,
23             dimnames = list(NULL, c("from", "to")))
24
25 plot(hc(asia, whitelist = wl), main ="HC with initial constraint")
26 plot(hc(asia, whitelist = NULL),main ="HC without constraint")
27
28 mygraph<-cpdag(hc(asia, whitelist = NULL, restart = 10))
29 cmygraph<-cpdag(hc(asia, whitelist = wl))
30 mygraph
31 cmygraph
32 all.equal(mygraph, cmygraph)
33
34 ################Alternative
35 set.seed(12345)
36 countinue <- TRUE
37 while(countinue){
38   hc1 <- hc(alarm, restart = 10)
39   hc2 <- hc(alarm, restart = 10)
40   continue <- ifelse(all.equal(vstructs(hc1), vstructs(hc2)) == TRUE, TRUE, FALSE)
41   if (continue !=TRUE){
42     par(mfrow = c(1,2))
43     graphviz.compare(hc1, hc2)
44     par(mfrow = c(1,1))
45     break
46   }
47 }
48
49 ###2
50 data("alarm")
51 n<-4
52 iss<- c(1,25,50,100)
53 mydag<-list()
54 par(mfrow=c(2,2))
55 for(i in 1:n){
56   mydag[[i]]<- hc(alarm, restart = n,
57                   score = "bde", iss = iss[i])
58   plot(mydag[[i]], sub = paste0("iss = ", iss[i] ))
59   print(alpha.star(mydag[[i]], alarm, debug = FALSE))
60 }
61
62
63 par(mfrow=c(1,1))
64
65 ####3
66 data(learning.test)
67 pdag = iamb(learning.test)
68 pdag
69 plot(pdag)
70 dag = set.arc(pdag, from = "B", to = "A")
71 dag = pdag2dag(pdag, ordering = c("A", "B", "C", "D", "E", "F"))
72 plot(dag)
73 LS<- as.grain(fit)## it creates LS
74 plot(LS)
75 fit = bn.fit(dag, learning.test)##It creates all conditional tables from one edge with another
       with his possible outcomes
76
77 MM<- compile(LS)##it triangulates and moralizes
78 plot(MM)
79 basic<-querygrain(MM, nodes = c("A","E"))
80
```

```r
81  ###################
82  ##################comparison of results for one change
83  ## I want the conditional distribution of A, E given B = b
84
85  ##Exact one
86  ChangeEvidence<-setEvidence(MM, c("B"), c("b"))
87  finalstate<-querygrain(ChangeEvidence, nodes = c("A","E"))
88
89  ##Approximate algorithm
90  rsample<-cpdist(fit, nodes = c("A", "E"), evidence= (B=="b"), method = "ls")
91  condprob<-lapply(rsample, FUN = function(x){table(x)/length(x)})
92
93  res1<-cbind(A_LS= as.data.frame(finalstate)[1], A_Greedy= condprob$A,
94        E_LS= as.data.frame(finalstate)[2], E_Greedy= condprob$E)
95  res1<-res1[,c(1,3,4,6)]
96  colnames(res1)<- c("A_LS", "A_Greedy","E_LS", "E_Greedy" )
97  res1
98
99  ##################comparison of results for two changes
100 ## I want the conditional distribution of A, E given B = b and F = b
101 ChangeEvidence2<-setEvidence(MM, c("B","F"), c("b", "b"))
102 finalstate2<-querygrain(ChangeEvidence2, nodes = c("A","E"))
103
104
105 res2<- list()
106 for(i in 1:4){
107   ##Exact one
108   ChangeEvidence2<-setEvidence(MM, c("B","F"), c("b", "b"))
109   finalstate2<-querygrain(ChangeEvidence2, nodes = c("A","E"))
110   ##Greedy algorithm
111   rsample2<-cpdist(fit, nodes = c("A", "E"), evidence= (B=="b")&(F=="b"), method = "ls")
112   condprob2<-lapply(rsample2, FUN = function(x){table(x)/length(x)})
113
114   res2[[i]]<-cbind(A_LS= as.data.frame(finalstate2)[1], A_Approx= condprob2$A,
115             E_LS= as.data.frame(finalstate2)[2], E_Approx= condprob2$E)
116   res2[[i]]<-res2[[i]][,c(1,3,4,6)]
117
118 }
119 res2
120
121
122 ###wITHOUT EVIDENCE
123 ## I want the conditional distribution of A, E
124
125 ##Exact one
126 finalstate3<-querygrain(MM, nodes = c("A","E"))
127
128
129 ##################comparison of results for two changes
130 ## I want the conditional distribution of A, E given B = b and F = b, D=d
131 ChangeEvidence4<-setEvidence(MM, c("B","F","D", "A"), c("b", "b","b", "b"))
132 finalstate4<-querygrain(ChangeEvidence4, nodes = c("A","E"))
133
134 i<-1
135 res5<- list()
136 for(i in 1:4){
137   ##Exact one
138   ChangeEvidence5<-setEvidence(MM, c("B","F","D", "C"), c("b", "b","b", "b"))
139   finalstate5<-querygrain(ChangeEvidence5, nodes = c("A","E"))
140   ##Greedy algorithm
141   rsample5<-cpdist(fit, nodes = c("A", "E"), evidence= (B=="b")&(F=="b")&(D=="b")&(C=="b"),
142         method = "ls")
142   condprob5<-lapply(rsample5, FUN = function(x){table(x)/length(x)})
143
144   res5[[i]]<-cbind(A_LS= data.frame(finalstate5)[1], A_Approx= data.frame(condprob5$A),
145                 E_LS= data.frame(finalstate5)[2], E_Approx= data.frame(condprob5$E))
146   res5[[i]]<-res5[[i]][,c(1,3,4,6)]
147
148 }
149 res5
150
151
152
153 ###4
154 burn_in<- c(1,100, 10000, 1000000)
155 every<-c(2,20,100, 200)
156
157 nodes<- LETTERS[1:5]
158 num <- 1000
159
160
161 checkingDags<- function(burnin, every, num, nodes){
162   mymat<- matrix(ncol =length(burnin), nrow=length(every))
163   for(i in 1:length(burnin)){
164     for(j in 1:length(every)){
165       z<-random.graph(nodes = nodes, num = num,  method = "ic-dag",
166                     every = every[j], burn.in =burnin[i])
```

```r
167
168         m<-unique(z)
169         print(length(m))# This gives the proportion the number of repeated graphs out my
170
171         dagGraph<-lapply(m, FUN= function(x){cpdag(x)})
172         mymat[j,i]<- length(unique(dagGraph))/length(m)
173         }
174   }
175   return(mymat)
176 }
177
178 result<-checkingDags(burnin = burn_in, every = every, num = num, nodes = nodes)
179 colnames(result)<- paste0("burn in = " , burn_in)
180 rownames(result)<- paste0("every = " , every)
181 result
```