



UNIVERSITÄT REGENSBURG

**Institute of Genomics & Practical Bioinformatics**

Master of Science Computational Science

# The max-min-hill-climbing algorithm an implementation in Rcpp

Report

in Practical Bioinformatics

by

Michael Bauer

Matrikelnummer: 152 8558

**Tutor:** Dr. Giusi Moffa

**Adviser:** Prof. Dr. Rainer Spang

**Date:** August 26, 2014

# Contents

List of Figures	3
1 Introduction	5
2 Background	7
3 How MMHC works	14
3.1 Introduction to an example . . . . .	14
3.2 General workflow of <i>MMHC</i> . . . . .	15
4 The max-min-hill-climbing algorithm	25
4.1 max-min parents and children (MMPC)) . . . . .	26
4.2 The scoring function . . . . .	28
4.3 Computational optimization . . . . .	28
5 A comparison to the bnlearn package	31
6 Conclusion	36
Bibliography	37

# List of Figures

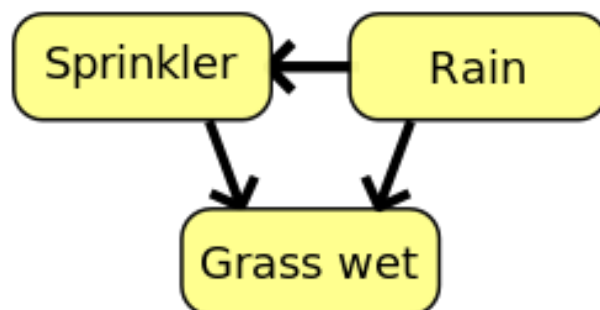
1.1	source: <a href="http://en.wikipedia.org/wiki/Bayesian_network#mediaviewer/File:SimpleBayesNetNodes.svg">http://en.wikipedia.org/wiki/Bayesian_network#mediaviewer/File:SimpleBayesNetNodes.svg</a> . . . . .	5
3.1	The graph of the underlying example. . . . .	15
3.2	First iteration step of the <i>MMPC</i> algorithm. . . . .	16
3.3	The first variable joined <i>PC</i> . . . . .	17
3.4	The second variable joined <i>PC</i> . . . . .	18
3.5	The third variable joined <i>PC</i> . . . . .	19
3.6	The skeleton of the graph. . . . .	20
3.7	The empty graph. . . . .	21
3.8	Add edges. . . . .	22
3.9	Reverse edges. . . . .	23
3.10	Delete edges. . . . .	24
4.1	source: The max-min hill-climbing Bayesian network structure learning algorithm, page 14. . . . .	25
4.2	source: The max-min hill-climbing Bayesian network structure learning algorithm, page 12. . . . .	26
4.3	source: The max-min hill-climbing Bayesian network structure learning algorithm, page 8. . . . .	27
4.4	Resulting graph after running the algorithm. . . . .	30
5.1	Comparing our implementation with bnlearn for <i>MMPC</i> . . . . .	33
5.2	Comparing our implementation with bnlearn for <i>MMPC</i> . . . . .	35

# Abstract

In this report we present an implementation of the max-min-hill-climbing algorithm (MMHC), first stated in [TBA], for R. It combines both: greedy search and constraint-based learning techniques. We will discuss those two methods separately. The main goal of this algorithm is to reconstruct Bayesian networks (BN) from estimated data. Bayesian networks play a great role in science, economics, sports and many more fields. Reconstructing them with the help of data is ongoing research and that's why there already exists a package in R which provides this algorithm. Our purpose was to write code which runs faster than the existing one. We used RCPP (C++ interface for R) for our implementation. This report tells you how our algorithm works and if we succeeded with our aim.

# 1 Introduction

The max-min-hill-climbing algorithm is one of the state of the art algorithms in statistical computing. The primal aim of this algorithm is reconstructing BN from estimated data. A Bayesian network is a Directed Acyclic Graph (DAG) whose nodes are random variables and edges represent conditional dependencies. If two random variables are connected they are said to be dependent. If there is no connection they are said to be conditional independent. BNs are more important than one can imagine. They play a great role in everdays life. For example [NBBCW] used them to predict the effect of missense mutations on the protein function. But not only medical scientists used Bayesian networks. Another example we found was football. In [PKA] they refer to an article where european football clubs tried to predict injuries of their players depending on BNs. With a simple Google search you may also find the prediction of stock exchanges and many more. Wikipedia provides a simple but descriptive example which illustrates BNs in a smooth way.



**Figure 1.1:** A simple example of a Bayesian network where Sprinkler & Rain, Sprinkler & Grass Wet and Grass Wet & Rain are conditionally dependent.

Reconstructing those networks is not easy, more precisely it is a np-hard problem. It is not only the amount of data which leads to a bad running time, also the dependencies between nodes can slow the code down.

Once we observe the Bayesian network from our data we then can look at the running time of the algorithm. But more important for us was to beat the existing algorithm for R (part of the "bnlearn" package). Our aim was to be faster. After an introduction to the mathematical

basis and our implementation we will see it it was possible to optimize our code with RCPP to beat the "bnlearn" algorithm.

## 2 Background

Before we are able to analyze our implementation and talk about it in detail we need some mathematical background. The notations, definitions, lemmas and theorems (given without proof) in this section are quoted verbatim from [TBA].

### Notation

We introduce our notation which is completely consistent to [TBA].

We denote

1. variables with an upper-case letter (e.g.,  $A, V_i$ ),
2. a state or a value of that variable by the same lower-case letter (e.g.,  $a, v_i$ ),
3. a set of variables by upper-case bold face (e.g.,  $\mathbf{Z}, \mathbf{Pa}_i$ ),
4. an assignment of state or value to each variable in the given set with the corresponding lower-case bold-face letter (e.g.,  $\mathbf{y}, \mathbf{pa}_i$ ),
5. special sets of variables (e.g. the set of all variables  $\mathcal{V}$ ) with calligraphic fonts.

### Definition (conditional independence)

Two variables  $X$  and  $Y$  are conditionally independent given  $\mathbf{Z}$  with respect to a probability distribution  $P$ , denoted as  $Ind_P(X; Y | \mathbf{Z})$ , if for all  $x, y, \mathbf{z}$  where  $P(\mathbf{Z} = \mathbf{z}) > 0$ ,

$$P(X = x, Y = y | \mathbf{Z} = \mathbf{z}) = P(X = x | \mathbf{Z} = \mathbf{z})P(Y = y | \mathbf{Z} = \mathbf{z}) \quad (2.1)$$

or

$$P(X, Y | \mathbf{Z}) = P(X | \mathbf{Z})P(Y | \mathbf{Z}) \quad (2.2)$$

for short. If  $X, Y$  are dependent given  $\mathbf{Z}$  we denote  $Dep_P(X; Y | \mathbf{Z})$ .

**Definition (Bayesian network)**

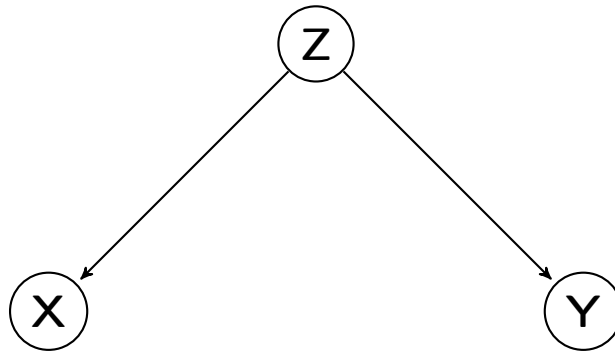
Let  $P$  be a discrete joint probability distribution of the random variables in some set  $\mathcal{V}$  and  $\mathcal{G} = \langle \mathcal{V}, E \rangle$  be a Directed Acyclic Graph (DAG). We call  $\langle \mathcal{G}, P \rangle$  a (discrete) *Bayesian network* if  $\langle \mathcal{G}, P \rangle$  satisfies the Markov condition.

**Definition (Markov condition)**

Any node in a Bayesian network is conditionally independent of its non-descendants, given its parents.

**Explanation**

With those definitions we have our first concept we will discuss briefly. We will explain the definitions by using the following picture:



The Markov condition states that  $X$  and  $Y$  given  $Z$  must be conditionally independent. This comes from the fact that  $X$  is a non-descendant of  $Y$  and vice versa and  $Z$  is a parent of both nodes. By fulfilling this condition we get from section 2 that this graph is a Bayesian network and with section 2 we have:  $P(X, Y|Z) = P(X|Z)P(Y|Z)$ .

**Defintion (collider)**

A node  $W$  of a path  $p$  is a *collider* if  $p$  contains two incoming edges into  $W$ .



**Definition (blocked path)**

A path  $p$  from node  $X$  to node  $Y$  is *blocked* by a set of nodes  $\mathbf{Z}$ , if there is a node  $W$  on  $p$  for which one of the following two conditions hold:

1.  $W$  is not a collider and  $W \in \mathbf{Z}$ , or
2.  $W$  is a collider and neither  $W$  or its descendants are in  $\mathbf{Z}$  [P88]

**Definition (d-seperated)**

Two nodes  $X$  and  $Y$  are *d-seperated* by  $\mathbf{Z}$  in graph  $\mathcal{G}$  (denoted as  $Dsep_{\mathcal{G}}(X;Y|\mathbf{Z})$ ) if and only if every path from  $X$  to  $Y$  is blocked by  $\mathbf{Z}$ . Two nodes are *d-connected* if they are not *d-seperated*.

**Definition (faithful)**

If all and only the conditional independencies true in the distribution  $P$  are entailed by the Markov condition applied to  $\mathcal{G}$ , we will say that  $P$  and  $\mathcal{G}$  are *faithful to each other* ([SGSN]). Furthermore, a distribution  $P$  is *faithful* if there exists a graph  $\mathcal{G}$ , to which it is faithful.

**Definition (faithfulness condition)**

A Bayesian network  $\langle \mathcal{G}, P \rangle$  satisfies the *faithfulness condition* if  $P$  embodies only independencies that can be represented in the DAG  $\mathcal{G}$  ([SGSN]). We will call such a Bayesian network a *faithful network*.

**Theorem**

In a faithful Bayesian network  $\langle \mathcal{G}, P \rangle$  the following equivalence holds ([P88])

$$Dsep_{\mathcal{G}}(X;Y|\mathbf{Z}) \iff Ind_P(X;Y|\mathbf{Z}) \quad (2.3)$$

## Remark and Explanation

**Remark:** For the rest of this report we assume faithfulness of the networks to learn. For this reason we do not explain the corresponding definitions in detail. Just note, that they are necessary for mathematical correctness.

**Explanation:** The definition of a collider already tells everything about it. To illustrate a collider, we have:



Here  $Z$  is a *collider* because it has two incoming edges. In this case if we just look for  $P(X; Y | \{\})$ , the path between  $X$  and  $Y$  would be blocked and for this  $X$  and  $Y$  are *d-separated*. If we look for  $P(X; Y | Z)$  with collider  $Z$  and obviously  $z \in Z, \forall z \in Z$ , then this path is not blocked and we state that  $X$  and  $Y$  are *d-connected*.

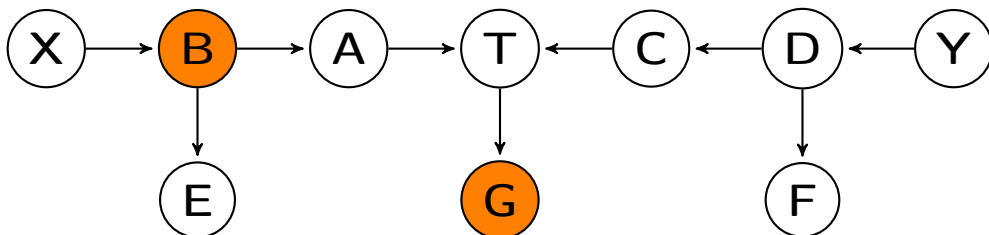
Because of section 2 and the faithfulness assumptions we state for the rest of our report that the terms d-separation and conditional independence are equivalent. With this we already know that  $X$  and  $Y$  are conditional dependent given  $Z$  in the example above. We want to give you two other examples for d-separation of variables.



If we are looking for  $Ind_P(X; Y | Z)$  with  $Z = \{B, D\}$  we learn that  $X$  and  $Y$  are conditionally independent given  $Z$ . In other words they are d-separated in the path because of the following reasons:

- $T$  is a collider with  $T \notin Z$  and blocks the path between  $X$  and  $Y$ .
- The nodes  $B$  and  $D$  are no colliders but they are elements of  $Z$ .

The situation becomes a bit more difficult if we take a look at the next example:



We learn that the path between  $X$  and  $Y$  remains blocked by looking for  $Ind_P(X; Y | \mathbf{Z})$  with  $\mathbf{Z} = \{B, G\}$ , i.e.  $X$  and  $Y$  are conditionally independent given  $\mathbf{Z}$ . If we would look at the path between  $A$  and  $Y$  we would learn that  $A$  and  $Y$  are d-connected. This comes from:

- $T$  is a collider but its descendant  $G \in \mathbf{Z}$ , i.e.  $T$  would not block the path.
- The node  $B$  is no collider but it is an element of  $\mathbf{Z}$ . For that it blocks the path.

So there is no element which blocks the path between  $A$  and  $Y$ , but for  $X$  and  $Y$ ,  $B$  blocks the path.

As we could see, detecting conditional independence of two nodes is quite difficult in small graphs. Since we normally observe large data sets with a couple of nodes, a concept for this is needed. As we will see, statistical methods, such as hypothesis testing will be a useful friend for this task.

### Definition

We define the minimum association of  $X$  and  $T$  relative to a feature subset  $\mathbf{Z}$ , denoted as  $MinAssoc(X; T | \mathbf{Z})$ , as

$$MinAssoc(X; T | \mathbf{Z}) = \min_{\mathbf{S} \subseteq \mathbf{Z}} Assoc(X; T | \mathbf{S}) \quad (2.4)$$

i.e., as the minimum association achieved between  $X$  and  $T$  over all subsets of  $\mathbf{Z}$ .

#### Remark:

1. In Figure 4.1 we will see this  $MinAssoc$  statement again as a function. With that we measure conditional independence of  $X$  and  $T$  given  $\mathbf{Z}$ .
2. For that the following equivalence holds (providing without a proof):

$$Ind(X; T | \mathbf{Z}) \iff (Assoc(X; T | \mathbf{Z}) = 0). \quad (2.5)$$

### Lemma ( $G^2$ value)

Because this algorithm is based on conditional independence testing, i.o.w. measuring the strength of association between two variables, we need some formulas for implementation to get this measure. We followed [SGSN] and calculated the  $G^2$  statistic, under the null hypothesis of the conditional independence holding. For that we have:

## 2 Background

Let  $S_{ijk}^{abc}$  be the number of times in the data where  $X_i = a$ ,  $X_j = b$  and  $X_k = c$ . We define in a similar fashion,  $S_{ik}^{ac}$ ,  $S_{jk}^{bc}$  and  $S_k^c$ , then the  $G^2$  statistic is defined as (c.f. [SGSN]):

$$G^2 := 2 * \sum_{a,b,c} S_{ijk}^{abc} * \ln \left( \frac{S_{ijk}^{abc} * S_k^c}{S_{ik}^{ac} * S_{jk}^{bc}} \right), \quad (2.6)$$

The  $G^2$  statistic is asymptotically distributed as  $\chi^2$  with appropriate degrees of freedom. To compute these degrees of freedom we use:

$$df = (|D(X_i)| - 1)(|D(X_j)| - 1) \prod_{X_l \in \mathbf{X}_k} |D(X_l)|, \quad (2.7)$$

where  $D(X_i)$  is the domain (number of distinct values) of variable  $X_i$ . After bringing this to code, we can work with the output of this function.

### Remark

section 2 is the basis of our calculations. To know whether two variables  $X$  and  $Y$  given a set  $Z$  are conditionally independent we run this statistical test and seek if it is likely that they are independent or not. In the next chapter you get an introduction on how to we realized this. To complete this chapter we give you the definition of the Bayesian Dirichlet likelihood-equivalence uniform (BDeu) score. It is also a statistical test which later helps us to set the edges in our graph such that we observe the graph which is most likely to be the right graph depending on the data. For that we quoted verbatim from [Ca06].

### The BDeu score

The Bayesian Dirichlet likelihood-equivalence uniform score is defined as:

$$g_{BDeu}(D, G) = \sum_{i=1}^n \left[ \sum_{j=1}^{q_i} \left[ \log \left( \frac{\Gamma(\frac{\eta}{q_i})}{\Gamma(N_{ij} + \frac{\eta}{q_i})} \right) + \sum_{k=1}^{r_i} \log \left( \frac{\Gamma(N_{ijk} + \frac{\eta}{r_i q_i})}{\Gamma(\frac{\eta}{r_i q_i})} \right) \right] \right], \quad (2.8)$$

where  $G$  is a graph,  $D$  is the underlying data, the number of states of the variable  $X_i$  is  $r_i$ , the number of possible configurations of the parent set  $Pa_G(X_i)$  of  $X_i$  is  $q_i$ , with  $q_i = \prod_{X_j \in Pa_G(X_i)} r_j$ ,  $w_{ij}, j = 1, \dots, q_i$  represents a configuration of  $Pa_G(X_i)$ ,  $N_{ijk}$  is the number of instances in the data set  $D$  where the variable  $X_i$  takes the value  $x_{ik}$  and the set of variables  $Pa_G(X_i)$  take the value  $w_{ij}$ ,  $N_{ij}$  is the number of instances in the data set where the variable

in  $Pa_G(X_i)$  take their  $j - th$  configuration  $w_{ij}$ , with  $N_{ij} = \sum_{k=1}^{r_i} N_{ijk}$ ,  $N_{ik}$  is the number of instances in  $D$  where the variable  $X_i$  takes its  $k - th$  value  $x_{ik}$  and therefor  $N_{ik} = \sum_{j=1}^{q_i} N_{ijk}$  and the total number of instances in  $D$  is  $n$ .

### Explanation

This formula takes as an input a graph and the observational data, it returns a score. Later, we are seeking for the graph with the highest score. Back to the formula, to compute the sums you just have to know how to get all those values explained above. Here is a short explanation of those variables in the formula:

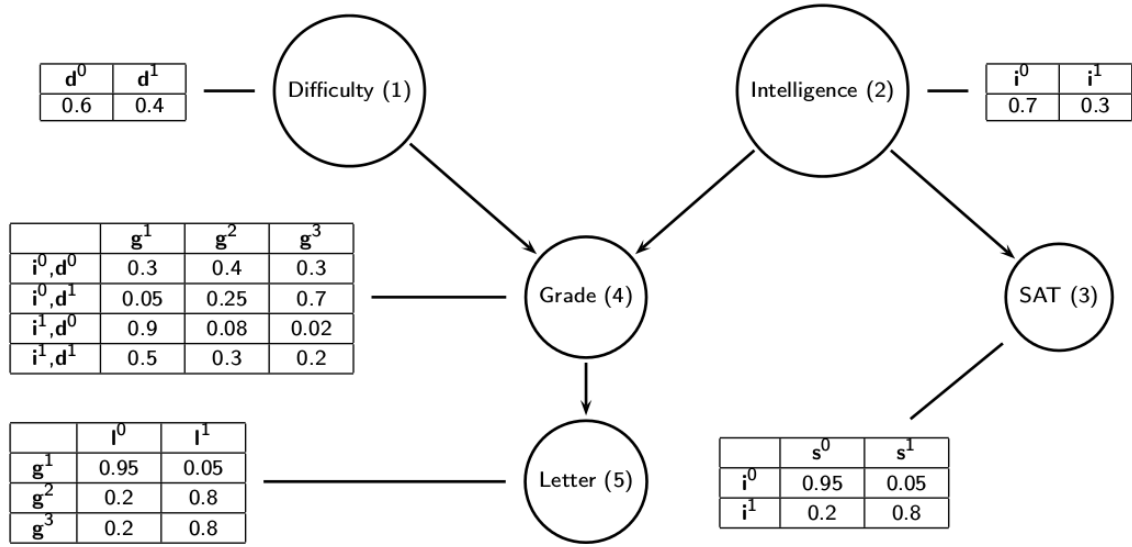
- **eta**: We set this value to 1 - the most common setting. There is no rule which value is the best and one could find many papers which try to find the "best" eta. If you are interested in that we recommend you [1]
- **n**: The number of variables in our BN.
- **$r_i$** : How many states the  $i - th$  variable can take.
- **$q_i$** : The product of all states of the parent's variables of the  $i - th$  variable, e.g. if node  $X$  has two parents  $Y$  and  $Z$   $q_X = r_Y \cdot r_Z$ .
- **$w_{ij}$** : This is a configuration of the parents of a certain variable, i.e. recall  $X$  with parents  $Y, Z$  and we assume that the states of  $Y, Z$  are binary. Then  $w_{ij}$  can have four different configurations:  $(0;0), (0;1), (1;0), (1;1)$ .
- **$N_{ijk}$** : Counts how often a variable  $X_i$  takes the value  $x_{ik}$  and its parents take the configuration  $w_{ij}$ .
- **$N_{ij}, N_{ik}$** : Is the calculation of the sums over the depending  $N_{ijk}$  values from above.

## 3 How MMHC works

In this chapter we want to introduce an example first. By this we will explain the concept behind the algorithm. We also used this example for our computations and for comparison with the bnlearn implementation.

### 3.1 Introduction to an example

We built a data set whose underlying structure is taken from [KoFr]. It is a BN with 5 nodes. In this we have that the intelligence (high or low) and difficulty of an exam (difficult or easy) affect the students grade (good, middle, bad). The intelligence also affects the SAT (the european equivalence to this is the PISA test), which can be good or bad. The professor's letter of recommendation can then be good or bad. This depends on the grade. In the following picture also the probabilities of an event are displayed, where the smallest upper case numbers stand for high, good and easy and the higher numbers for low, bad and difficult.



**Figure 3.1:** The graph of our example with 5 nodes and some conditional independencies.

In science you normally do not start with a graph, reconstruct data and then again reconstruct the graph. You start with (random) data and try to find out which dependencies you find within your data. For our purpose this example helped to check if the code we wrote is valid. Because of randomness in the data the output can vary. If we talk about getting the "right" graph, we mean that it looks like our example. By statistical testing there is no right or wrong.

### 3.2 General workflow of MMHC

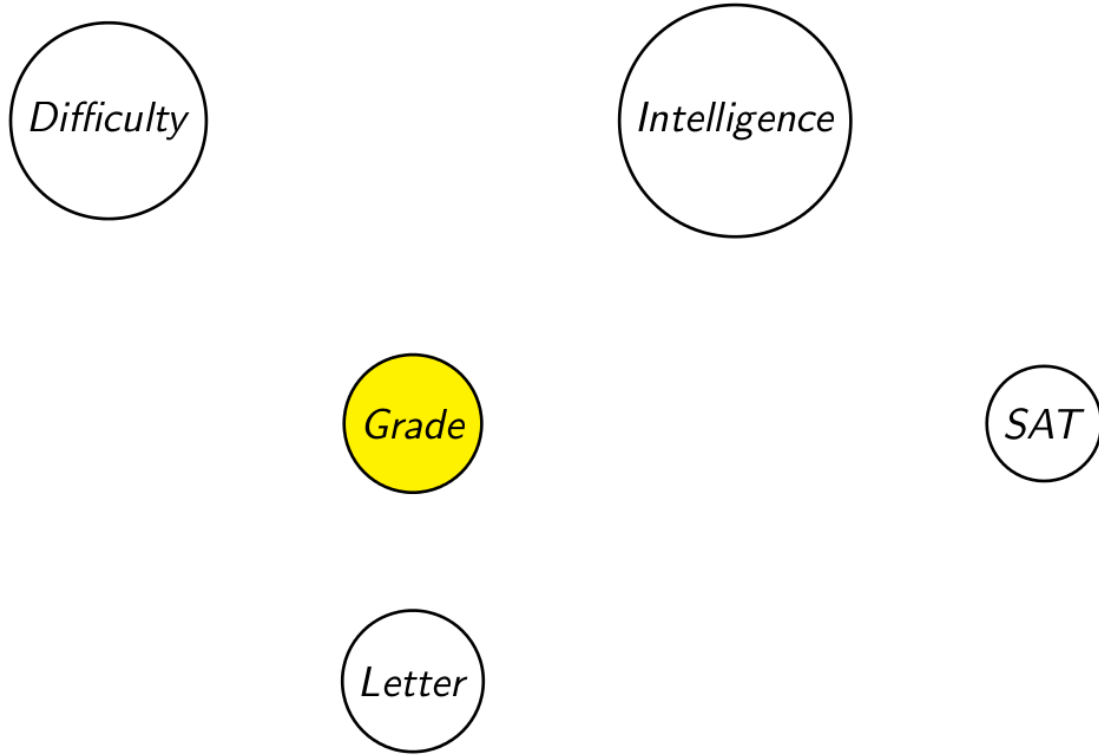
The max-min-hill-climbing algorithm generally takes your underlying data and firstly finds the skeleton of the graph. Since, it knows the structure of the graph it starts to add edges between nodes. It also reverses existing edges or deletes them, until it finds the graph which is most likely to fit best to the data.

#### Finding the skeleton (MMPC)

Finding the skeleton is the first part of the algorithm. This search is constraint-based and is called the max-min parents and children (MMPC) algorithm. The name refers to the fact that we find for a specific variable  $T$  (target) all variables which are conditional dependent on  $T$ .

All we then know is that those variables can be a parent or a children of  $T$ .

We start with an empty graph and iterate over all variables. In each step - a target  $T$  is selected - we try to find all the variables which are conditional dependent on  $T$ . We now assume - depending on our example - that the algorithm picked  $T = \text{"grade"}$  to find its parents and children.



**Figure 3.2:** The first iterations step where "grade" is selected.

To find the Parents and children (PC) we first test conditional independence of  $T$  with all nodes  $\mathcal{V} = \mathcal{D} \setminus \{\text{"grade"}\}$ . Since  $PC = \emptyset$  we are seeking for  $Ind(T; X|\emptyset)$  for all  $X \in \mathcal{V}$ .

The procedure of the MMPC algorithm then works as follows:

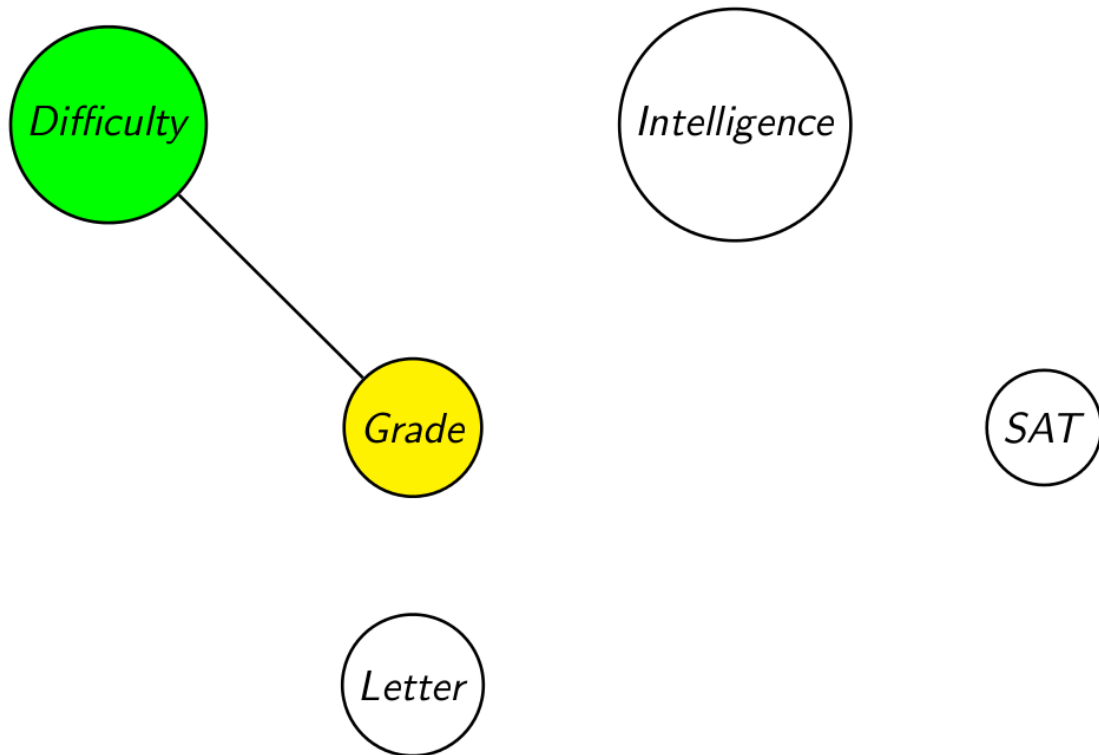
- Compute the  $G^2$  value.
- Compute the degrees of freedom  $df$ .
- Calculate the  $pvalue$  - we used the  $pchisq(G^2, df)$  in Rcpp.
- Test if the  $pvalue$  is smaller than a threshold of 0.05, if yes then keep this  $pvalue$ , it's corresponding  $G^2$  and  $X$ . If it is bigger than 0.05 then you know that  $T$  and  $X$  are conditionally independent and  $X$  can be crossed out from  $\mathcal{V}$ , i.e. you don't have to consider this  $X$  in your calculations again.



- At the end let the  $X$  join the **PC** set which has the smallest  $pvalue$ .
- If two  $pvalues$  are equaled then take the variable  $X$  with the higher  $G^2$ .
- If  $X$  joined  $PC$  then it can be crossed out from  $\mathcal{V}$ , since you already know that  $X$  and  $T$  are conditionally dependent.
- Repeat this calculation but now with all subsets of  $PC$ . Terminate the algorithm when  $\mathcal{V} = \emptyset$ .

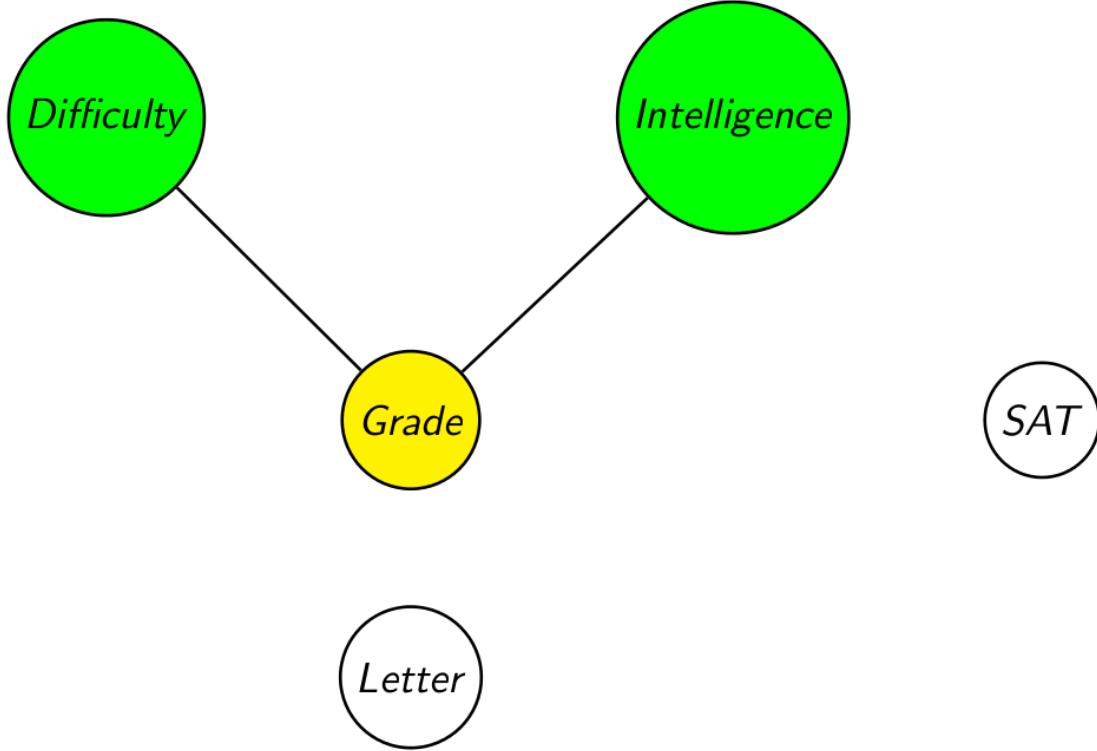
**Remark:** From section 2 we also know that two independent variables have an association value of zero. Later when we discuss the pseudo code we compute the association value which is nothing but the  $G^2$  value.

To come back to our example we assume that the strongest association between "grade" was found for "difficulty" and we found out that "grade" and "SAT" are conditionally independent. Then "difficulty" joins the  $PC$  set ( $PC = \{\text{"difficulty"}\}$ ) and we won't consider "SAT" for the next calculations, i.e.  $\mathcal{V} = \{\text{"intelligence"}, \text{"letter"}\}$ .



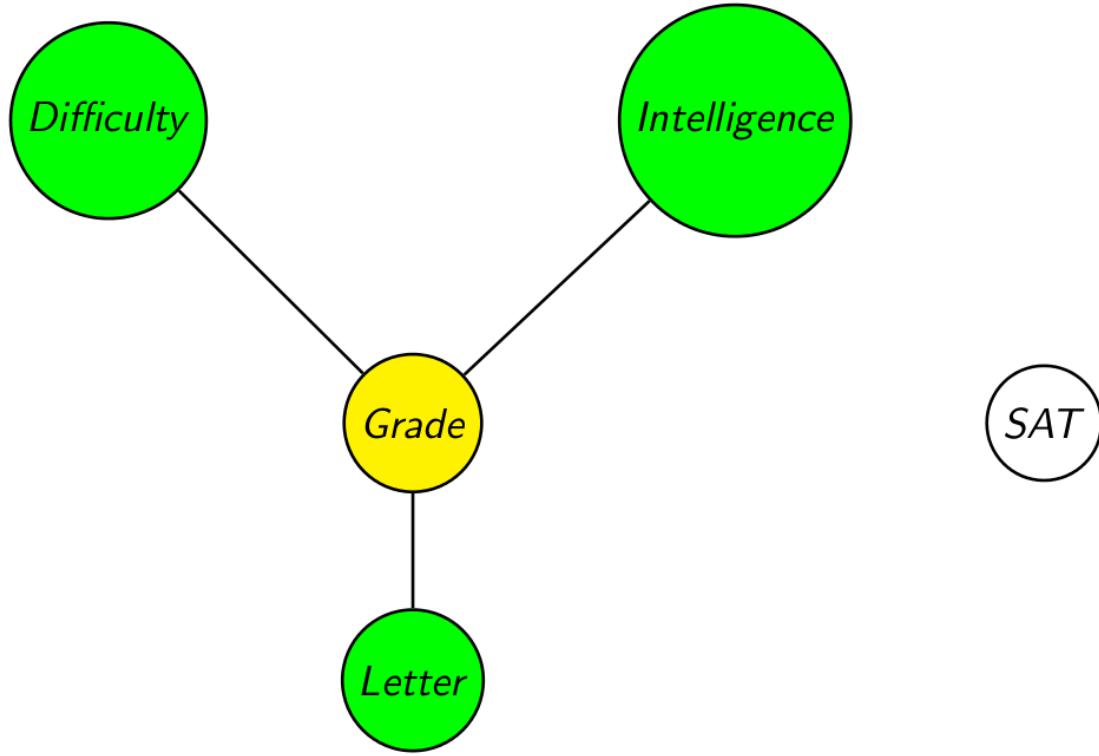
**Figure 3.3:** Here "difficulty" joined  $PC$  and is a parent or a children of "grade".

By testing  $Ind("grade"; "intelligence" | "difficulty")$  and  $Ind("grade"; "letter" | "difficulty")$  we assume that the stronger association lies between "grade" and "intelligence". Since there is an association between "grade" and "letter", too, i.e. they are not conditionally independent, we have:  $\mathcal{V} = \{"letter"\}$  and  $PC = \{"difficulty", "intelligence"\}$ .



**Figure 3.4:** Here "intelligence" joined  $PC$  and is a parent or a children of "grade".

Now we have to calculate the association given all subsets of  $PC$ , i.e. We test  $Ind(T; X|Z) \forall X \in \mathcal{V}, Z \in PC^2$ . Some of those calculations in the power set of  $PC$  are done before (e.g.  $Ind("grade"; "letter" | "difficulty")$  c.f. above), so we just saved the corresponding values and used it again. But since the  $PC$  grows you have to do more and more calculations. We now assume that "letter" also joined our set and we terminate the algorithm. Now we have all possible parents and children of "grade", which we already knew.



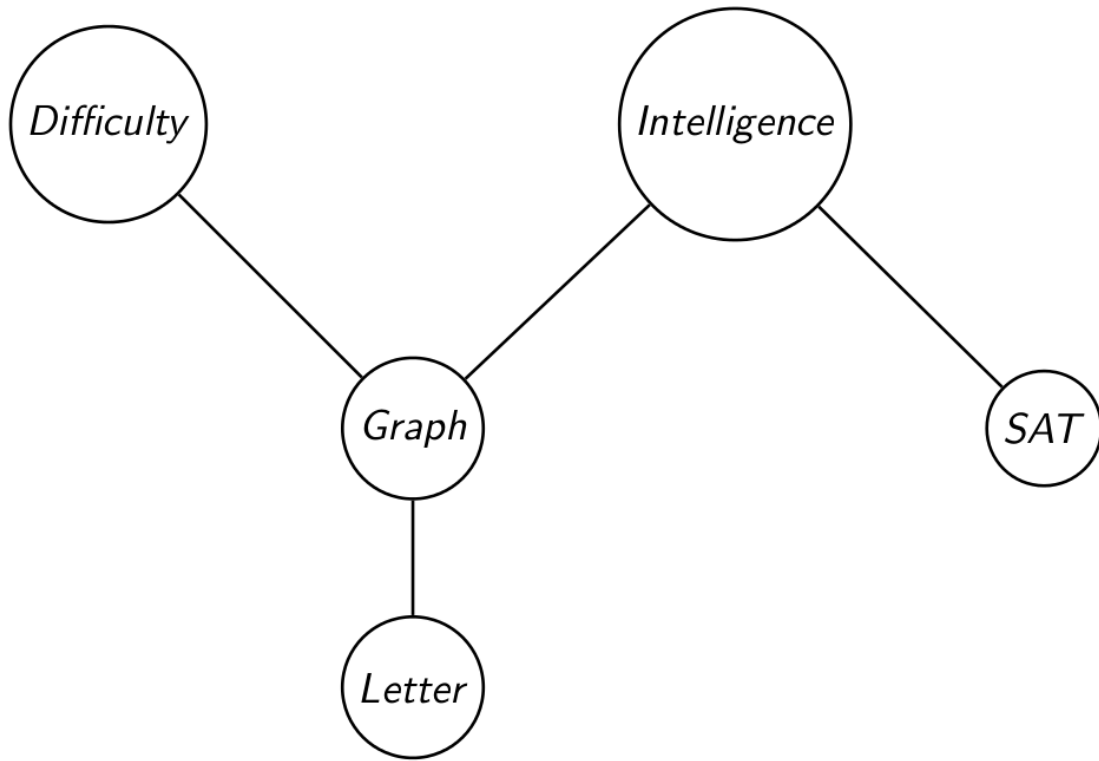
**Figure 3.5:** Here "letter" joined  $PC$  and is a parent or a children of "grade".

It may happen that there are false positive variables in the  $PC$  set. For this reason we have to check if the parents/children in the set really belong to our selected variable  $T$ . The relation between  $T$  and its parents/children is symmetric. That means for a target  $T$  and  $A \in PC_T$  it follows:

$$T \longleftrightarrow A \iff A \longleftrightarrow T \quad (3.1)$$

for a target  $A$  and  $T \in PC_A$ . Since this relation holds we have to check in a second step if this symmetrie is fullfild. For that we check for every  $X \in PC_T$  if  $T \in PC_X$ .

As we have finished all those calculations we have the structure/the skeleton of our graph.



**Figure 3.6:** The skeleton we should observe after our calculations.

Though we now the structure and that there are connection we do not know how they affect each other. For this we have to direct the edges which is done by the BDeu score.

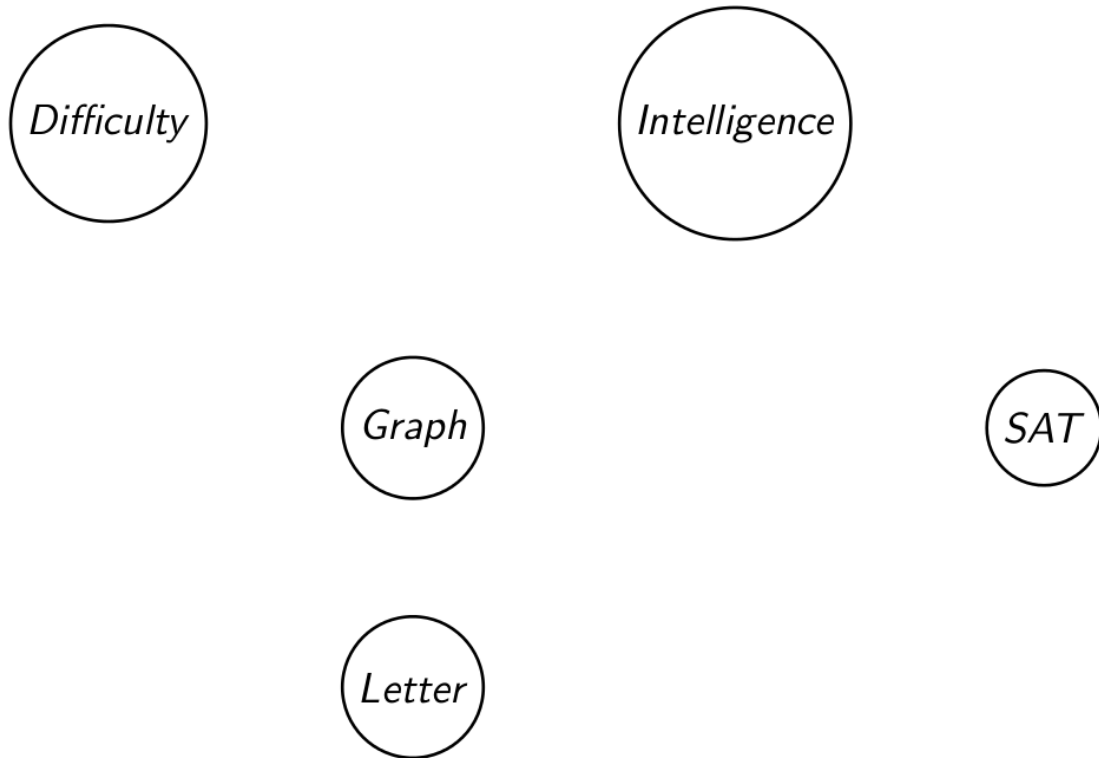
### 3.2.1 General workflow of *BDeu*

The Bayesian Dirichlet likelihood-equivalence uniform score takes as seen in section 2 the data set  $\mathcal{D}$  and a graph  $G$  as an input. What we do in one iteration step is:

- Calculate the *BDeu* score of the current graph. In the first iteration score the empty graph, i.e. the graph without any edges.
- Add an edge randomly between two conditional dependent nodes  $X$  and  $Y$ . From *MMPC* we already know those variables which are conditionally dependent.
- Calculate the *BDeu* score with the new set edge.
- If the score is higher than before the edge stays, else the edge is removed. This is the greedy search part. We are seeking for an local optimum.

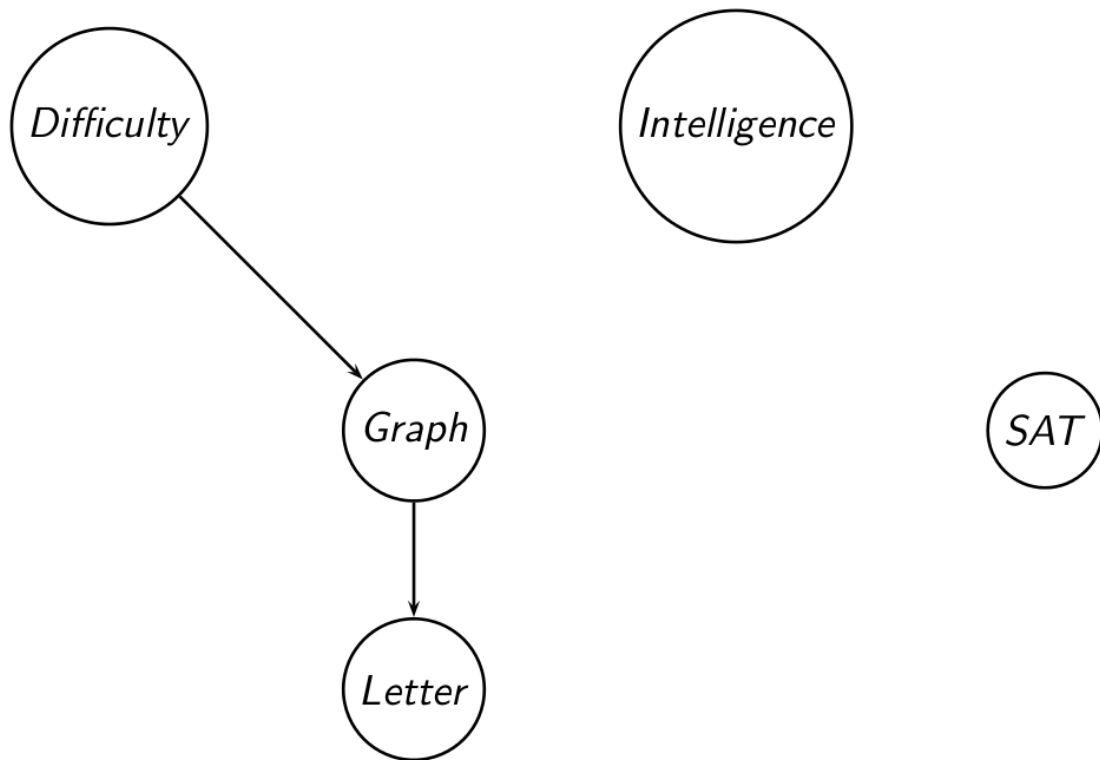
- Do this as long as the score stays in its (possible) maximum for at least 10 times in a row.
- Also apply reversing edges and deleting edges to the possible operations on the graph during one iteration step.

The following four pictures show what happens during computations:



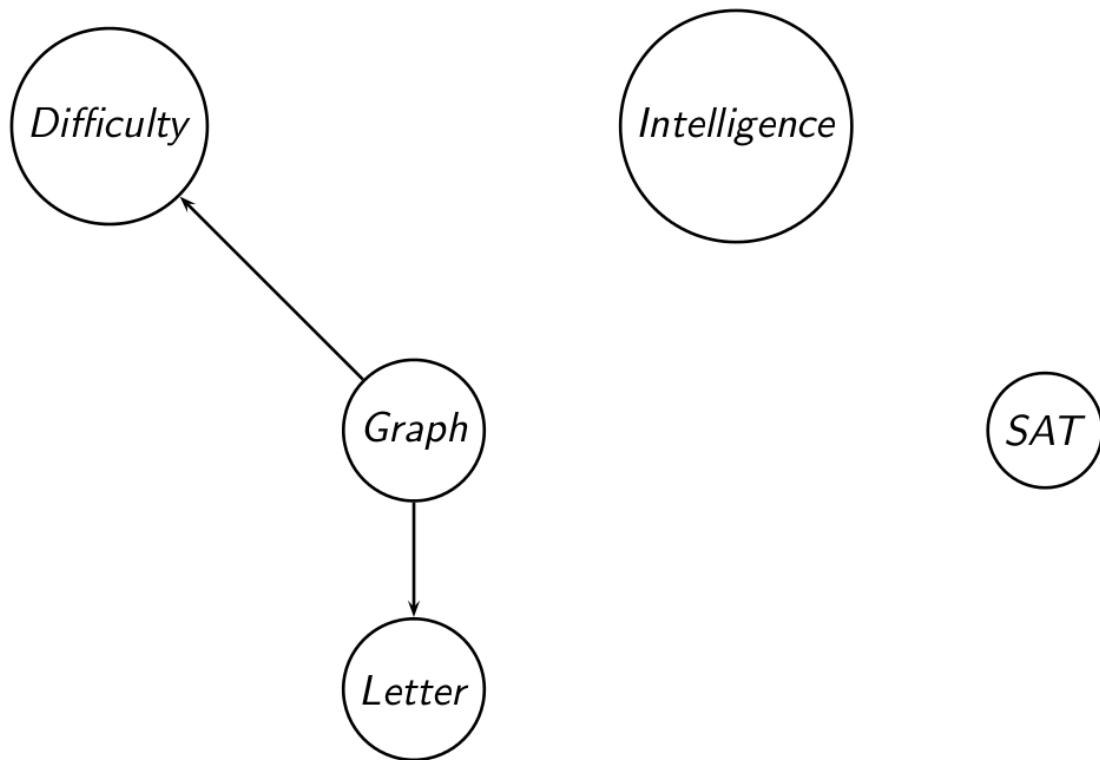
**Figure 3.7:** Starting with the empty graph.

After scoring the empty graph edges where applied to our graph.



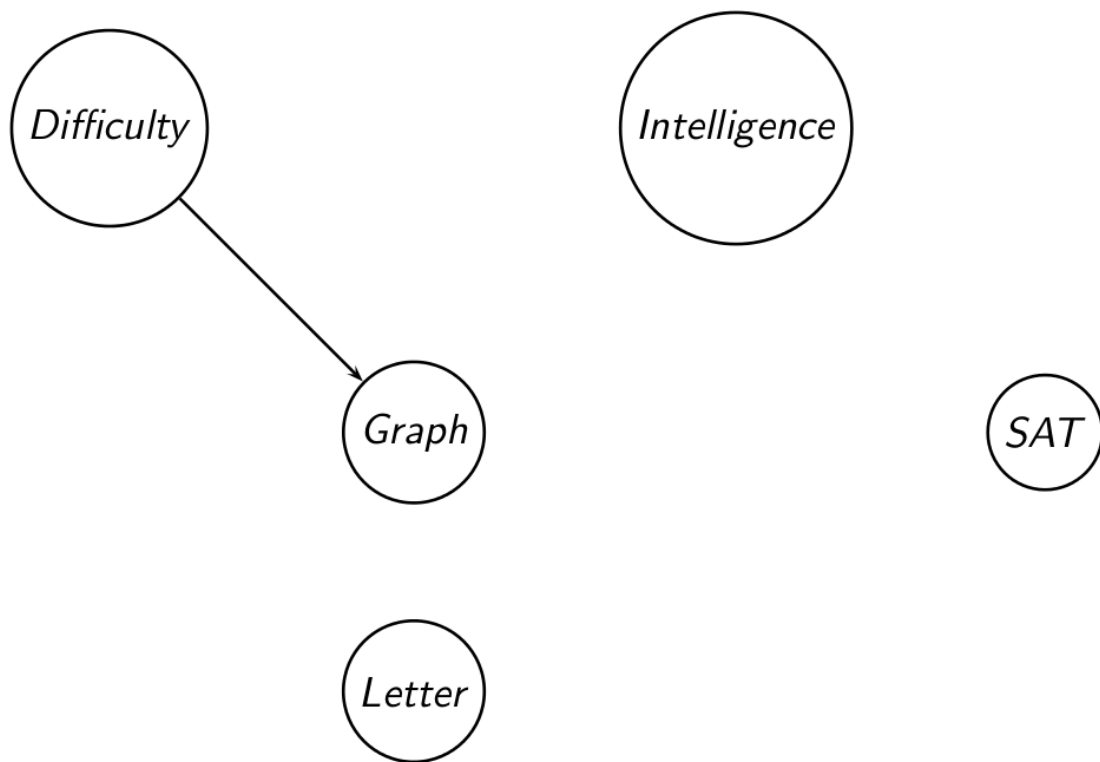
**Figure 3.8:** Adding edges between nodes.

At a certain point one edge was reversed. We assume that this reversion does not hold, because after scoring it turned out that the score didn't increase.



**Figure 3.9:** Reverse an edge.

With the back reversed edge we tried to delete one edge. As it would turn out, this edge would be set again.



**Figure 3.10:** Delete an edge.



## 4 The max-min-hill-climbing algorithm

As we stated before our main aim was to beat the existing algorithm provided by the bnlearn package in R. There are a lot of bottlenecks which slow the code down. To have the chance to beat an existing algorithm which runs quite fast, we needed to optimize the code as much as possible. For that we first want to present the pseudo code of the algorithm to be able to talk about optimization.

---

**Algorithm**    *MMHC* Algorithm

---

```
1: procedure MMHC( $\mathcal{D}$ )
   Input: data  $\mathcal{D}$ 
   Output: a DAG on the variables in  $\mathcal{D}$ 
   % Restrict
2:   for every variable  $X \in \mathcal{V}$  do
3:      $\mathbf{PC}_X = \text{MMPC}(X, \mathcal{D})$ 
4:   end for
   % Search
5:   Starting from an empty graph perform Greedy Hill-Climbing
      with operators add-edge, delete-edge, reverse-edge. Only try
      operator add-edge  $Y \rightarrow X$  if  $Y \in \mathbf{PC}_X$ .
6:   Return the highest scoring DAG found
7: end procedure
```

---

**Figure 4.1:** The pseudo code of the MMHC algorithm. Line 2-4 represent the loop over all nodes calling the MMPC function. At line 5 the scoring starts.

As mentioned before, the MMHC function has two parts. First we find the skeleton and then direct the edges of the skeleton graph. We first discuss the *MMPC* (max-min parents and children) function which is executed with the observational data  $\mathcal{D}$  in the for loop (over all possible nodes in the graph) and then we take a closer look at the scoring function starting at line 5 (*BDeu* score).

## 4.1 max-min parents and children (MMPC))

The max-min parents and children (MMPC) is the algorithm which reconstructs the graph skeleton from data. To bring you the concept of this closer, let us first present the pseudo code of it:

---

**Algorithm**    Algorithm *MMPC*

---

```

1: procedure MMPC( $T, \mathcal{D}$ )
2:    $\mathbf{CPC} = \overline{MMPC}(T, \mathcal{D})$ 
3:   for every variable  $X \in \mathbf{CPC}$  do
4:     if  $T \notin \overline{MMPC}(X, \mathcal{D})$  then
5:        $\mathbf{CPC} = \mathbf{CPC} \setminus X$ 
6:     end if
7:   end for

8:   return  $\mathbf{CPC}$ 
9: end procedure

```

---

**Figure 4.2:** The pseudo code of the MMPC function. First the **CPC** set is computed by the  $\overline{MMPC}$  function. Afterwards the false positives are erased.

In this algorithm another function  $\overline{MMPC}$  is executed. In Figure 4.1 this function - including the *MaxMinHeuristic* function - is shown. Back to Figure 4.1, we see, that  $\overline{MMPC}$  is executed twice. Firstly at line 2 and secondly in the if-statement at line 4. As we stated in Equation 3.1, for two variables they are connected, the symmetrie holds. At line 2 we only now that some  $X$  are in the set **CPC** for a target variable  $T$ , i.e. the  $X$  is a parent or a child of  $T$ . By testing if  $T$  is also a parent or child of  $X$ , we see if the symmetrie holds, if not,  $X$  is removed from the **CPC** set. The more interesting part here is the  $\overline{MMPC}$  function.

---

**Algorithm**  $\overline{MMPC}$  Algorithm

---

```

1: procedure  $\overline{MMPC}$  ( $T, \mathcal{D}$ )
   Input: target variable  $T$ ; data  $\mathcal{D}$ 
   Output: the parents and children of  $T$  in any Bayesian
   network faithfully representing the data distribution
   %Phase I: Forward
2:   CPC =  $\emptyset$ 
3:   repeat
4:      $\langle F, assocF \rangle = \text{MaxMinHeuristic}(T; \mathbf{CPC})$ 
5:     if  $assocF \neq 0$  then
6:       CPC =  $\mathbf{CPC} \cup F$ 
7:     end if
8:   until CPC has not changed

   %Phase II: Backward
9:   for all  $X \in \mathbf{CPC}$  do
10:    if  $\exists \mathbf{S} \subseteq \mathbf{CPC}$ , s.t.  $Ind(X; T | \mathbf{S})$  then
11:      CPC =  $\mathbf{CPC} \setminus \{X\}$ 
12:    end if
13:  end for

14:  return CPC
15: end procedure

16: procedure  $\text{MAXMINHEURISTIC}(T, \mathbf{CPC})$ 
   Input: target variable  $T$ ; subset of variables CPC
   Output: the maximum over all variables of the minimum asso-
   ciation with  $T$  relative to CPC, and the variable that achieves
   the maximum
17:   $assocF = \max_{X \in V} \text{MinAssoc}(X; T | \mathbf{CPC})$ 
18:   $F = \arg \max_{X \in V} \text{MinAssoc}(X; T | \mathbf{CPC})$ 
19:  return  $\langle F, assocF \rangle$ 
20: end procedure

```

---

**Figure 4.3:** Here you see both: the  $\overline{MMPC}$  function and the calculation of the association between  $X$  and  $T$  by the *MaxMinHeuristic* function.

In this subfunction we have one interesting routine: the *MaxMinHeuristic*.

#### 4.1.1 The *MaxMinHeuristic* function

This function finds the  $X$  which maximizes the measure of association between  $X$  and  $T$  given the current **CPC** set. The current **CPC** set is therefor passed into the function. The return of the function is the  $X$  and the value of association between  $X$  and  $T$  given **CPC**. At this point, we left one important question: How do we measure this association.

## 4.2 The scoring function

The scoring works as we stated earlier in this report. The important thing here is that we only add an edge from  $X$  to  $Y$  iff  $Y$  is a member of the parents/children set of  $X$ . In the case of scoring most interesting part was *eta*. There is no "perfect" value one could take. There is a lot of research done to find the optimal *eta* for an algorithm. In our case there are two possibilities: we could set  $eta = 1$  or we could calculate it with  $eta = \frac{1}{N} \sum_{i=1}^N |X_i|$  where  $N$  is the number of variables and  $|X_i|$  is the cardinality of the variable  $X_i$ . In our tests it did not affect the results or the running time. For that reason we decided to set it to one.

## 4.3 Computational optimization

We decided to write our algorithm object oriented. To save running time we instantiate our class. The constructor of the class will do some precomputations like getting the cardinality of the variables. Since the class holds all information the algorithm needs, we have our information in every method of our class available. This has the effect that we do not pass the objects (vectors, lists, etc.) from function to function. We work with the right in our class. Another thing we needed to know was that the underlying data matrix is stored differently in R and C++. R saves a matrix firstly by column and then by row, i.e. you have a pointer to a column which holds the pointer to the row. In C++ it is vice versa. We just had to skip the normal way of iterating over matrices in C++.

To compute the  $G^2$  you need to count the  $S_{ijk}^{abc}$ , etc. values. For this we implemented a powerful application, hash tables. At first side we tried to use the `unordered_map` of C++. This is a hash table within you can look up your elements in constant time. We precomputed all possible combinations of our variables, i.e. we tested  $Ind(X; Y|Z)$  for all possible  $X$ ,  $Y$  and sets  $Z$ . Our idea was that we only need to look up the values we need in our methods. Of course, this did not improve running time since you had to compute all permutations. With 5 variables and all possible  $X, Y, Z$  you have  $2! + 3! + 4! + 5! = 152$  calculations.

We then decided to use n-dimensional arrays hash tables where we convert the values of a variable into integer from  $1, \dots, n$  where  $n = |X|$  are the possible states of one variable. We then use those integer values as indices for our arrays and increment the value of the array everytime it is accessed. Afterwards we just needed to iterate over it again and use the values in the arrays for our calculations. The biggest array we allocate is 5-dimensional. For small problems with less nodes in the graph and less dependencies within those nodes this code runs extremely fast. For bigger problems there also exists an implementation to compute the  $G^2$  value, but it is quite slow and we try to avoid using it. The difficulty behind our arrays is that we use pointer arrays - otherwise we could not allocate 5-dimensional arrays. Here we

needed to be careful to free the memory after usage.

To present the validity of our code we recall the example from the previous section. For this example we computed 10,000 observations, i.e. we have a R data frame  $df \in \mathbb{N}^{10000 \times 5}$ , where all nodes have binary values 1,2 except the node holding the grade. This has three states: 1,2,3. The time our algorithm needs to compute the skeleton of the graph, i.e. using the *MMPC* function, is as follows:

```
df <- student(10000)
bm <- benchmark(C$mmpc(),
               columns = c("test", "elapsed", "relative"),
               replications = 1)
      test elapsed relative
1      C$mmpc()   0.043    1.000
```

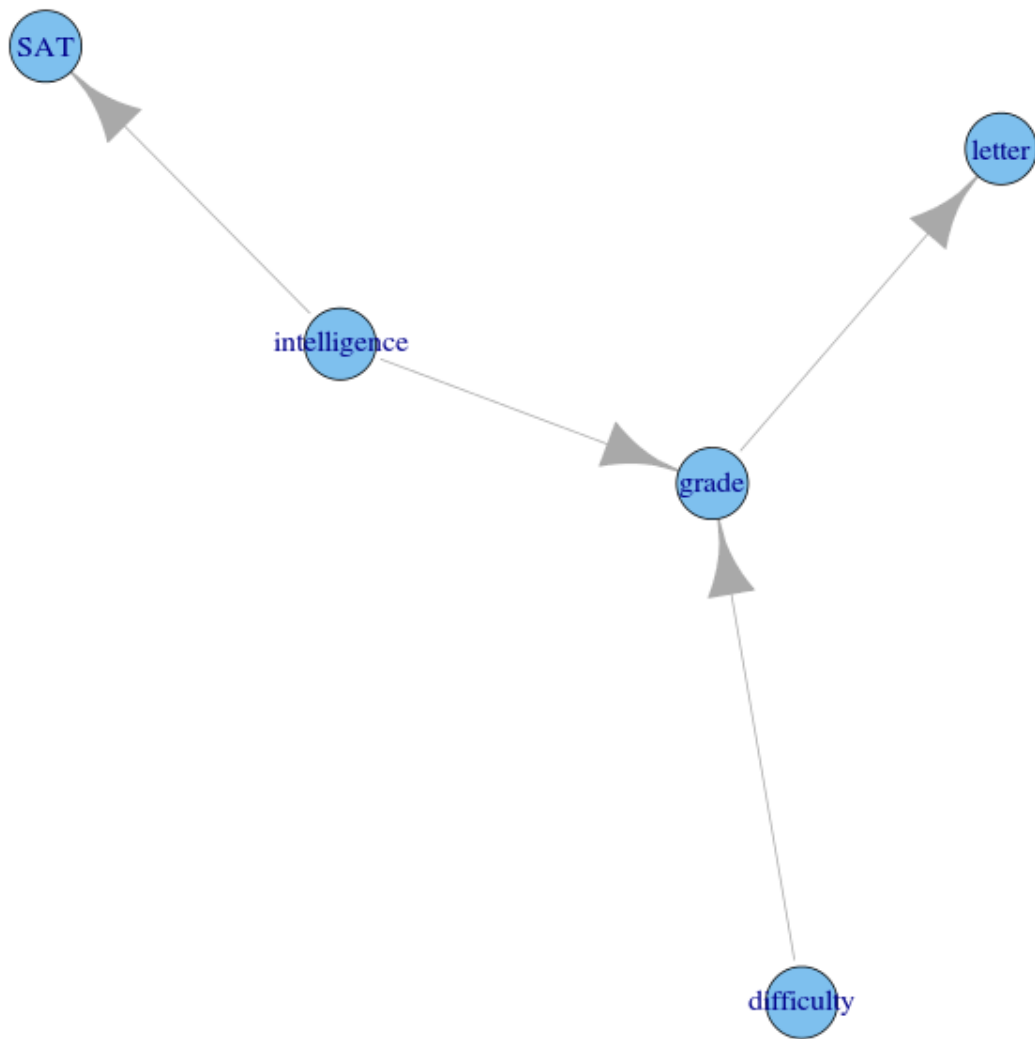
where the elapsed time is measured in seconds. We see that our implementation is extremely fast, even for a high number of observations. In the next chapter we will see if this time is good enough to be faster than the *bnlearn* function.

By looking at the whole algorithm we have the following running time:

```
df <- student(10000)
bm <- benchmark(C$mmhc()
               columns = c("test", "elapsed", "relative"),
               replications = 1)
-27812.9 # score returned by BDeu(...)
-27812.9 # score returned by C$mmpc()
print(bm)
      test elapsed relative
1      C$mmhc()    1.000
```

again with the elapsed time in seconds. As we see the most time we loos was in the *BDeu* scoring. At the end of this work we did not find a way to make this function faster. If you are interested in the progress we recommend you to follow our Github account, stated at the end of this report.

The "proof" of the validity now follows from the plot we present. Here you can see that it looks like the one we started with. To plot our results we used the "igraph" package in R.



**Figure 4.4:** The graph returned looks like the one we started with.

## 5 A comparison to the bnlearn package

Beating the bnlearn was a very challenging and difficult task. It was a PhD thesis project which went over at least one year. We only had three months for this task. But with some tricks and computational knowledge we were able to implement an algorithm which is fast and valid. The question now is, if our version is faster and the answer is no. There are several reasons we want to discuss in this section. As we will see our running time is also quite good and developing goes on. For the discussion we consider our example from section 3.1. First of all we want to compare the running times of the *MMPC* function. For that we run our implementation and the one of bnlearn with the example having 1,000, 5,000 and 10,000 numbers of observations. For the rest of this section we measure the elapsed time in seconds. The R code looks as follows, where we used the R package "rbenchmark" to be able to run the benchmark function:

```
MMPC <- function(myDataFrame) {  
  C <- new(MMHC, myDataFrame)  
  C$mmpc()  
}  
  
myDataFrame <- student(1000)  
C <- new(MMHC, myDataFrame)  
  
bm <- benchmark(MMPC(myDataFrame), mmpc(myDataFrame),  
  columns = c("test", "elapsed", "relative"),  
  replications = 1)
```

The parameter of the *MMPC* function is a R data frame. This function allocates our class *C*. The method *mmpc()* is executed. To have access to the *PC* set we would need to type to the console *C\$pc()*. This method returns the *PC* set. If we compare our implementation with the bnlearn function with 1,000 observations we get:

```
test elapsed relative
```

```
2 mmpc(myDataFrame)    0.018    2.571
1 MMPC(myDataFrame)    0.007    1.000
```

If we repeat this for 5,000 observations we have:

```
test elapsed relative
2 mmpc(myDataFrame)    0.029    1.074
1 MMPC(myDataFrame)    0.027    1.000
```

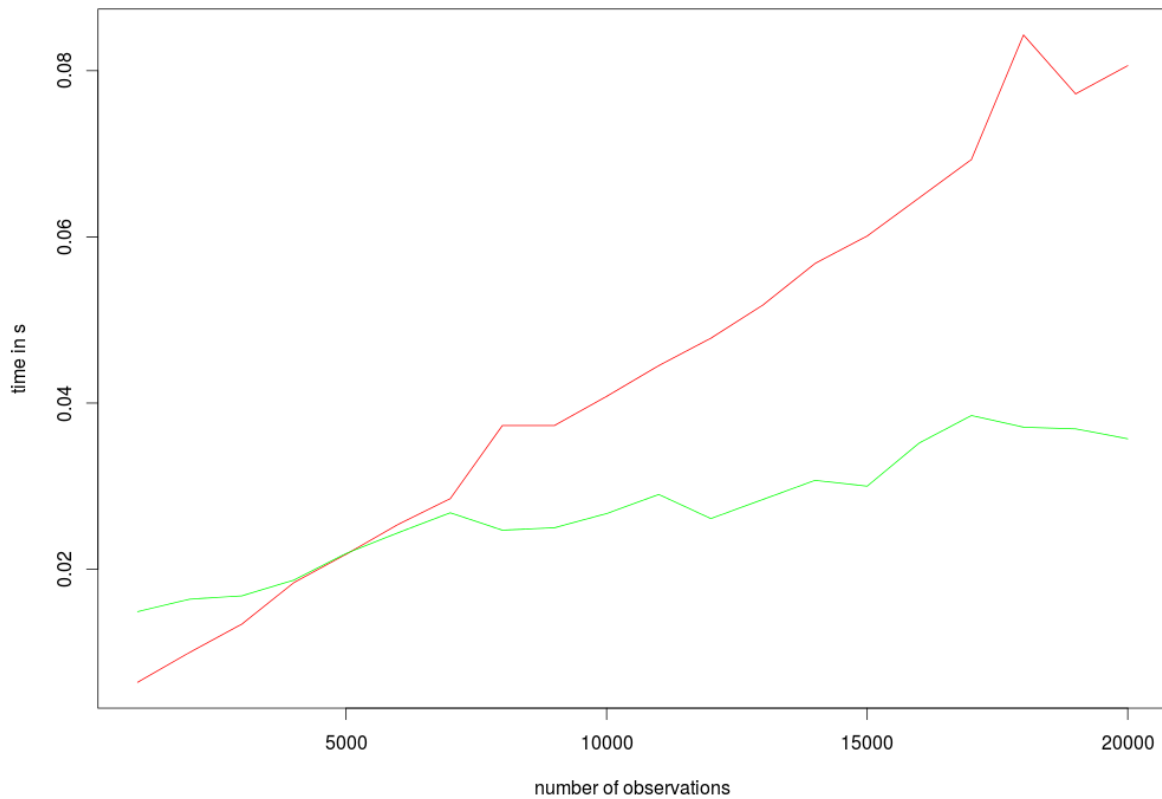
and for 10,000 observations:

```
test elapsed relative
2 mmpc(myDataFrame)    0.021    1.000
1 MMPC(myDataFrame)    0.050    2.381
```

As we can see, our implementation is much faster for a smaller amount of data. The bigger the data gets, the slower our code runs. The interesting thing is, that the bnlearn package seems to have constant running time. Another thing is that the tests this algorithm uses are not exactly the same as ours. We tried to set the parameters of the function rightly but it did not work with our data (setting test to "x2" and score to "bde"). For other data it did not make big differences. At the end of this chapter we will see one comparison where the setting of the right parameters worked.

Before we compare the whole function we want to present a plot where we benchmarked both function 20 times in the range of 1,000 to 20,000 numbers of observations.





**Figure 5.1:** From 1,000 to 7,500 observations our algorithm (red line) is quite fast. Above this number bnlearn one (green line) is much faster.

It seems that the break even point of those two implementations is at about 7,500 numbers of observations.

You may ask why we do not compare the scoring function, only the *MMPC* and the complete *MMHC*. The reason why we do this is, that we can not run the scoring explicitly in the bnlearn function - only within our package. But before we compare the running time of the *MMHC* function we want to present again the R code where *C* is again a class:

```
MMHC <- function(myDataFrame) {
  C <- new(MMHC, myDataFrame)
  C$mmpc()
  C$mmhc()
}
```

```
myDataFrame <- student(1000)
C <- new(MMHC, myDataFrame)

bm <- benchmark(MMHC(myDataFrame), mmhc(myDataFrame),
  columns = c("test", "elapsed", "relative"),
  replications = 1)
```

The running time varies a lot because our loop is terminated if the score did not increase at least 10 times in a row. That means we apply all our operations to the graph quite often even if we already found the best fit with the data. So for the comparison we took the most common outcomes to compare:

	test	elapsed	relative
2 mmhc(myDataFrame)		0.024	1.000
1 MMHCs(myDataFrame)		0.038	1.583

with 1,000 observations

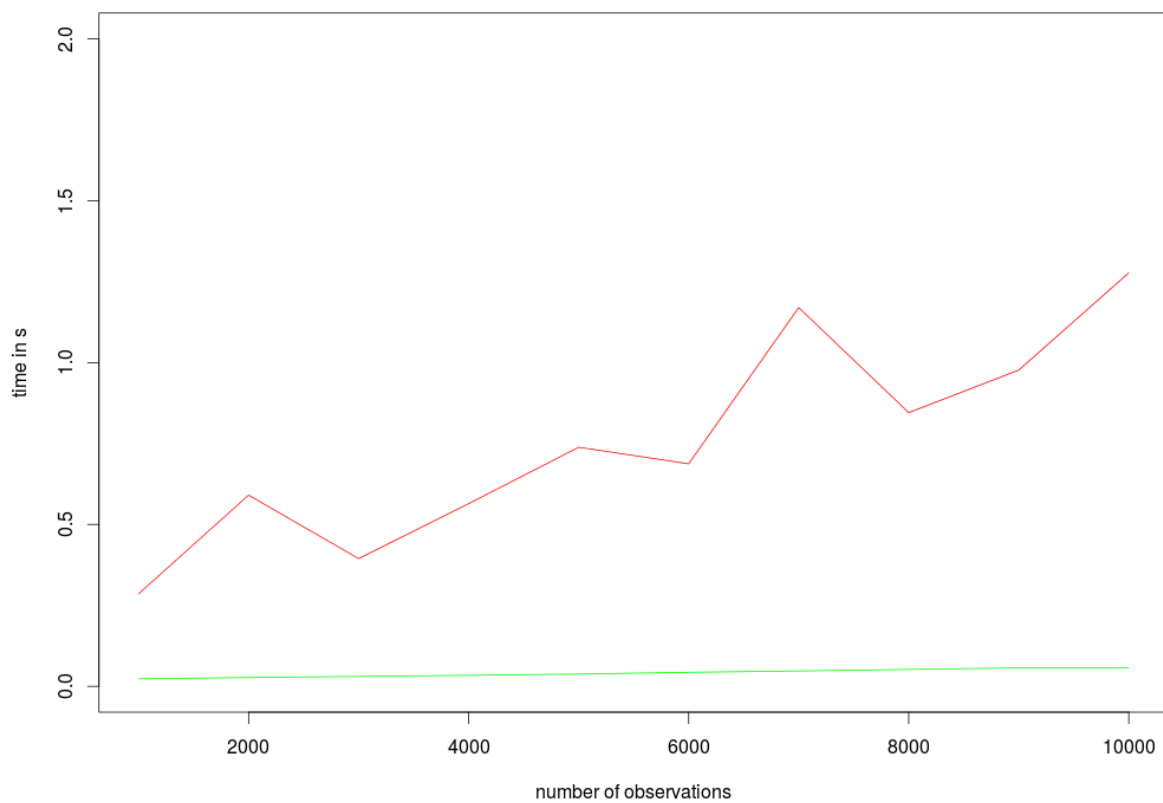
	test	elapsed	relative
2 mmhc(myDataFrame)		0.039	1.000
1 MMHCs(myDataFrame)		0.228	5.846

with 5,000 observations

	test	elapsed	relative
2 mmhc(myDataFrame)		0.068	1.000
1 MMHCs(myDataFrame)		0.335	4.926

with 10,000 observations

Again we see that our algorithm grows as the data grows. For this part we ran both functions 10 times from 1,000 to 10,000 numbers of observations.



**Figure 5.2:** Hear you see that our algorithm (red line) grows linearly where as the bnlearn one (green one) is constant.

As the running time of our code grows linearly the bnlearn package seems to stay about constant. We could not figure out how and why the bnlearn package as no growth in its running time but with no result.

But our running time is maximally 1.5 seconds for a huge data frame. This is still a very good outcome. Maybe in the next months we can improve our algorithm more and more.

## 6 Conclusion

Despite the fact that we could not be faster than the `bnlearn` package, it was worth to implement this algorithm. In the future there will be the chance to improve running time and make our algorithm faster. But for now we provide the chance to construct BNs from a smaller amount of data which is faster than existing implementations. We really hope that research goes on and maybe there will be a faster implementation or even a new algorithm which reconstructs BNs and runs faster than existing ones.

# Bibliography

- [TBA] Ioannis Tsamardinos, Laura E. Brown, Constantin F. Aliferis,  
The max-min hill-climbing Bayesian network structure learning algorithm,  
Springer Science + Business Media, Inc. 2006
- [NBBCW] Chris J. Needham<sup>1</sup>, James R. Bradford, Andrew J. Bulpitt, Matthew A. Care and  
David R. Westhead,  
Predicting the effect of missense mutations on protein function: analysis with Bayesian  
networks,  
<http://www.biomedcentral.com/1471-2105/7/405>
- [PKA] [http://www-ekp.physik.uni-karlsruhe.de/~zupanc/WS1011/docs/Datenanalyse2010\\_3.pdf](http://www-ekp.physik.uni-karlsruhe.de/~zupanc/WS1011/docs/Datenanalyse2010_3.pdf)
- [P88] Pearl, 1988
- [SGSN] Spirtes, Glymour & Scheines, 1993, 2000;  
Neapolitan, 2003
- [Ca06] Luis M. de Campos,  
A Scoring Function for Learning Bayesian Networks based on Mutual Information and  
Conditional Independence Tests,  
Journal of Machine Learning Research 7, 2006
- [1] [SKM] On Sensitivity of the MAP Bayesian Network Structure to the Equivalent Sample  
Size Parameter,  
Tomi Silander and Petri Kontkanen and Petri Myllymäki,  
UAI, 2007
- [KoFr] Daphne Koller, Nir Friedman,  
Probabilistic Graphical Models: Principles and Techniques (Adaptive Computation and  
Machine Learning),  
The MIT Press,  
First edition (16. November 2009)