



university of  
 groningen

# Evaluation of a General Bayesian Method using the Cholesky Decomposition to model the trans- mission of intelligence

## Master's Thesis

Faculty of Science and Engineering

Program: Science and Communication

D.M. Heeg (s2577135)

First supervisor: dr. M.A. Grzegorzcyk

Second supervisor: dr. W.P. Krijnen

## Abstract

The classical twin design has proven to be a powerful method in order to study the transmission of intelligence from parent to child. However, it comes with a few unrealistic assumptions such as random mating. In this paper the classical twin design is extended with cultural transmission, phenotypic assortment, social homogamy, the dominance effect or a combination of these extensions. In order to avoid analyzing sum scores, the Rasch Model is used as a measurement model and both the genetic and measurement model are put in a Bayesian framework, in order to estimate both simultaneously. A Cholesky decomposition is applied to make the models more easy to access and finally a Metropolis-Hasting algorithm was implemented to keep the model transparent. Six different models are compared by calculating the DIC scores, to find the genetic model that represents the real-world data the best.

## Preface

This thesis is made as a part of the master Science and Communication. Instead of starting with courses and finishing with a thesis, the Science and Communication master starts with a thesis as an opportunity for the students to show that they are capable of conducting research in their own scientific field, in this case Mathematics. So, even though this thesis is presented as a master's thesis, it is written without any knowledge from master courses. The goal of this research is therefore not to show that obtained knowledge during the master can be applied in a relevant research, but to show that the student is capable of making herself familiar with a scientific subject, conduct research about this subject and by means of writing a thesis communicate the findings clearly to others. This thesis is intended for undergraduate students with a basic knowledge of Bayesian Statistics.

# Contents

<b>1</b>	<b>Introduction</b>	<b>4</b>
<b>2</b>	<b>Data</b>	<b>5</b>
<b>3</b>	<b>Methodology</b>	<b>6</b>
3.1	Measurement Model . . . . .	6
3.1.1	Item Response Theory . . . . .	6
3.1.2	The Rasch Model . . . . .	7
3.2	The Genetic model . . . . .	8
3.2.1	The classical twin design . . . . .	8
3.2.2	Cultural transmission . . . . .	9
3.2.3	Assortative mating . . . . .	11
3.2.4	The dominance effect . . . . .	13
3.2.5	The six models . . . . .	14
3.3	The Cholesky Decomposition . . . . .	16
3.4	A Bayesian Framework . . . . .	18
3.4.1	The Bayesian Framework . . . . .	18
3.4.2	Model Fit . . . . .	19
<b>4</b>	<b>Simulated data</b>	<b>19</b>
<b>5</b>	<b>Results</b>	<b>20</b>
5.1	Simulations . . . . .	20
5.2	Model Fit . . . . .	22
<b>6</b>	<b>Conclusion/Discussion</b>	<b>23</b>
<b>7</b>	<b>Bibliography</b>	<b>26</b>
<b>A</b>	<b>Genetic Models</b>	<b>27</b>
A.1	Model 1 . . . . .	27
A.2	Model 2 . . . . .	28
A.3	Model 3 . . . . .	29
A.4	Model 4 . . . . .	30
A.5	Model 5 . . . . .	32
A.6	Model 6 . . . . .	32
<b>B</b>	<b>R Codes</b>	<b>34</b>
<b>C</b>	<b>Convergence Plots</b>	<b>36</b>
<b>D</b>	<b>Results: Parameter Estimates</b>	<b>37</b>

# 1 Introduction

Individual differences in intelligence tend to cluster within families. This similarity is due to a mixture of different mechanisms such as genetic relatedness, environmental similarities and cultural transmission between family members (Van Leeuwen et al., 2008). The classical twin design, in which monozygotic (MZ) and dizygotic (DZ) twin pairs are studied, is a powerful method to detect genetic influences on such resemblance. Since MZ twin pairs share nearly all of their DNA and DZ twin pairs share only 50 percent, a distinction can be made between shared genetics and shared environmental effects. A larger resemblance of MZ twins than of DZ twins would suggest that genetics is a more dominant factor than environment on the intelligence of an individual. In adoption studies the lack of genetic transmission from the parents to the adopted offspring gives the opportunity to clearly distinguish effects of the environment from effects of the genes.

Although the classical twin design has turned out to be very effective, research has shown that it can be extended in various ways to make it more fitting to the real world. Four of such extensions will be discussed in this article. Starting with assortative (non-random) mating. In the classical twin design random mating is assumed. Meaning that people randomly choose a partner to mate with. However, significant evidence has been found for a clear resemblance in intelligence between spouses (Fulker & DeFries, 1983; Eaves et al., 1989; Van Leeuwen et al., 2008), which would imply assortative mating. When random mating is assumed while assortative mating actually occurs this could lead to biased results. The effect of genes would be interpreted too high and effects from the environment would be too low. Assortative mating might be caused by phenotypic assortment (1) or social homogamy (2). The latter defines the phenomena of spouse selection based on social stratification. The former defines spouse selection based on the preference of similar characteristics. Moreover, the classical twin design can be extended by assuming the presence of cultural transmission (3). Cultural transmission means that the phenotype of the offspring depends on not only their environmental circumstances but also on the ability of their parents to pass information to them. One could argue that parents with a higher IQ could provide their offspring with a more stimulating environment than parents with a lower IQ. Furthermore, the dominance effect (4) can be included into the classical twin design. This effect is the influence of the dominant genes on the phenotype of an individual. For example, say that a bird has black feathers when its genotype is A1A1 and white feathers when its genotype is A2A2. If there is no dominance interaction between these alleles, then a bird with genotype A1A2 would have grey feathers. However, if there is a dominance interaction, then genotype A1A2 would be more like either genotype A1A1 or A2A2, depending on which allele is dominant. If A1 was dominant to the A2 allele, then the bird would have had black feathers. Likewise, the bird would have had white feathers if A2 was dominant to the A1 allele.

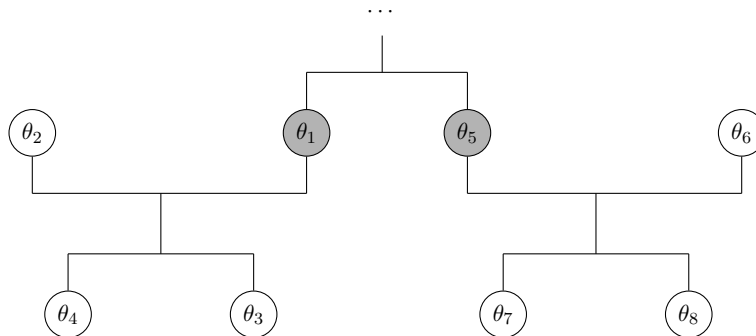
To avoid working with sum scores, which could lead to biased estimates (Van den Berg, Glas & Boomsma, 2007), the genetic model will be combined with a measurement model from the Item Response Theory (IRT). To simultaneously estimate both models, a Bayesian framework will be used (Van den Berg et al., 2007; Van Leeuwen et al., 2008; Schwabe & Van den Berg, 2014). In order to use off-the-shelf (OTS) Bayesian software for this new methodology, the model has to be re-specified into a directed acyclic graph (DAG). Veldkamp et al. (unpublished) proposed to use a Cholesky Decomposition to accomplish this and illustrated how this could be done for one explicit model. Since this method makes models more easily accessible for applied researchers, this method will be used in this thesis as well. However, instead of using

off-the-shelf Bayesian software such as openBUGS and JAGS, an R package is written which includes a Metropolis-Hasting algorithm to avoid using the OTS software as black boxes.

The question that is central in this thesis is: which genetic model represents the transmission of intelligence within families the most accurate with respect to the real world? This question will be investigated by directly comparing six different models using the deviance information criterion (DIC) of Spiegelhalter et al. (2002). The six models are all extended versions of the classical twin design, extended with phenotypic assortment, social homogamy, cultural transmission, the dominance effect or combinations of these four phenomena. In the upcoming chapters the implementation of the measurement model will be firstly explained. Next, the four types of extension of the classical twin design will be elaborated on and the six models will be defined. Thereafter, it will be explained how the Cholesky Decomposition can be applied and how the complete model is put into a Bayesian framework. When the methodology has been discussed, the outcomes of simulations will be given. Convergence plots are included, followed by an overview of all DIC values for each model and a discussion of the results.

## 2 Data

In order to test the various models for the transmission of intelligence, their fit to intelligence data is tested. The fit of the model is tested through eight family members in one family: a twin pair, their spouses and their offspring. Therefore, intelligence data is needed from families with eight family members in this particular structure. Figure 1 is a family tree, showing the structure of the families that participated in this research.  $\theta_i$  is used to represent the intelligence level of family member  $i$ , with  $i=1, \dots, 8$ . Table 1 explains the relationships between  $\theta_i$  more explicit.



**Figure 1:** The structure within a participating family, consisting of a twin pair (grey), their spouses ( $i = 2, 6$ ) and their offspring ( $i = 3, 4, 7, 8$ ).  $\theta_i$  represents the level of intelligence for the  $i$ th family member.

A genetic model including all of these eight family members and the relationships between those members is quite complicated. To simplify, only visual representations of the genetic models are given for one co-family and for twins. With these two figures, a representation can be made for all eight family members, by pasting a second co-family ( $i = 5, \dots, 8$ ) next to the first co-family ( $i = 1, \dots, 4$ ). And connecting  $\theta_1$  with  $\theta_5$  using the paths shown in the visual representation of the genetic model for twins.

The methods used in this research are based on an unpublished paper from Veldkamp, Schwabe and Van den Berg in which they use such families as well to conduct

**Table 1:** Descriptions of the  $\theta_i$  as seen in Figure 1.  $\theta_i$  represents the intelligence level of individual  $i$ .

Parameter	Description
<b>Co-family 1</b>	
$\theta_1$	Mother/Father, twins with $\theta_5$ from co-family 2
$\theta_2$	Spouse of $\theta_1$
$\theta_3$	First offspring from $\theta_1$ and $\theta_2$
$\theta_4$	Second offspring from $\theta_1$ and $\theta_2$ , sibling of $\theta_3$ .
<b>Co-family 2</b>	
$\theta_5$	Mother/Father, twins with $\theta_1$ from co-family 1
$\theta_6$	Spouse of $\theta_5$
$\theta_7$	First offspring from $\theta_5$ and $\theta_6$
$\theta_8$	Second offspring from $\theta_5$ and $\theta_6$ , sibling of $\theta_7$

data. They used the intelligence data of Reynolds et al. (1996), consisting of item data obtained by the Raven Progressive Matrices (RPM) test (Raven, 2000). This test contains 60 test items, distributed over 5 sets ( $A - E$ ). In each set items increase in difficulty. Every item is dichotomously scored, 0 for an incorrect response, 1 for a correct response.

Unfortunately, it was not possible to get this exact same data set. Prof. dr. Van den Berg from the University of Twente was asked for a synthetic data set, which could have been used for an application of the models. However, it was not possible to receive such a synthetic data set in time for this research, and since there were some time restrictions, only simulated data was used in this research. In Section 4.1 it is explained how the data was simulated.

## 3 Methodology

### 3.1 Measurement Model

#### 3.1.1 Item Response Theory

Item Response Theory (IRT) provides a class of models describing the relationship between individual item responses and the characteristics of a test. In these models, the interaction between test items and the responses of an individual on the test is modelled as the probability of the expected response (0 or 1) in terms of an unobserved hypothetical variable (a latent trait,  $\theta$ ) and items specific parameters ( $\alpha_i$  and  $\beta_i$ ):

$$P_i(\theta) = P(\beta_i, \alpha_i, \theta).$$

The latent trait  $\theta$  can, among other things, represent intelligence. In general this variable is referred to as *ability*.  $\beta_i$  is the *location parameter*, expressed in unit of  $\theta$ , that indicates the point on the ability scale at which the probability of correct response is 0.50 and  $-\infty \leq \beta_i \leq \infty$ . This parameter is also referred to as the *difficulty parameter*.  $\alpha_i$  is the *discrimination parameter*. This parameter indicates the extent to which an item  $i$  loads onto the latent trait (Baker & Kim, 2004; Van den Berg et al., 2007).

Item responses are usually given in a dichotomous (two-category) data set; the item is either answered correct or incorrect. A well-known IRT model for dichotomous data is the *Rasch model* (Rasch, 1960).

### 3.1.2 The Rasch Model

The Rasch model is one of the most widely recognized models for dichotomous data by practitioners (Baker & Kim, 2004). In this model every person  $i$  gets its own ability  $\eta_i$  and every  $j$ th item has its own difficulty parameter  $\delta_j$ . Rasch defined these symbols in such a way that when examinee 1 was twice as able as examinee 2, then  $\eta_1 = 2 \cdot \eta_2$ . In similar way, if item 1 was twice as difficult as item 2, then  $\delta_1 = 2 \cdot \delta_2$ . When this condition of Rasch is satisfied the ratio  $\frac{\eta_1}{\delta_1} = \frac{L\eta_2}{L\delta_2} = \frac{\eta_2}{\delta_2}$ , where  $L$  could be any real values constant. Thus, one should be able to establish similar ratios using arbitrary sets of people and items as long as the same value of  $L$  was involved. As a result, it is necessary to impart meaning simultaneously to the terms ability,  $\eta_i$ , and item difficulty,  $\delta_i$ . In other words,  $\eta_i$  and  $\delta_j$  need to be on the same scale (Baker & Kim, 2004). To establish similar ratios, Rasch considered the ratio  $\xi_{ij} = \frac{\eta_i}{\delta_j}$ . Next, Rasch needed to find a function of  $\xi_{ij}$  that takes on only values from 0 to 1 while  $\xi_{ij}$  goes from 0 to  $\infty$ . Rasch choose one of the most simple functions with this trait:

$$f(\xi_{ij}) = \frac{\xi_{ij}}{1 + \xi_{ij}}.$$

Which is a logistic function, a type of function known for the range between 0 and 1. It can be easily seen that substituting  $\xi_{ij} = \frac{\eta_i}{\delta_j}$ , the general term for the probability of response is given by

$$P(y_{ij}|\eta_i, \delta_j) = \frac{(\eta_i/\delta_j)^{y_{ij}}}{1 + (\eta_i/\delta_j)}. \quad (1)$$

Applying a logarithmic transformation to both the ability and item difficulty scale, the model is brought into alignment with the existing models used in IRT. This results in:

$$\begin{aligned} \theta_i &= \log(\eta_i), \\ \beta_j &= \log(\delta_j). \end{aligned}$$

The inverse transformation yields

$$\begin{aligned} \eta_i &= e_i^\theta \\ \delta_j &= e_j^\beta. \end{aligned}$$

Then, from the general term for the probability of response (Equation 1),

$$P(y_{ij}|\theta_i, \beta_j) = \frac{e^{(\theta_i - \beta_j)y_{ij}}}{1 + e^{(\theta_i - \beta_j)}}, \quad (2)$$

which is exactly the two-parameter logistic Item Characteristic Curve model with difficulty parameter  $\beta_j$ , a discrimination parameter  $\alpha_j$  fixed at one and an ability variable  $\theta_i$ . The probability of a correct response of examinee  $i$  on item  $j$  is:

$$P(y_{ij} = 1|\theta_i, \beta_j) = \frac{e^{(\theta_i - \beta_j)}}{1 + e^{(\theta_i - \beta_j)}}. \quad (3)$$

So the Rasch model uses a logistic function of the difference between the ability of a person and the difficulty of an item. When  $\theta_i$  equals  $\beta_j$ , the probability of a correct response of examinee  $i$  is 50% on item  $j$ . When  $\theta_i$  is higher than  $\beta_j$ , the probability of a correct response will be higher than 50% and vice versa when  $\theta_j$  is lower than  $\beta_i$ . Since the Rasch model requires that  $\theta_j$  and  $\beta_k$  are on the same scale, it is possible to define a scale for the latent traits and compare individuals. A scale can be defined by ordering the location parameter values in the ability scale. The lowest difficulty

parameter is then associated with the beginning of the ability scale and likewise the end of the ability scale is associated with the highest difficulty parameter.

However, before the Rasch model can be used to start comparing individuals, one must take into account the assumptions that comes with Rasch's work (Rasch, 1966):

1. The probability of correct response of examinee  $i$  to a dichotomously scored item  $j$  is given by  $P(y_{ij} = 1|\theta_i, \beta_j) = \xi_{ij}/(1 + \xi_{ij})$ , where  $\xi_{ij} = \theta_i/\beta_j$ .
2. Given the values of the parameters, all answers are stochastically independent.

The Rasch model can not be used, if those two assumptions are not satisfied in the model and the data. The logistic function in the first assumption makes sure that the function takes on only values between 0 and 1 and  $\xi_{ij}$  ensures that the ability and difficulty parameters are on the same scale. The second assumption includes the local independence assumption of item response theory where, for examinees having a given ability  $\theta_i$ , responses to the  $J$  items are independent. It also includes local independence of the responses of the examinees to a given item (Baker & Kim, 2004). Hence, if it is about items with the same difficulty level and examinees with the same abilities, then the probability of giving a correct response should be equal. If this assumption can not be made based on the items and the examinees, then the Rasch model is not an appropriate model to analyse the data.

Using the Rasch model has multiple advantages. It prevents dealing with sum scores, which could lead to biased estimates. Furthermore, the Rasch model gives the possibility to separate the influences of item difficulty and ability level on responses (Baker & Kim, 2004). Differences between persons can be assessed independent of what specific items are in the test, so response data from individuals that were tested with different test versions can be analysed in one analysis (Van den Berg et al., 2014). However, in order to do that, one needs to first estimate the differences in difficulty for all items in the test versions. This is called test linking. This is only necessary when different tests are used to measure the intelligence of the participants. The data that will be used in this study will all be gathered from the same test. Therefore there is no need for test linking and the concept of test linking will not be explained any further. The details can be found in the article of Van den Berg et al. (2004).

## 3.2 The Genetic model

### 3.2.1 The classical twin design

The most basic model to describe the phenotype of an individual is the linear model:

$$P = G + E,$$

originally introduced by Johanssen (1909).  $P$  represents the phenotype,  $G$  the genotype and  $E$  the environmental factors that influence the individuals phenotype. In this research, the phenotype is the level of intelligence of an individual and is given by  $\theta$ .

In later research, it became clear there are two different types of genes: additive and non-additive gene pairs. The former are genes that code for the same trait and both genes have an equally large effect on the phenotype. The latter are genes in which the effect of only one gene, the dominant one, has an effect on the phenotype. This phenomena is also known as the *dominance effect*. With this new knowledge, the genetic effects, as Johanssen introduced, were reduced to the additive genetics ( $A$ ) and it was concluded that each individual has a different level of sensitivity to



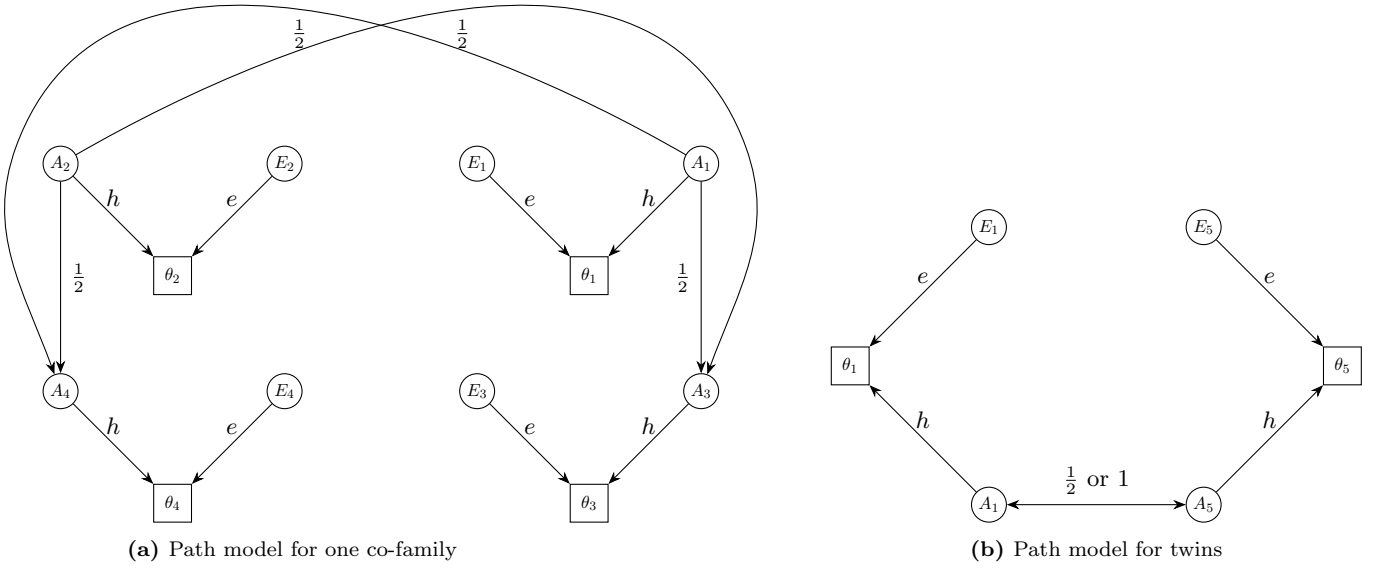
environmental effects and additive genetic effects, hence two loading factors were added to the model,  $h$  and  $e$  (Eaves et al., 1989). This leads to the following model:

$$P = h \cdot A + e \cdot E.$$

All parameters are defined relative to a total variance of unity and  $A$  and  $E$  are standardized to have unit variance (Cloninger, 1980; Eaves et al., 1989). Leading to the following equality constraint on the parameters:

$$\begin{aligned} \text{Var}(P) &= \text{Var}(hA) + \text{Var}(eE) \\ &= h^2 \text{Var}(A) + e^2 \text{Var}(E) \\ &= h^2 + e^2 = 1. \end{aligned}$$

A visual representation of this model is given in Figure 2a, where the transmission of intelligence is presented for the parents,  $\theta_1$  and  $\theta_2$  and their offspring,  $\theta_3$  and  $\theta_4$ . Both parents give half of their genes to their offspring. Figure 2b shows the relationships between twins using this model.



**Figure 2:** Path model of the transmission of intelligence.  $A_i$  represents the additive genetic value with factor loading  $h$  of person  $i$ ,  $E_i$  is the environmental value and  $e$  its factor loading.  $\theta_i$  represents the intelligence of person  $i$ . In case of monozygotic twins the expected correlations of  $A_1$  and  $A_5$  is 1 in case of dizygotic twins it equals  $\frac{1}{2}$ .

### 3.2.2 Cultural transmission

Eaves, Eysinck and Martin (1989) also found that the maternal genotype has a direct environmental impact on the phenotype of her offspring. This is called *cultural transmission* and is denoted by  $z$ . Intuitively this makes sense. Parents with a high IQ may be able to provide a more stimulating environment for their children than parents with a lower IQ. This might have an indirect effect on the intelligence level of their offspring. Since the genotype of the parents do not only effect the genotype of

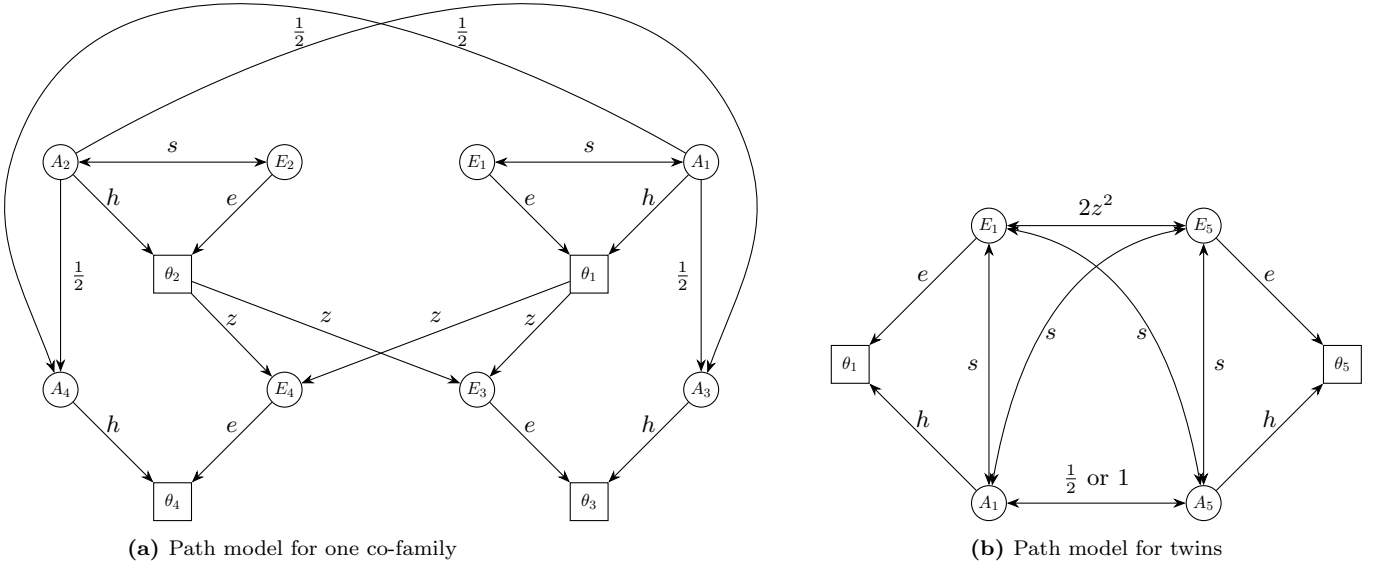
the children, but also indirectly influences the environmental effects on the children, a correlation between  $A$  and  $E$  is present and denoted by  $Cov(A, E) = s$ . With the presence of this correlation, the variance of the phenotype becomes:

$$Var(P) = h^2 Var(A) + e^2 Var(E) + 2hes.$$

And the equality constraint on the parameters changes into:

$$h^2 + e^2 + 2hes = 1.$$

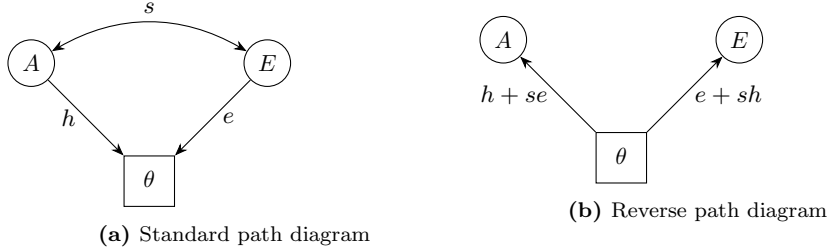
See Figure 3 for a visual representation.



**Figure 3:** Path model of the transmission of intelligence assuming the presence of cultural transmission ( $z$ ), causing a correlation between the additive genes ( $A$ ) and the environmental effects ( $E$ ), which is denoted by  $s$ .  $A_i$  represents the additive genetic value with factor loading  $h$  of person  $i$ ,  $E_i$  is the environmental value and  $e$  its factor loading.  $\theta_i$  represents the intelligence of person  $i$ . In case of monozygotic twins the expected correlations of  $A_1$  and  $A_2$  is 1 in case of dizygotic twins it equals  $\frac{1}{2}$ .

Figure 3b visualizes the relationships between twins. The correlation between  $A_1$  and  $A_2$  remains the same. However, due to cultural transmission, a correlation between the environment effects of twins arises. This correlation is due to the fact that the effect of the parents phenotype on the environmental value of twin 1 are for obvious reasons correlated to the effects of the parents phenotype on the environmental effect of twin 2. The expected correlation between the environmental effects of the twins is set to  $2z^2$ , as can be seen in Figure 3b. This correlation is derived by employing the technique of reverse path analysis, previously discussed by Li (1975), Wright (1978) and Cloninger et al. (1979). The rules of this technique require that we trace every connecting pathway between the variables in question by going backwards along paths, then forwards, but not vice versa, without going through the same arrow twice in a single pathway. A pathway thus comprises several paths. A small example is given in Figure 4.

In Figure 4a we can go from  $\theta$  to  $A$  by following the path  $h$ , and we can get back from  $A$  to  $\theta$  by following  $se$ . Hence, the expected correlation between  $\theta$  and  $A$  is  $h + se$  as is shown in Figure 4b. Likewise, the correlation between  $\theta$  and  $E$  can be found as  $e + sh$ .



**Figure 4:** Standard and reverse path diagrams with  $A$  representing additive genes,  $E$  environmental effects,  $h$  and  $e$  corresponding factor loadings and  $s$  the correlations between  $A$  and  $E$ .

Applying these techniques to Figure 3, it can be seen that the only pathway from  $E_3$  to  $E_4$ , without passing one edge more than once, is by following  $z^2 + z^2 = 2z^2$ . Note that  $\theta_3$  and  $\theta_4$  are siblings, hence they have the same relationships as dizygotic twins. However, since no paths have been used that would differ when dealing with monozygotic twins, the correlation between  $E_3$  and  $E_4$  (siblings/dizygotic twins) is the same as the correlation between  $E_1$  and  $E_5$  (monozygotic twins).

### 3.2.3 Assortative mating

Eaves, Eysenck and Martin (1989) found spousal correlations in such an extend that a general model for the transmission of intelligence in families, cannot ignore the idea of non-random mating. Therefore, the next extension of the classical twin design is a model including this non-random mating, also known as assortative mating. There are various hypothesis about the causes of this type mating. It may be due to marital interaction, phenotypic assortment, or social homogamy (Van Leeuwen et al., 2008). All three will be discussed.

#### *Marital interaction*

Marital interaction yields that when two people are married, the time spend together eventually causes a spousal correlation, since spouses would become more similar after spending a lot of time together. However, a few studies have investigated this cause of assortative mating and no significant evidence of marital interaction has been found (Van Leeuwen et al., 2008).

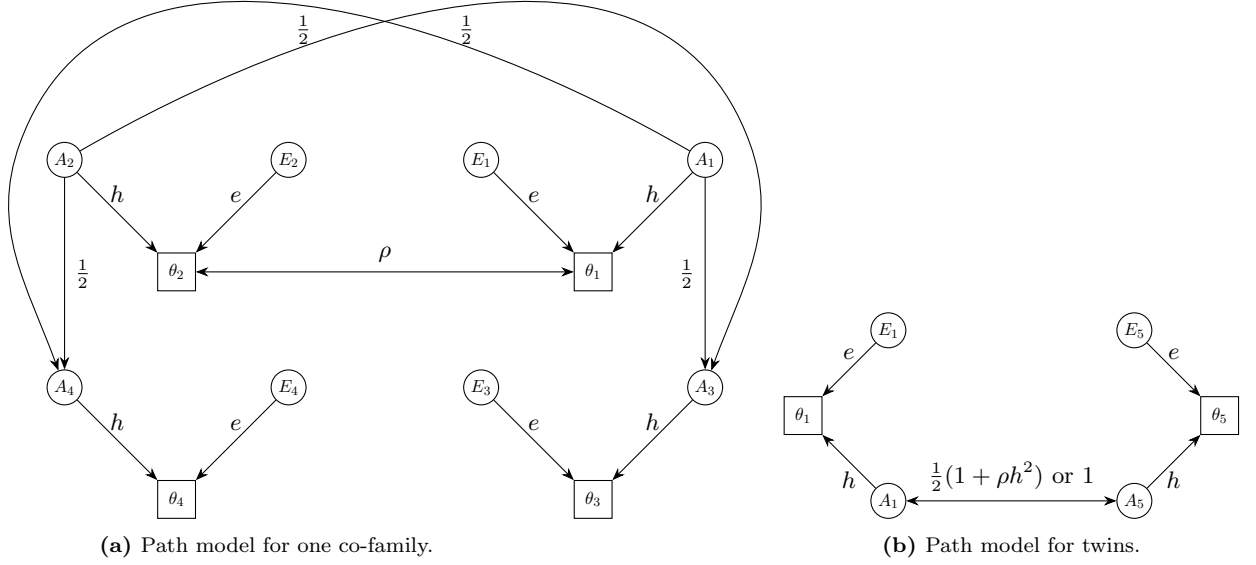
#### *Phenotypic assortment*

Another hypotheses of assortative mating is phenotypic assortment. This is the phenomena that individuals tend to select spouses on the basis of observed similar phenotype (Fulker & deFries, 1983). In more recent research van Leeuwen et al. (2008) found a spouse correlation for IQ of 0.3 confirming the findings of Eaves, Eysenck and Martin. This spouse similarity causes a correlation between the phenotype of spouses. Cloninger (1980) defined this correlation as  $\rho$ . The presence of this correlation  $\rho$ , leads to an increased genetic variance in the next generations (Eaves, et al., 1989). A visual representation of this model is very similar to the one given in Figure 2, only one extra path is added between spouses  $\theta_1$  and  $\theta_2$ . The path has a factor of  $\rho$ . Applying

the technique of reverse path analysis again, a different expected correlation between the additive genetics of dizygotic twins or siblings is found:  $\frac{1}{2}(1 + \rho^2 h)$ . Both are visualized in Figure 5. This additional path does not cause any correlation between  $A$  and  $E$ , hence

$$h^2 + e^2 = 1,$$

is the equality constraint on the parameters for this model.



**Figure 5:** Path models of the transmission of intelligence assuming the presence of phenotypic assortment  $\rho$ .  $A_i$  represents the additive genetic value with factor loading  $h$  of person  $i$ ,  $E_i$  is the environmental value and  $e$  its factor loading.  $\theta_i$  represents the intelligence of person  $i$ . In case of monozygotic twins the expected correlations of  $A_1$  and  $A_5$  is 1 in case of dizygotic twins it equals  $\frac{1}{2}(1 + \rho h^2)$ .

#### *Social homogamy*

The last hypothesis of assortative mating is social homogamy. This hypothesis states that people with the same intelligence level live in the same social environment (Van Leeuwen et al., 2008). When an individual selects a partner from the same social environment, it automatically means that both individuals have similar IQ, based on the assumption of social homogamy. When the social environment is of any influence on the intelligence, spousal correlations will occur. This social environment is denoted as  $C$  and the influence of the social environment on the level of intelligence of an individual is denoted with factor loading  $c$ .

Since  $C$  is of direct influence on the phenotype, the linear model to describe the phenotype of an individual changes:

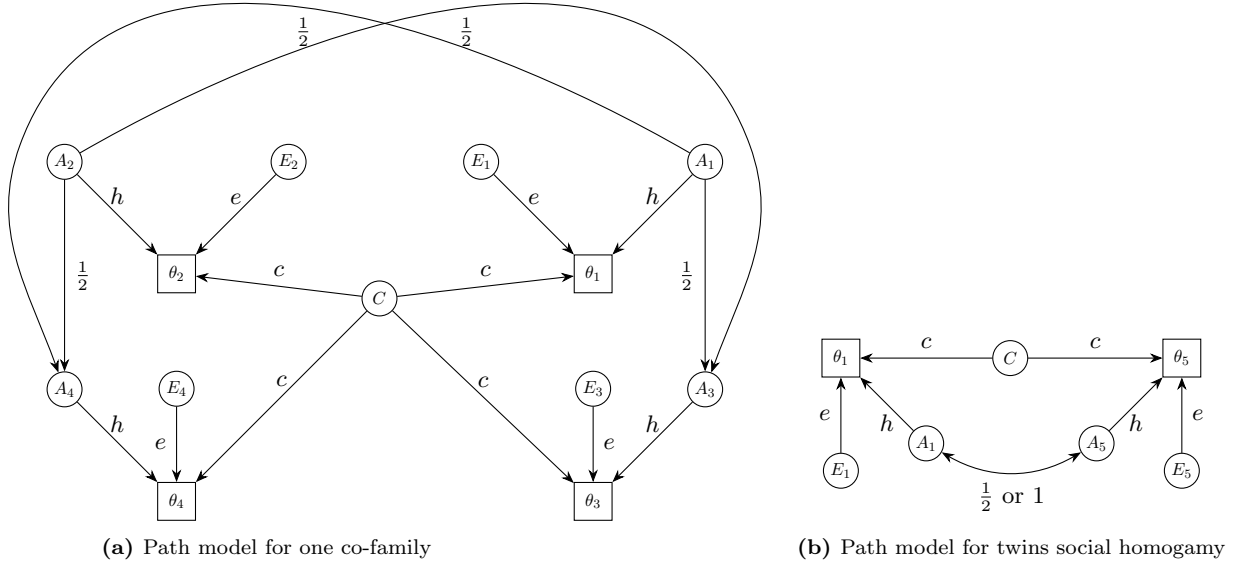
$$P = h \cdot A + e \cdot E + c \cdot C,$$

where  $C$  is standardized to have unit variance and no correlations between parameters

$A$ ,  $E$  and  $C$  occur, hence:

$$\begin{aligned}\text{Var}(\mathbf{P}) &= \text{Var}(hA) + \text{Var}(eE) + \text{Var}(cE) \\ &= h^2\text{Var}(A) + e^2\text{Var}(E) + c^2\text{Var}(C) \\ &= h^2 + e^2 + c^2 = 1.\end{aligned}$$

Figure 7 represents a visualization of this extension. Please note that in Figure 7b the twins live in the same environment. However, if the twins live in different co-families, this is not necessarily the case. If they don't,  $c$  must be fixed to zero.



**Figure 6:** Path models of the transmission of intelligence assuming social homogeneity  $C$  with an effect of  $c$ .  $A_i$  represents the additive genetic value with factor loading  $h$  of person  $i$ ,  $E_i$  is the environmental value and  $e$  its factor loading.  $\theta_i$  represents the intelligence of person  $i$ . In case of monozygotic twins the expected correlations of  $A_1$  and  $A_5$  is 1 in case of dizygotic twins it equals  $\frac{1}{2}$ .

### 3.2.4 The dominance effect

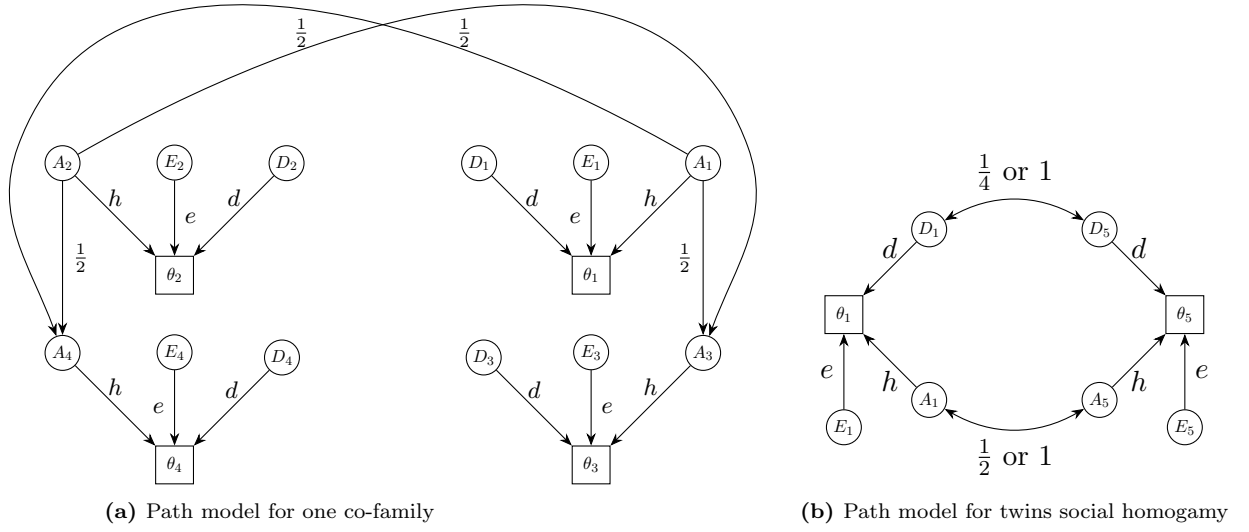
The dominance effect is the influence of the dominant genes on the phenotype of an individual. Since the dominance effects is only visible when genes are combined, parents do not give dominance effects present in their phenotype to their offspring since they only give half of their genes. Therefore, the dominance effect in the offspring does not have any correlation with the dominance effect in the phenotype of the parents. Between twins and siblings, the dominance effect are correlated. In fact, they are perfectly correlated between monozygotic twins, since they share all of their DNA and there exist a correlation of  $\frac{1}{4}$  for dizygotic twins and siblings (Van Leeuwen et al, 2008). Again, the dominance effect is of direct influence on the phenotype and therefore the linear model of the phenotype of an individual is adjusted to:

$$P = h \cdot A + e \cdot E + d \cdot D,$$

where  $D$  is standardized to have unit variance and the factor loading is  $d$ . No correlations between parameters  $A$ ,  $E$  and  $D$  occur. This implies the following equality constraint on the parameters:

$$h^2 + e^2 + d^2 = 1.$$

Figure 7 is the visual representation of this model.



**Figure 7:** Path models of the transmission of intelligence assuming the presence of the dominance effect  $D$  with an effect of  $d$ .  $A_i$  represents the additive genetic value with factor loading  $h$  of person  $i$ ,  $E_i$  is the environmental value and  $e$  its factor loading.  $\theta_i$  represents the intelligence of person  $i$ . In case of monozygotic twins the expected correlations of  $A_1$  and  $A_5$  is 1 in case of dizygotic twins it equals  $\frac{1}{2}$ .

### 3.2.5 The six models

Four extensions of the classical twin design have been discussed:

1. Presence of cultural transmission ( $z$ )
2. Presence of assortative mating:
  - (a) phenotypic assortment ( $\rho$ )
  - (b) social homogamy ( $C$ )
3. Presence of the dominance effect ( $D$ )

The four extensions are presented separately to show the influences of the extensions clearly. However, multiple of these extensions can and might be present at the same time. In this research six different models will be tested, which are all based on the classical twin design and extended with one or multiple extensions. An overview of the models and their included extensions can be found in Table 2.

The first model is the most basic model presented in Figure 2. This is the least complex model one can define for the transmission of intelligence. It uses the genetic model from the classical twin design without adding any extensions. In the next model

only the presence of phenotypic assortment was assumed (i.e., Model 2), see Figure 5. This was the preferred model according to Van Leeuwen et al., (2008). Model 3 assumes the presence of both phenotypic assortment and cultural transmission. Veldkamp et al. (unpublished) found negative values for the cultural transmission in model 3. This would mean that parents with a lower intelligence level would be better in providing a stimulating environment than parents with a higher intelligence level. Intuitively this does not make sense and they tested if this issue could be resolved by adding the influences of the dominance effect. Therefore, it was decided to test model 4 as well, which includes phenotypic assortment, cultural transmission and the dominance effect. Model 5 assumes instead of phenotypic assortment, social homogamy as a cause of assortative mating. In model 6, model 5 is extended with the presence of the dominance effect.

More details on the defined models are all summarized in Appendix A. Visual representations of the models can there be found, together with tables in which the expected correlations between the eight family members are given, based on the genetic model associated with the defined model. Deriving the visual representation of the six models is very straightforward. The visual representation from Section 3.2 can be used and additional paths can be added. However, the derivation of the expected correlation is quite technical. For every model the expected correlations are derived by using the technique of reverse path analysis and compared to the expected correlation used in previous research. However, this family structure has not been used in previous research and therefore not all expected correlations could be taken from previous literature and can not be checked with works of other researchers. Furthermore, The constraints per model on the parameters are given, with technical details if needed. Thus, for more details on the six models, please see Appendix A.

Only one interesting phenomena related to model 3 and model 4 is discussed in this section, to show that combining extensions might be more complex then it seems. As stated before, by assuming the presence of phenotypic assortment, the phenotypes of the parents become correlated. This yields a higher genetic correlation in their offspring. Moreover, assuming cultural transmission a correlation,  $s$ , arises between an individuals additive genes  $A$  and environmental influences  $E$ . When both of these extensions are constant, research has shown that over several generations, the correlation between  $A$  and  $E$  approaches an equilibrium (Eaves, Eysinck & Martin, 1989). In other words, the  $Cov(A,E)$  will reach a stable value. Therefore, it is assumed that this correlation between  $A$  and  $E$  is equal across generations:

$$s = Cov_{i=1,2}(A_i, E_i) = Cov_{j=3,4}(A_j, E_j).$$

Using reverse path analysis on model 3, the following equation can be found:

$$\begin{aligned} s = Cov_{j=3,4}(A_j, E_j) &= z(h + eCov_{i=1,2}(A_i, E_i))(1 + \rho) \\ &= z(h + es)(1 + \rho) \\ &= \dots \\ &= \frac{(1 + \rho)zh}{1 - (1 + \rho)ze}. \end{aligned}$$

This expression is important when we come to fit the model, because it represents a second constraint that must be imposed on the parameter estimates if the requirements of equilibrium are to be satisfied. Second to the first constraint which stated that all parameters are defined relative to a total variance of unity.

Due to the phenotypic correlation between parents,  $\rho$ , the genotypes and environmental factors of the parents become correlated as well. The correlation between the

genotypic values in the parents is defined as  $\gamma$ , the correlation between environmental factors in the parents as  $\epsilon$  and the correlation between genotype parent 1 and environmental factors parent 2 as  $\delta$ . The same rules for path analysis are employed, which lead to the following equations:

$$\begin{aligned}\gamma &= \rho(h + se)^2 \\ \epsilon &= \rho(e + sh)^2 \\ \delta &= \rho(h + se)(e + sh)\end{aligned}$$

(Wright, 1968; Cloninger, 1980; Fulker & deFries, 1983). Note that these correlations only apply to the additive genetic and environmental effects of the parents and not of the offspring.

In conclusion, for all models a constraint on the parameter values is implied by the fact that the parameters are defined relative to a total variance of unity and  $A$  and  $E$  are standardized to have unity variance (Cloninger, 1980; Eaves, Eysinck & Martin, 1989). Hence:

$$Var(P) = 1$$

is a general constraint on the genetic model. When dealing with both cultural transmission and phenotypic assortment, the correlation  $s$  reaches an equilibrium which forms a second constraint:

$$s = \frac{(1 + \rho)zh}{1 - (1 + \rho)ze}.$$

In Appendix A the implementation of these two constraints is discussed.

**Table 2:** Overview of the six models. The crosses indicate which extensions are included to the classical twin design and thus how the genetic model is defined.

	Cultural transmission	Phenotypic assortment	Social homogeneity	Dominance effect
Model 1				
Model 2		x		
Model 3	x	x		
Model 4	x	x		x
Model 5			x	
Model 6			x	x

### 3.3 The Cholesky Decomposition

The genetic models are all examples of structural equation models (SEM). A key property to such models is the possibility to return to the start node when moving along paths. Veldkamp, Schwabe and Van den Berg (unpublished) showed that a re-specification of the genetic models into directed acyclic graphs (DAG), make the models more easy to access for applied researchers. Rewriting them as DAG models gives the opportunity to use off-the-shelf software to estimate the models. A DAG is a model which, in contrast to a SEM, has as key property the impossibility to return to the start node when moving along paths. They proposed using a Cholesky decomposition to rewrite the SEM models as DAG models.

Any positive definite and symmetric matrix  $A$  can be expressed in the form  $A = X^T X$  for a non-singular matrix  $X$ . The Cholesky factorization is a particular form of this factorization in which  $X$  is upper triangular with positive diagonal elements, defined as  $\Lambda$ ;  $A = \Lambda^T \Lambda$ . The Cholesky decomposition can be obtained using the



following general formulas:

$$\Lambda_{i,j} = \begin{cases} \lambda_{i,j} = 0 & \text{for } i > j, \\ \lambda_{i,i} = \sqrt{A_{j,j} - \sum_{k=1}^{j-1} \lambda_{j,k}^2} & \text{for } i = j, \\ \lambda_{i,j} = \frac{1}{\lambda_{j,j}} \left( A_{i,j} - \sum_{k=1}^{j-1} \lambda_{i,k} \lambda_{j,k} \right) & \text{for } i < j, \end{cases}$$

where  $i$  and  $j$  represent the rows and columns from matrices  $\Lambda$  and  $A$ . In this research the Cholesky decomposition is applied to the standardized covariance matrix,  $\Sigma$ .  $\Sigma$  is defined for all participating families and is a  $J \times J$  matrix, representing the correlations between all  $J$  family members. For each genetic model, expected correlations between family members were defined.  $\Sigma_{ij}$  represents the correlation between the level of intelligence of family members  $i$  and  $j$ .  $\Sigma$  is defined as:

$$\Sigma = \begin{bmatrix} 1 & \sigma_1 & \sigma_4 & \sigma_4 & \sigma_5 & \sigma_6 & \sigma_8 & \sigma_8 \\ \sigma_1 & 1 & \sigma_4 & \sigma_4 & \sigma_{3*} & \sigma_5 & \sigma_7 & \sigma_7 \\ \sigma_4 & \sigma_4 & 1 & \sigma_2 & \sigma_7 & \sigma_8 & \sigma_9 & \sigma_9 \\ \sigma_4 & \sigma_4 & \sigma_2 & 1 & \sigma_7 & \sigma_8 & \sigma_9 & \sigma_9 \\ \sigma_5 & \sigma_3 & \sigma_7 & \sigma_7 & 1 & \sigma_1 & \sigma_4 & \sigma_4 \\ \sigma_6 & \sigma_5 & \sigma_8 & \sigma_8 & \sigma_1 & 1 & \sigma_4 & \sigma_4 \\ \sigma_8 & \sigma_7 & \sigma_9 & \sigma_9 & \sigma_4 & \sigma_4 & 1 & \sigma_2 \\ \sigma_8 & \sigma_7 & \sigma_9 & \sigma_9 & \sigma_4 & \sigma_4 & \sigma_2 & 1 \end{bmatrix}.$$

As  $\Sigma$  shows, every  $\Sigma_{ij}$  is defined as  $\sigma_k$  with  $k = 1, \dots, 9$ .  $\sigma_k$ , where  $\sigma_k$  shows what kind of relationship there is between family member  $i$  and  $j$ . The relationships are defined in Table 3.

**Table 3:** Definitions of the type of relationships between family members represented by  $\sigma_k$ .

Abbreviations	Relationship
$\sigma_1$	Spouses
$\sigma_2$	Siblings/Dizygotic twins
$\sigma_3$	Monozygotic twins
$\sigma_4$	Parent-Offspring
$\sigma_5$	Sibling in law
$\sigma_6$	Spouses of twins
$\sigma_7$	Nephew/Niece with co-twin of twin parent
$\sigma_8$	Nephew/Niece with spouse of co-twin
$\sigma_9$	Cousins

The level of intelligence for family  $i$  can be defined as an  $J \times 1$  vector  $\theta_i \sim MVN(\mathbf{0}, \Sigma)$ .  $\theta_{ij}$  represents the level of intelligence of family member  $j$  from family  $i$ .  $\Sigma$  ensures that the correlations between the level of intelligence of all  $J$  family members, defined based on a particular genetic model, are taken into account. If the covariance matrix is symmetric and positive definite, there exists an upper triangular matrix  $\Lambda$  such that:

$$\Sigma = \Lambda^T \Lambda.$$

By definition:

$$\begin{aligned}
Cov(\theta_i) &= \Sigma \\
&= \Lambda^T \Lambda \\
&= \Lambda^T I \Lambda \\
&= \Lambda^T Cov(\mathbf{v}_i) \Lambda, \quad \text{if } \mathbf{v}_i \sim MVN(\mathbf{0}, I) \\
&= Cov(\Lambda \mathbf{v}_i).
\end{aligned}$$

Therefore, for each  $j$ th family member, a standard normal distributed auxiliary variable,  $v_{ij}$ , is specified, yielding:

$$\mathbf{v}_i \sim MVN(\mathbf{0}, I),$$

where  $i$  represents the family ( $i = 1, 2, \dots, N$ ),  $\mathbf{0}$  is a  $J \times 1$  zero vector and  $I$  is a  $J \times J$  identity matrix. The phenotype of every family  $i$  is then defined as:

$$\theta_i = \Lambda \mathbf{v}_i.$$

### 3.4 A Bayesian Framework

#### 3.4.1 The Bayesian Framework

In order to estimate the measurement model and the genetic model simultaneously, Bayesian statistical modelling is used. In this framework, parameter spaces are obtained from the posterior density of the model parameters,  $P(\delta|Y)$ , where  $\delta$  is the set of model parameters and  $Y$  the observed data. By Bayes' theorem, this density is proportional to the product of the likelihood function,  $P(Y|\delta)$  and the prior distributions,  $P(\delta)$ :

$$P(\delta|Y) \propto P(Y|\delta) \cdot P(\delta).$$

The prior distributions, for all possible parameters in the genetic models, are defined as in Table 4 and the likelihood function is Bernoulli distributed with a probability of having a correct answer, drawn from the Rasch model, illustrated in Equation 3. The boundaries set on parameter  $z$  are from the research of Veldkamp et al. (unpublished). In previous research they concluded that parameter  $z$  did not always converge, leading to biased estimates of  $h^2$  and  $s$ . Therefore, they proposed two inequality constraints:  $h^2 < 1$  and  $s < 1$ . These two equations were solved for  $z$  leading to the lower and upper boundary given in Table 4.

**Table 4:** Prior distributions used in the six models defined in Section 3.2.5.

Parameter	Prior Distribution
$e$	$Beta(-1, 1)$
$\rho$	$Beta(-1, 1)$
$z$	$U(0, 1)$ with boundaries $[\frac{e}{(e^2-2)(1+\rho)}, \frac{1}{1+\rho}]$
$d$	$Beta(-1, 1)$
$c$	$Beta(-1, 1)$

Van den Berg et al. (2007), Van den Berg (2009) and Otermann and Van den Berg (2016) have used off-the-shelf Bayesian software packages such as JAGS and openBUGS. Such packages however, tend to work like black-boxes and therefore a Metropolis Hastings (MH) algorithm was implemented. This was previously done by Nolle (2018) during her bachelor project. To avoid using black-boxes and not being

fully aware of how the models behave, this research will also make use of a MH-algorithm. An R package was written containing two functions. The first function is to simulate data. The second function executes a MH-algorithm in which the fit of a specified genetic model is measured with respect to the given data set. In Appendix B the information about the package when using the `help()` function in R is given to show how the functions can be used.

### 3.4.2 Model Fit

To compare the six different models the deviance information criterion (DIC) of Spiegelhalter et al. (2002) was calculated. This criterion takes both the fit of the data to the model and the model complexity into account and it is defined as:

$$DIC = \text{'goodness of fit'} + \text{'complexity'}.$$

The model fit is measured via the deviance, also known as the log likelihood statistic:

$$D(\delta) = -2\log(P(y|\delta)),$$

where  $y$  is the observed data and  $\delta$  the set of parameters. The complexity is measured by estimating the effective number of parameters which is the difference between the posterior mean of the deviance and the deviance evaluated at the posterior mean of the parameters:

$$P_D = \overline{D(\delta)} - D(\bar{\delta}).$$

The DIC is then defined as:

$$DIC = \overline{D(\delta)} + P_D.$$

## 4 Simulated data

Unfortunately, no real-world data was available for this study. The data set used by Reynolds et al. (1996) contains families with the structure needed for this study. However, it was not possible to get access to this data set. An attempt was made to get access to this data through Stephanie van den Berg, from the university of Twente, who has used this data set herself in her own studies. But due to time restrictions, there was no option to wait for those any longer. Therefore the six models were only compared to each other based on simulated data sets.

In total, six different data sets have been simulated, all based on one of the genetic models. Item responses were simulated for 400 families containing 8 family members each. In total 45 item responses were simulated per individual. A difficulty parameter for each of the 45 items is sampled from the normal distribution. Intelligence levels were per family drawn from a multivariate normal distribution with  $\Sigma$  as variance.  $\Sigma$  ensures the expected correlations between family members are taken into account. Item responses were drawn from a binomial distribution, with a probability of having a success equal to Equation 3.

In order to define  $\Sigma$ , the parameters must be fixed. Reasonable parameter values are chosen based on the results of previous simulation studies from Veldkamp et al. (unpublished) and are denoted in Table 5.

**Table 5:** An overview of the fixed parameter values used to simulate 6 data sets. Each data set is simulated using the covariance matrix  $\Sigma$  associated with one of the six models defined in section 3.2.5. — means that the parameter does not apply to the  $\Sigma$  used for this data set.

	$\Sigma$ from model	$e$	$\rho$	$z$	$d$	$c$	$N$	$K$
Data set 1	1	0.6	-	-	-	-	400	45
Data set 2	2	0.76	0.33	-	-	-	400	45
Data set 3	3	0.63	0.33	0.05	-	-	400	45
Data set 4	4	0.31	0.37	-0.03	0.83	-	400	45
Data set 5	5	0.1	-	-	-	0.79	400	45
Data set 6	6	0.05	-	-	0.85	0.21	400	45

## 5 Results

### 5.1 Simulations

The Metropolis Hasting algorithm inside the R package generates samples from the marginal posterior distributions of the parameters of a model. With a trace plot the sampled values of a parameter are visualized over a number of iterations and it can be easily seen if and when the Markov Chain has converged to the stationary distribution. For calculating the DIC only samples from the stationary distribution are used. A burn-in phase is initialized, which withdraws the first states who have not yet reached the stationary distribution. The burn-in phase and the number of iterations are determined based on analyzing model 3. Model 3 contains three parameters that need to be estimated:  $e$  (environmental effects),  $\rho$  (phenotypic correlations) and  $z$  (cultural transmission).

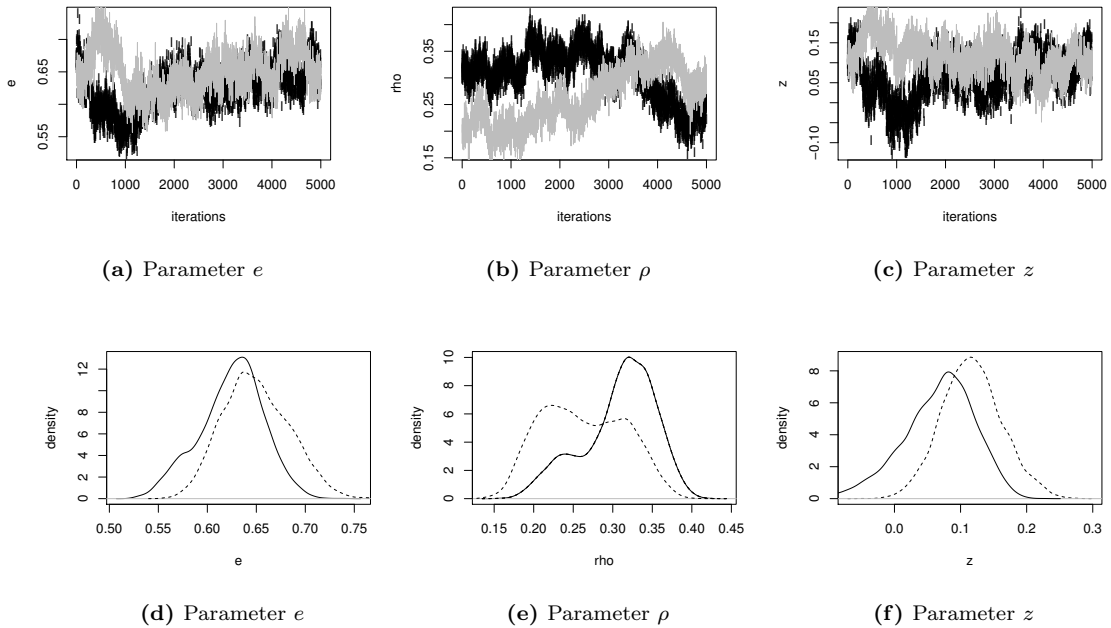
Figures 8-10 show multiple trace and density plots for parameters  $e$ ,  $\rho$  and  $z$ . The simulations are performed based on the simulated data set 3 described in Section 4 and Table 5. In total 10 Markov Chains are conducted, where each chain started with randomly drawn starting values between 0 and 1, while satisfying the two constraints described in Section 3.2.5 and Appendix A.3. The first two chains are shown in Figure 8, where 1000.000 iterations are performed and a burn-in phase of 500.000 is used. Saving the parameter values for all these iterations would take too much memory, therefore a thinning factor of 100 is introduced. Meaning only every 100-th state is saved. One simulation of 1000.000 iterations takes about 20 hours. Fortunately the Peregrine HPC cluster of the University of Groningen was made available, which was used to run multiple simulations simultaneously.

The first two chains in Figure 8 do not show very convincing convergence. Especially the two chains for parameter  $\rho$  seem to follow different paths, hence the convergence is quite weak. This might be due to the fact that only 5000 states are given in the figure, as a consequence of both the burn-in phase and the thinning factor. Increasing the length of the chain might show a more convincing convergence and results in a better well-balanced coverage of the posterior parameter space. Therefore, the iterations were raised to 2000.000 iterations, with still a thinning factor of a 100. The results are shown in Figures 9, 10 and 15 (Appendix C). These figures show seven different simulations and their corresponding estimates of the parameters. Six from the seven chains show very similar behavior, which argues for a strong convergence. The similar behavior can be seen from the similar paths the chains show in Figure 9 and 15 (Appendix C) and the almost identical density plots. The shape of a normal distribution can be recognized in the density plots. However, some of them are left or right skewed. This is a result of deleting the burn-in phase. In the trace plots it can be seen that before the stationary distribution is reached, the parameters first tend

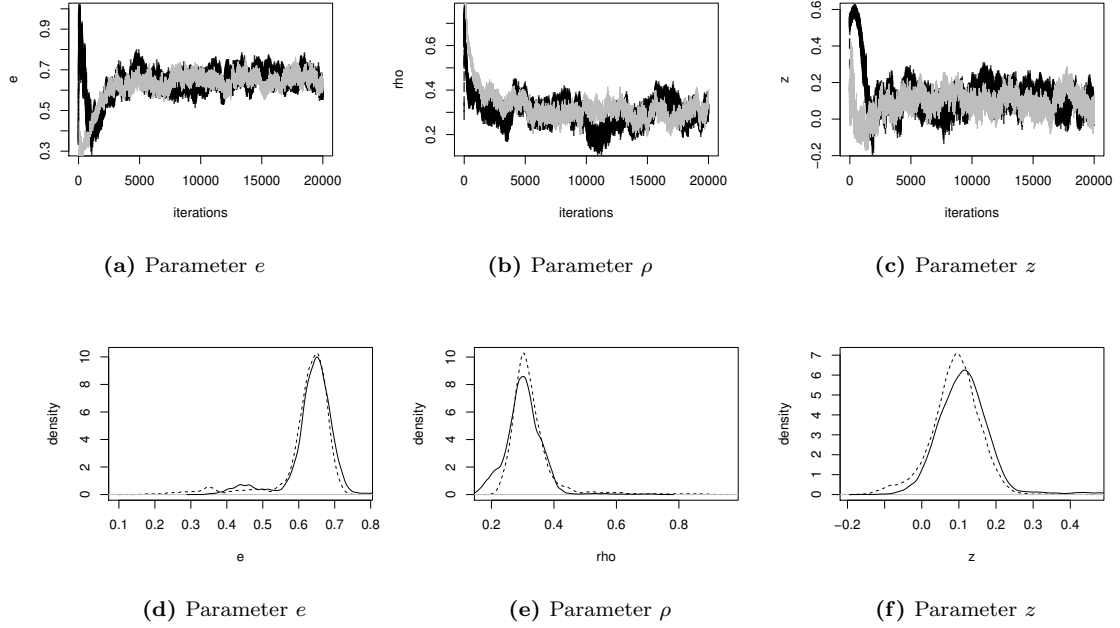
to some extreme values (0 or 1), which results in a slightly biased density plot. In Table 6 the posterior means of each chain are given and it can be seen that the chains converged closely to the values that were used to simulate the data.

Figure 10 shows the only chain with different behavior than the others. The three parameters seem to have converged, but at very different values than the other chains. The estimated parameter values for Run 2 in Table 6 clearly show different values than the other runs. However, the likelihood of the data based on these parameter values is very similar to the likelihood found using the parameter values from the other runs. This might indicate an identification problem. Identification problems arises when different sets of parameter values fit the model equally well. If this occurs, preference of one set of parameter values over other ones is completely arbitrary and can lead to wrongfully determined conclusions (Bekker, P.A., Merckens, A. & Wansbeek, T.J., 1994). In this research this particular problem will not be targeted, but one must be aware of the possible existence of this phenomena during the simulations.

In conclusion, the simulations on model 3 has shown that 2000.0000 iterations are needed to get a well-balanced coverage of the posterior parameter space, with a burn-in of 500.000 iterations and for limiting the memory space that is needed a thinning factor of 100 is introduced. These decisions are made based on 10 different Markov Chains. Ideally this amount of chains would be much higher and it would have been tested for all models individually. However, 1 simulation takes about approximately 48 hours and since this research has a time restriction the number of simulations has to be limited. Hence, the conclusions based on model 3 will be applied to all six models.



**Figure 8:** Trace plots (8a-8c) and density plots (8d-8f) from two MH-MCMC chains with 1000.000 iterations, with a thinning factor of 100 and a burn-in phase of 500.000 resulting in 5000 states given in the plots. The solid line in the density plots corresponds to the dark chain in the trace plots. The dotted line corresponds to the lighter chain in the trace plots.



**Figure 9:** Trace plots (9a-9c) and density plots (9d-9f) from two MH-MCMC chains with 2000.000 iterations, with a thinning factor of 100 and no burn-in resulting in 20.000 states given in the plots. The solid line in the density plots corresponds to the dark chain in the trace plots. The dotted lines corresponds to the lighter chains in the trace plots.

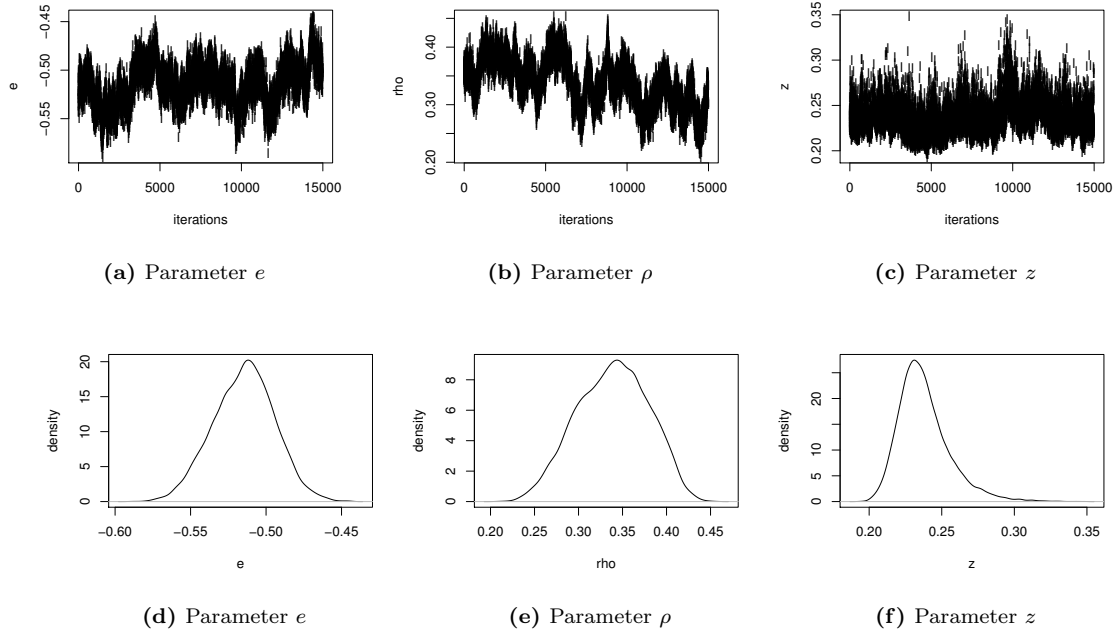
**Table 6:** Parameter estimates (posterior mean) for model 3 assuming cultural transmission and phenotypic assortment. Results from 5 different runs are shown and an average of all 10 runs is given together with the simulated values that were used to simulate the data.

Parameter	Simulated values	All chains (10)	chain 1	chain 2	chain 3	chain 4	chain 5
$h$	0.74	0.73	0.72	1.0	0.69	0.70	0.67
$e$	0.63	0.64	0.64	-0.51	0.65	0.65	0.63
$\rho$	0.33	0.30	0.30	0.33	0.28	0.30	0.32
$z$	0.05	0.11	0.08	0.24	0.12	0.11	0.12
$s$	0.05	0.12	0.08	0.28	0.12	0.11	0.12

## 5.2 Model Fit

In Table 7 a summary is given of the DIC values for the different models representing their fit to the six different data sets. The table shows very similar DIC values, making it quite difficult to make statements about the fit of the models. In the work of Otermann (Master's Thesis) and Van Leeuwen et al. (2008) similar DIC values for different models were found as well. This might indicate that the DIC is not an appropriate information criterion for these models or that the six models all fit the data equally well.

In Appendix D a table per analyzed data set is given in which the estimated parameter values and the standard deviations are given for all parameters according to the six different models. Almost all of the chains from the simulation study have converged, which was checked by plotting trace plots and can be confirmed by low



**Figure 10:** Trace plots (10a-10c) and density plots (10d-10f) from one MH-MCMC chain with different behavior. The chain did run for 2000.000 iterations, with a thinning factor of 100 and a burn-in phase of 500.000 resulting in 15.000 states given in the plots.

standard deviation values. Only the simulations in which model 4 and 6 were used, models including the dominance effect, showed a weaker convergence in their parameter estimates. This can also be seen in the higher values of the standard deviations associated with the parameter estimates from these models. This suggests that for these models more iterations are needed, with a higher burn-in phase as well in order to get a well-balanced coverage of the stationary distribution.

Appendix D shows that the parameter estimates are close to the true values, provided that no identification problem occurs. Parameter  $e$  seems to occasionally converge to a high negative value, which would implicate that the environmental effects are of negative influence on an individuals intelligence. This is very unlikely and seems to be the result of an identificatin problem. Therefore in further research  $e^2$  should be estimated instead of  $e$  to ensure that  $e$  stays positive. When  $e$  is estimated to be positive, it seems that it is close to the true parameter values used to simulate the data. The results in Appendix D are all based on one chain per parameter. To obtain a better coverage of the parameter space, more chains per parameter are preferred, all with different starting values.

## 6 Conclusion/Discussion

In this thesis, a comparison study has been performed to test which genetic model represents the transmission of intelligence within families the most accurate with respect to the real-world. In order to measure the fit of the genetic models to the data a new improved methodology is used. This improved methodology is based on a literature study that investigated the methodology of the classical twin design. The new method-

**Table 7:** DIC values for all six models based on the six data sets described in Section 4. For each data set the lowest DIC value is printed in bold.

	Data 1	Data 2	Data 3	Data 4	Data 5	Data 6
Model 1	<b>508897</b>	<b>503892</b>	<b>514226</b>	527188	<b>495217</b>	507433
Model 2	508937	503901	514370	527216	495547	507387
Model 3	508947	504243	514607	<b>527139</b>	495455	507197
Model 4	509138	503899	514276	527438	507483	<b>495785</b>
Model 5	508966	503929	514322	527230	495433	507355
Model 6	508940	503918	514242	527166	495348	507518

ology includes three improvements. Firstly, a measurement model is introduced, which avoids the use of sum scores that could lead to biased estimates (Van den Berg, Glas & Boomsma, 2007). Secondly, a Bayesian framework is used to simultaneously estimate the measurement model and the genetic model (Van den Berg et al., 2007; Van Leeuwen et al., 2008; Schwabe et al., 2014). Finally, the methodology is improved by using a Cholesky decomposition to redefine the genetic model as a Directed Acyclic Graph, which gives applied researches the opportunity to use off-the-shelf Bayesian software programs. This method was introduced by Veldkamp et al. (unpublished) with the aim to make the models and their analysis easily accessible.

This research does not make use of the off-the-shelf software programs. Instead, an R package is written to execute the simulations. Veldkamp et al. (unpublished) emphasize that their method is flexible; only expected correlations have to be changed if one wants to estimate an alternative model. However, deriving those expected correlations for alternative models can be quite a challenge. Moreover, such off-the-shelf software programs tend to work like black boxes. Therefore, the R package written for this research contains two function: one to simulate data and another in which a Metropolis Hastings algorithm is executed in which the fit of a genetic model is measured with respect to the given data set. The genetic model that is used in the simulations can be chosen from six different models, defined in Section 4.1. This gives applied researchers the opportunity to use the R package for simulations, without the need of defining the genetic model and the expected correlations between family members first. The codes can however still be adjusted to different genetic models when needed. The R package prevents using off-the-shelf software programs as black boxes and is completely transparent.

The R package was used to perform a simulation study in which the six models are compared. The simulation results show that the models are able to estimate the simulated values closely (Appendix D). Only models containing the dominance effect showed some imprecise estimates, with a higher standard deviation than was seen in other simulations. Therefore it would be recommended to increase the number of iterations for the models including the dominance effect, to obtain a better coverage of the parameter space. For analysing real-world data, it would be recommended to increase the number of iterations for the genetic models in general. Not all estimations of the parameters were close to the true values. An identification problem appears to arise which especially influences the values of parameter  $e$ , causing it to obtain negative values. Negative values for  $e$  are very unrealistic; it would imply a negative correlation between the environmental effects and an individuals intelligence level. Therefore  $e$  could be restricted to only positive values by estimating parameter  $e^2$  instead of  $e$  in future research. Veldkamp et al. (unpublished) used this trick and seem to have solved the identification problem with it.

In order to directly compare the fit of the six different models, DIC values were



calculated (Table 7). These DIC values are all gathered from one simulation per DIC value. In order to get a better insight in the DIC values of the models the DIC values should be calculated multiple times for each model and an average could be taken. Since Table 7 took about three weeks to obtain, this was not a realistic goal for this research. But for further research it would be recommended to base the DIC values on multiple simulations instead of one. The DIC values conducted from the simulations in this study show very similar values. Two reasons for these similar values could be that either the DIC is not an appropriate information criterion for these models or that the differences in the model fits are too small to detect. Similar DIC values have been seen in earlier comparison studies as well, in which other methods were used but similar genetic models are compared. Van Leeuwen et al. (2008) tested two different hypothesis of the cause of spousal correlation. They directly compared the basic classical twin design (model 1), the model assuming phenotypic assortment (model 2) and the model assuming social homogamy (model 5), by calculating the AIC values and the likelihood. The model assuming phenotypic assortment appeared to be the superior model, but this was only based on very small difference in likelihood and therefore the superiority was not very convincing. Otermann and Van den Berg (2014) directly compared model 2 (phenotypic assortment), 4 (phenotypic assortment, cultural transmission, dominance effect) and a genetic model including both phenotypic assortment and the dominance effect. In their outcomes only small differences in DIC values were seen as well. Based on an analysis of real-world data, Otermann and Van den Berg could conclude that the most simple model, only including phenotypic assortment (model 2), was the best fit. Unfortunately, due to the lack of a real-world data set, no further distinction between the fit of the genetic models can be made in this research. But suggestions for further research can be made.

In future research, firstly the fit of the models could be further investigated by performing a similar comparison study on a real-world data set in which  $e^2$  is estimated instead of  $e$ . One must check if there is still an identification problem present after this adjustment. If so, it should be investigated how this problem can be avoided to prevent wrongful interpretations of the estimated parameters. Secondly, instead of using the DIC to compare the models, a different information criterion can be used or posterior predictive checks can be performed (Otermann & Van den Berg, 2014). In these checks a new data set is simulated under the fitted model and then compared to the original (old) data set. If the simulated data is similar to the observed data, the model is a good fit. Finally, Otermann and Van den Berg (2014) showed that the test circumstances for children can differ from the test circumstances for parents. Including such differences in the model by using different  $e$  values for the parents and offspring lead to a better fitting model. This could also be applied to the models used in this research, leading to better fitting models.

### Acknowledgements

I would like to thank dr. M.A. Grzegorzczuk for his guidance during this research project. He was willing to make time for meetings on a weekly basis and provided useful feedback during the process of this research. He helped me to decide which way this research had to go, for which I am very grateful. I would also like to thank my second supervisor dr. W.P. Krijnen for his thoughts on various topics during this research and taking the time to read and evaluate my work. Finally, I would like to thank the Center of Information Technology of the University of Groningen for providing access to the Peregrine high performance computing cluster which made the simulation study possible.

## 7 Bibliography

- Baker, F. & Kim, S.-H. (2004). *Item response theory: Parameter estimation techniques*. New York: Marcel Dekker.
- Bekker, P., Merckens, A., & Wansbeek, T. (1993). *Identification, equivalent models, and computer algebra (Statistical modeling and decision science)*. Boston, Mass., etc.: Academic Press.
- Cloninger, C. R. (1980). Interpretation of intrinsic and extrinsic structural relations by path analysis: Theory and applications to assortative mating. *Genetical Research*, 36(2), 133-145.
- Eaves, L. J., Eysenck, H. J., & Martin, N. G. (1989). *Genes, culture, and personality: An empirical approach*. New York: Academic Press.
- Fox, J.P., & Glas, C.A. (2001). Bayesian estimation of a multilevel IRT model using Gibbs sampling. *Psychometrika*, 66(2), 271-288.
- Fragoso, T.M., Giolo, S.R., Pereira, A.C., de Andrade, M., & Soler, J.M. (2014). Using item response theory to model multiple phenotypes and their joint heritability in family data. *Genetic epidemiology*, 38(2), 152-161.
- Fulker, D.W., & DeFries, J.C., (1983). Genetic and environmental transmission in the Colorado Adoption Project: Path analysis. *British Journal of Mathematical and Statistical Psychology*, 36(2), 175-188.
- Heath, A., & Eaves, L. J. (1985). Resolving the effects of phenotype and social background on mate selection. *Behavior Genetics*, 15(1), 15-30.
- Higham, N. (2009). Cholesky factorization. *Wiley Interdisciplinary Reviews: Computational Statistics*, 1(2), 251-254. doi:10.1002/wics.18
- Holland, P.W., & Wainer, H. (2012). *Differential item functioning*. Hoboken: Taylor and Francis.
- Neale, M.C., & Cardon, L.R. (1992). *Methodology for Genetic Studies of Twins and Families*. Dordrecht: Kluwer Academic.
- Nolle, G. (2018). Modeling the transmission of intelligence. *University of Groningen: Bachelor project*.
- Otermann, B., & Van den Berg, S.M. (2014). Linking the Standard and Advanced Raven Progressive Matrices tests to model intelligence covariance in twin families. *University of Twente*.
- Otermann, B., & Van den Berg, S. (2014). Bayesian Assessment of parent-offspring models using item data. *University of Twente: Master thesis*.
- Raven, J. (2000). The Raven's progressive matrices: Change and stability over culture and time. *Cognitive Psychology*, 41(1), 1-48.
- Rasch, G. (1960). *Probabilistic models for some intelligence and attainment tests*. Copenhagen: Danish Institute for Educational Research.
- Schwabe, I., & Van den Berg, S.M. (2014). Assessing genotype by environment interaction in case of heterogeneous measurement error. *Behavior Genetics*, 1-13.

- Spiegelhalter, D.J., Best, N.G., Carlin, B.P., & Van der Linde, A. (2002). Bayesian measures of model complexity and fit. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 64(4), 583-639.
- Van den Berg, S.M., Glas, C.A.W., & Boomsma, D.I. (2007) Variance Decomposition Using an IRT Measurement Model. *Behavior Genetics*, 37(4), 604-616.
- Van den Berg, S.M. (2009). Imposing Nonlinear Constraints When Estimating Genetic and Cultural Transmission Under Assortative Mating: A Simulation Study Using Mx and BUGS. *Behavior Genetics*, 39(1), 123-131.
- Van Leeuwen, M., Van den Berg, S.M., & Boomsma, D.I. (2008). A twin-family study of general IQ. *Learning and Individual Differences*, 18(1), 76-88.
- Veldkamp, S.A.M., Schwabe, I., & Van den Berg, S.M. (20XX). *A General Bayesian Method for Analysing Extended Twin Family Data Using Cholesky Factorization* (Unpublished article).
- Wright, S. (1968). *Evolution and the genetics of populations*. Chicago: the University of Chicago Press.

## A Genetic Models

The tables in this Appendix show the expected correlation between all eight family members. The expected correlations are derived by using the rules of reverse path analysis and were checked by comparing the equations with the ones derived by Eaves et al. (1989) and Van Leeuwen et al. (2008). For model 5 and 6 however, Van Leeuwen et al. and Eaves et al. did not show the same expected correlations for model 5 and 6. By using the rules of reverse path analysis it was decided to follow the expected correlations as described by Eaves et al. (1989). They did not define the correlations for all eight family members, the tables are shown to complete the formula sets.

### A.1 Model 1

The visual representation of this model is given in Figure 2. The expected correlations are given in Table 8. For this model there is only one parameter constraint, as discussed in Section 3.2.1:

$$h^2 + e^2 = 1.$$

In the MH-MCMC algorithm only parameter  $e$  is set to be a free variable and estimated. Parameter  $h^2$  is defined as:  $1 - e^2$ .

**Table 8:** Expected correlations for model 1 (Figure 2) derived by using reverse path analysis.

Abbreviations	Relationship	Expected correlations
$\sigma_1$	Spouses	0
$\sigma_2$	Siblings/Dyzogitic Twins	$0.5h^2$
$\sigma_3$	Monozygotic Twins	$h^2$
$\sigma_4$	Parent-Offspring	$0.5h^2$
$\sigma_5$	Sibling in law*	0
$\sigma_6$	Spouses of twins*	0
$\sigma_7$	Nephew/niece with co-twin of twin parent*	$0.5h^2$
$\sigma_8$	Nephew/Niece with spouse of co-twin*	0
$\sigma_9$	Cousins*	$0.25h^2$

## A.2 Model 2

The visual representation of this model is given in Figure 5. The expected correlations are given in Table 9. For this model there is only one parameter constraint, as discussed in Section 3.2.3:

$$h^2 + e^2 = 1.$$

In the MH-MCMC algorithm parameters  $e$  and  $\rho$  are set to be free variables and estimated. Parameter  $h^2$  is defined as:  $1 - e^2$ .

**Table 9:** Expected correlations for model 2 (Figure 5) derived by using reverse path analysis.

Relationship	Expected correlations
$\sigma_1$	Spouses $\rho$
$\sigma_2$	Siblings/Dyzogitic Twins $0.5h^2(1 + h^2\rho)$
$\sigma_3$	Monozygotic Twins $h^2$
$\sigma_4$	Parent-Offspring $0.5h^2(1 + \rho)$
$\sigma_5$	Sibling in law* $\tau\rho$
$\sigma_6$	Spouses of twins* $\tau\rho^2$
$\sigma_7$	Nephew/niece with co-twin of twin parent* $0.5h(\lambda h + \tau\rho h)$
$\sigma_8$	Nephew/Niece with spouse of co-twin* $0.5h\rho(\lambda h + \tau\rho h)$
$\sigma_9$	Cousins* $0.25h^2(\lambda + 2\rho\lambda h^2 + \tau\rho^2 h^2)$

\*In case of MZ and DZ twin pairs,  $\tau$  equals the correlation between the MZ twins or DZ twins. Likewise,  $\lambda$  equals the genetic correlation between MZ or DZ twins; 1 or  $0.5(1 + \rho h^2)$ .

### A.3 Model 3

The visual representation of this model is given in Figure 11. The expected correlations are given in Table 10. For this model there are two parameter constraints, as discussed in Section 3.2.5:

$$h^2 + e^2 + 2hse = 1,$$

$$s = \frac{(1 + \rho)zh}{1 - (1 + \rho)ze}.$$

Van den Berg (2009) showed by substituting  $s$  and solving for  $h^2$ , no equality constraints needed to be imposed for this model:

$$h^2 + e^2 + 2he \frac{(1+\rho)zh}{1-(1+\rho)ze} = 1$$

$$h^2 \left( 1 + \frac{2ze(1+\rho)}{1-(1+\rho)ze} \right) = 1 - e^2$$

$$h^2 = (1 - e^2) \frac{1 - (1 + \rho)ze}{1 - (1 + \rho)ze + 2ze(1 + \rho)}$$

$$h^2 = (1 - e^2) \frac{1 - (1 + \rho)ze}{1 + (1 + \rho)ze}.$$

Substituting this new value for  $h^2$  in  $s$  results in,

$$s = \frac{(1 + \rho)z\sqrt{1 - e^2}\sqrt{1 - (1 + \rho)ze}}{\sqrt{1 + (1 + \rho)ze}} \cdot \frac{1}{1 - (1 + \rho)ze}$$

$$= \frac{\sqrt{1 - e^2}z(1 + \rho)}{\sqrt{1 + (1 + \rho)ze}\sqrt{1 - (1 + \rho)ze}}$$

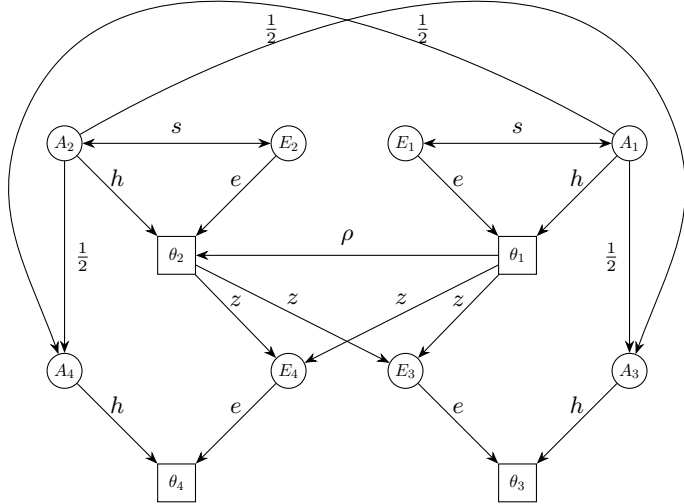
$$= \frac{\sqrt{1 - e^2}z(1 + \rho)}{\sqrt{1 - (1 + \rho)^2z^2e^2}}.$$

In the MH-MCMC algorithm parameters  $e$ ,  $\rho$  and  $z$  are set to be a free variables and estimated. With these estimations  $h^2$  and  $s$  are determined.

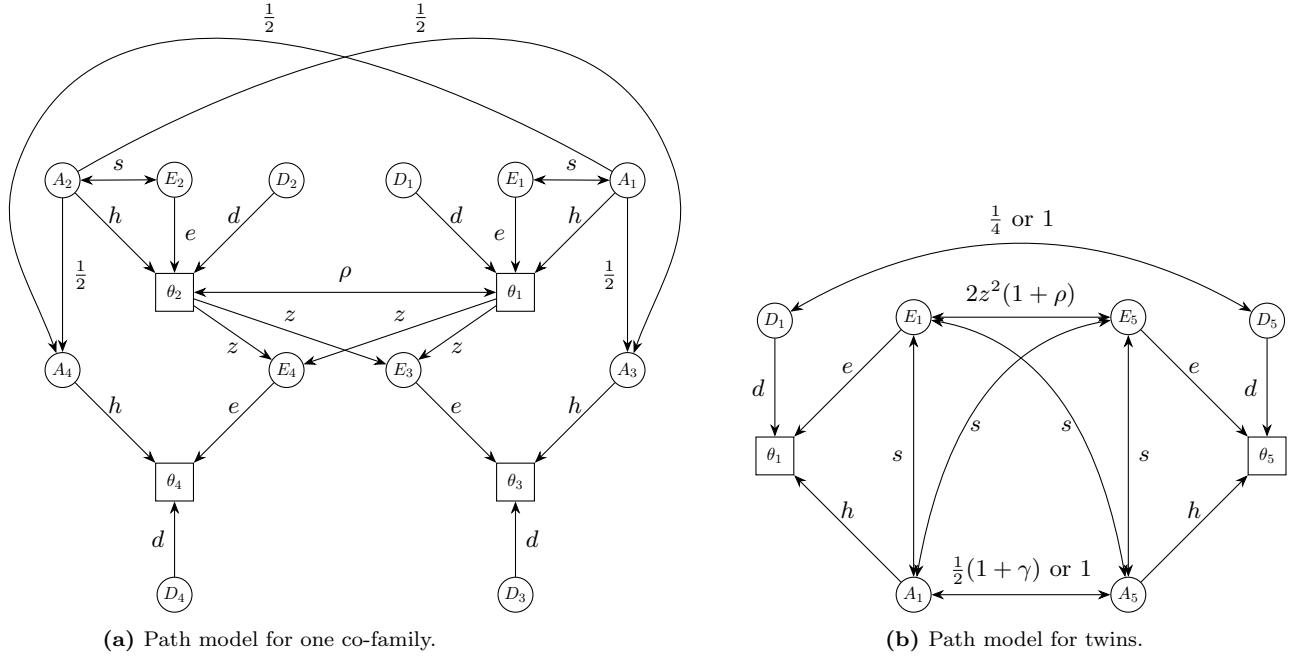
**Table 10:** Expected correlations for model 3 (Figure 11 ). Derived using reverse path analysis.

Abbreviations	Relationship	Expected correlations
$\sigma_1$	Spouses	$\rho$
$\sigma_2$	Siblings/Dyzogite Twins	$0.5h^2(1 + (h + se)^2\rho) + 2e^2z^2(1 + \rho) + 2hse$
$\sigma_3$	Monozygotic Twins	$h^2 + 2e^2z^2(1 + \rho) + 2hse$
$\sigma_4$	Parent-Offspring	$0.5h(h + se)(1 + \rho) + ez(1 + \rho)$
$\sigma_5$	Sibling in law*	$\tau\rho$
$\sigma_6$	Spouses of twins*	$\tau\rho^2$
$\sigma_7$	Nephew/niece with co-twin of twin parent*	$0.5h(\lambda h + es + \tau\rho(h + se)) + \tau ez(1 + \rho)$
$\sigma_8$	Nephew/Niece with spouse of co-twin*	$0.5h\rho(\lambda h + es + \tau\rho(h + se)) + \tau ez\rho(1 + \rho)$
$\sigma_9$	Cousins*	$0.25h^2(\lambda + 2\rho(h + se)(\lambda h + se) + \tau\rho^2(h + se)^2) + \tau z^2e^2(1 + rho)^2 + ehz(1 + \rho)(\lambda h + se + \tau\rho(h + se))$

\* In case of MZ and DZ twin pairs,  $\tau$  equals the correlation between the MZ twins or DZ twins. Likewise,  $\lambda$  equals the genetic correlation between MZ or DZ twins; 1 or  $0.5(1 + (h + se)^2\rho)$ .



The diagram shows a directed graph with five main nodes:  $E_1$ ,  $E_5$ ,  $A_1$ ,  $A_5$ , and two terminal nodes  $\theta_1$  and  $\theta_5$ . The nodes are arranged in a diamond shape with  $E_1$  and  $E_5$  at the top,  $A_1$  and  $A_5$  at the bottom, and  $\theta_1$  and  $\theta_5$  on the left and right respectively. Directed edges are labeled with parameters:  $E_1 \rightarrow \theta_1$  (e),  $E_5 \rightarrow \theta_5$  (e),  $A_1 \rightarrow \theta_1$  (h),  $A_5 \rightarrow \theta_5$  (h),  $E_1 \rightarrow A_1$  (s),  $E_5 \rightarrow A_5$  (s),  $A_1 \rightarrow E_5$  (s), and  $A_5 \rightarrow E_1$  (s). There are also horizontal edges between  $E_1$  and  $E_5$  labeled  $2z^2(1+\rho)$  and between  $A_1$  and  $A_5$  labeled  $\frac{1}{2}(1+\gamma)$  or 1.



**Figure 12:** Path model of the transmission of intelligence assuming the presence of phenotypic assortment  $\rho$ , cultural transmission  $z$  causing a correlation  $s$  between  $A$  and  $E$  and the dominance effect  $D$ .  $A_i$  represents the additive genetic value with factor loading  $h$  of person  $i$ ,  $E_i$  is the environmental value and  $e$  its factor loading and  $D_i$  is the dominance effect and  $d$  its factor loading.  $\theta_i$  represents the intelligence of person  $i$ . In case of monozygotic twins the expected correlations of  $A_1$  and  $A_5$  is 1 in case of dizygotic twins it equals  $\frac{1}{2}(1 + \gamma)$ , where  $\gamma$  equals  $(h + se)^2\rho$ . Likewise the correlation between  $D_1$  and  $D_5$  equals either 1 (MZ) or  $\frac{1}{4}$  (DZ/siblings).

**Table 11:** Expected correlations for model 4 (Figure 12). Derived using reverse path analysis.

Abbreviations	Relationship	Expected correlations
$\sigma_1$	Spouses	$\rho$
$\sigma_2$	Siblings/Dyzogitic Twins	$0.5h^2(1 + (h + se)^2\rho) + 2e^2z^2(1 + \rho) + 2hse + 0.25d^2$
$\sigma_3$	Monozygotic Twins	$h^2 + 2e^2z^2(1 + \rho) + 2hse + d^2$
$\sigma_4$	Parent-Offspring	$0.5h(h + se)(1 + \rho) + ez(1 + \rho)$
$\sigma_5$	Sibling in law*	$\tau\rho$
$\sigma_6$	Spouses of twins*	$\tau\rho^2$
$\sigma_7$	Nephew/niece with co-twin of twin parent*	$0.5h(\lambda h + es + \tau\rho(h + se)) + \tau ez(1 + \rho)$
$\sigma_8$	Nephew/Niece with spouse of co-twin*	$0.5h\rho(\lambda h + es + \tau\rho(h + se)) + \tau ez\rho(1 + \rho)$
$\sigma_9$	Cousins*	$0.25h^2(\lambda + 2\rho(h + se)(\lambda h + se) + \tau\rho^2(h + se)^2) + \tau z^2e^2(1 + rho)^2 + ehz(1 + \rho)(\lambda h + se + \tau\rho(h + se))$

\* In case of MZ and DZ twin pairs,  $\tau$  equals the correlation between the MZ twins or DZ twins. Likewise,  $\lambda$  equals the genetic correlation between MZ or DZ twins; 1 or  $0.5(1 + (h + se)^2\rho)$ .

## A.5 Model 5

The visual representation of this model is given in Figure 7. The expected correlations are given in Table 12. For this model there is only one parameter constraint, as discussed in Section 3.2.3:

$$h^2 + e^2 + c^2 = 1.$$

In the MH-MCMC algorithm parameters  $e$  and  $c$  are set to be free variables and estimated. Parameter  $h^2$  is defined as:  $1 - e^2 - c^2$ .

**Table 12:** Expected correlations for model 5 (Figure 7). Derived using reverse path analysis.

Abbreviations	Relationship	Expected correlations
$\sigma_1$	Spouses	$c^2$
$\sigma_2$	Siblings/Dyzogitic Twins	$0.5h^2 + c^2$
$\sigma_3$	Monozygotic Twins	$h^2$
$\sigma_4$	Parent-Offspring	$0.5h^2 + c^2$
$\sigma_5$	Sibling in law*	$\tau c^2$
$\sigma_6$	Spouses of twins*	$\tau c^4$
$\sigma_7$	Nephew/niece with co-twin of twin parent	$0.5h^2 + c^2h^2 + c^2$
$\sigma_8$	Nephew/Niece with spouse of co-twin	$c^2(0.5h^2 + c^2h^2 + c^2)$
$\sigma_9$	Cousins	$0.25h^2 + c^4h^2 + h^2c^2 + c^4$

\* In case of MZ and DZ twin pairs,  $\tau$  equals the correlation between the MZ twins or DZ twins.

## A.6 Model 6

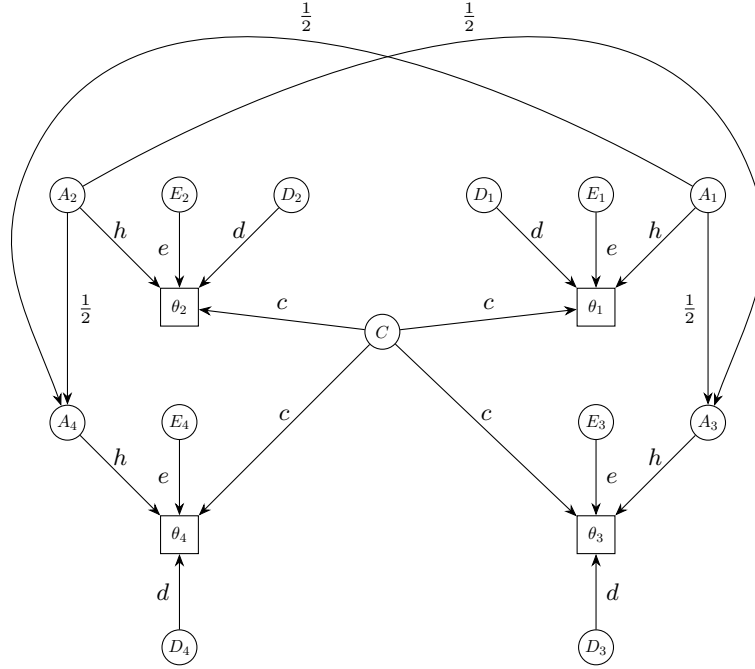
The visual representation of this model is given in Figure 13. The expected correlations are given in Table 13. For this model there is only one parameter constraint:

$$h^2 + e^2 + c^2 + d^2 = 1.$$

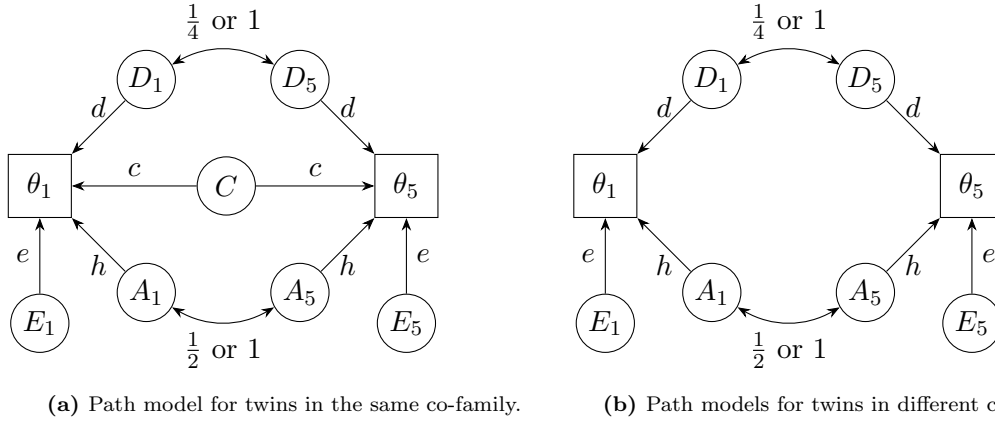
In the MH-MCMC algorithm parameters  $e$ ,  $c$  and  $d$  are set to be free variables and estimated. Parameter  $h^2$  is defined as:  $1 - e^2 - c^2 - d^2$ .

Please note that with the family structure used in this research, Figure 14b represents the structure between the twins the correct way.





**Figure 13:** Path models of the transmission of intelligence assuming social homogamy  $C$  with an effect of  $c$ .  $A_i$  represents the additive genetic value with factor loading  $h$  of person  $i$ ,  $E_i$  is the environmental value and  $e$  its factor loading.  $\theta_i$  represents the intelligence of person  $i$ . In case of monozygotic twins the expected correlations of  $A_1$  and  $A_5$  is 1 in case of dizygotic twins it equals  $\frac{1}{2}$ .



**Figure 14:** Path models of the transmission of intelligence assuming social homogamy  $C$  and the dominance effect  $D$ .  $C_i$  is the effect of social homogamy on person  $i$  with factor loading  $c$  and  $D_i$  is the dominance effect on person  $i$  with factor loading  $d$ . Two different situations are visualized. The first model assumes that twins are in one co-family, the second assumes they live in separate co-families.

**Table 13:** Expected correlations for model 6 (Figure 13 with Figure 14b). Derived by using reverse path analysis.

Abbreviations	Relationship	Expected correlations
$\sigma_1$	Spouses	$c^2$
$\sigma_2$	Siblings/Dyzogitic Twins	$0.5h^2 + c^2 + 0.25d^2$
$\sigma_3$	Monozygotic Twins	$h^2 + d^2$
$\sigma_4$	Parent-Offspring	$0.5h^2 + c^2$
$\sigma_5$	Sibling in law*	$\tau c^2$
$\sigma_6$	Spouses of twins*	$\tau c^4$
$\sigma_7$	Nephew/niece with co-twin of twin parent	$0.5h^2 + c^2h^2 + c^2d^2$
$\sigma_8$	Nephew/Niece with spouse of co-twin	$c^2(0.5h^2 + c^2h^2 + c^2d^2)$
$\sigma_9$	Cousins	$0.25h^2 + c^4h^2 + h^2c^2 + c^4d^2$

\*In case of MZ and DZ twin pairs,  $\tau$  equals the correlation between the MZ twins or DZ twins.

## B R Codes

Simulating Data {GeneticModels}

R Documentation

# Simulating data based on six different genetic models

## Description

Simulating a data set based on a one of the six genetic models. For specifications of the six different models we refer to the Masters thesis of D.M. Heeg

## Usage

```
SimulateData_function(beta, alpha, N, K, J, modelnumber, e, rho, z, d, c)
```

## Arguments

beta	Difficulty parameters of items, used to conduct the data
alpha	Discrimination parameter
N	Number of families
K	Number of items in questionnaire
J	Number of family members per family
modelnumber	Number of the model that needs to be used: 1= model 1, 2= model 2, ..., 6= model 6. For explanation of the different models, see Masters Thesis of D.M. Heeg, RUG.
e, rho, z, d, c	parameter values used to construct the covariance matrix of the genetic model. If the parameter does not apply to the model, fix it to 0.

## Value

A data file where the first list is data Y, the second beta and the third voff.

## Examples

```
SimulateData_function(beta, alpha, N, K, J, modelnumber, e, rho, z, d, c)
```

# MH-MCMC algorithm for 6 Genetic Models

## Description

Estimates the parameters from the genetic model in  $T$  iterations \ and calculates the DIC value of the model.

## Usage

```
geneticmodels_function(Y, beta, alpha, N, K, J, T, modelnumber, thinning, burnin)
```

## Arguments

<code>Y</code>	Data, the data should be a $N \times (J \times K)$ matrix, where each row represents a family \ and the columns $[(i-1)*K+1):(i*K)]$ are the item responses of family member $i$ .
<code>beta</code>	Difficulty parameters of items, used to conduct the data
<code>alpha</code>	Discrimination parameter
<code>N</code>	Number of families
<code>K</code>	Number of items in questionnaire
<code>J</code>	Number of family members per family
<code>T</code>	Number of iterations performed in MH-MCMC algorithm
<code>modelnumber</code>	Number of the model that needs to be used: \ 1= model 1, 2= model 2, ..., 6= model 6. For explanation of the different models, \ see Master's Thesis of D.M. Heeg, RUG.
<code>thinning</code>	Thinning factor
<code>burnin</code>	burn-in phase

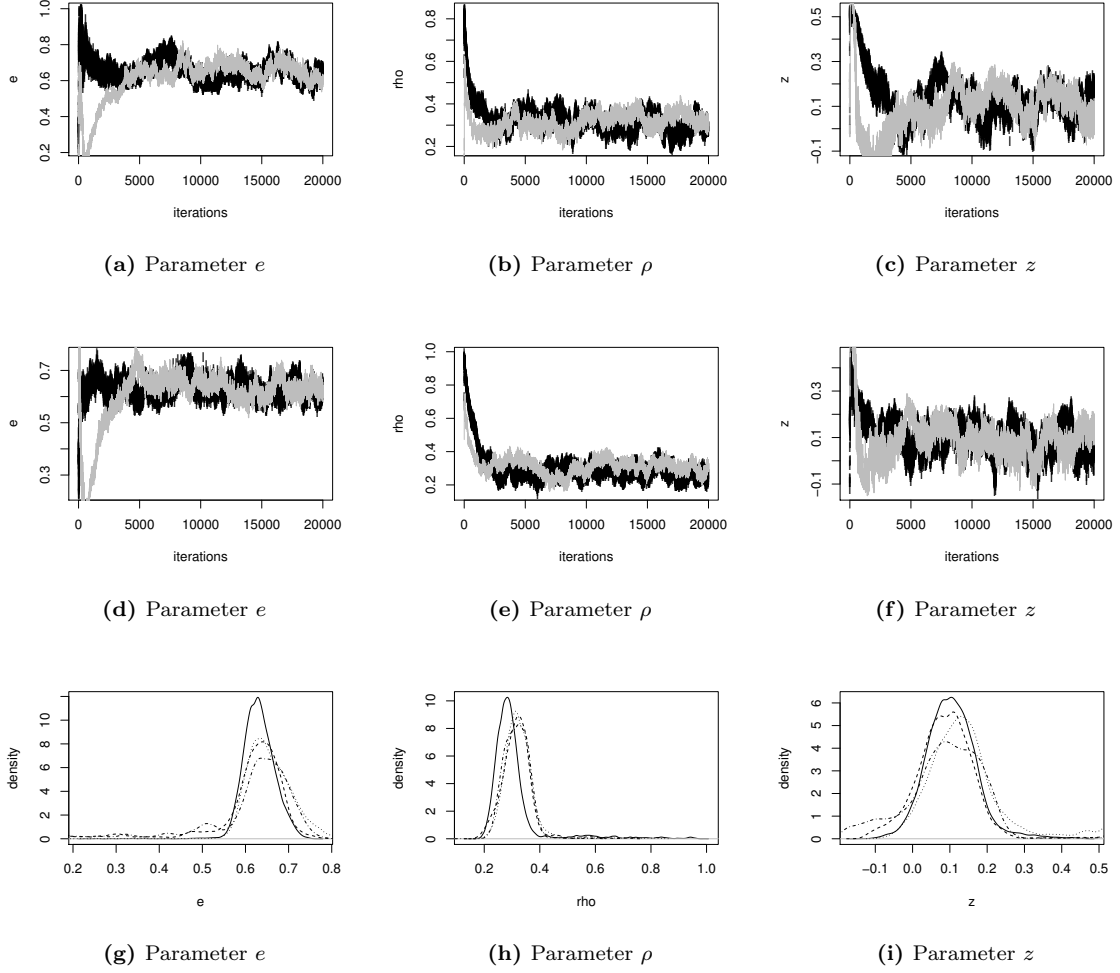
## Value

6 columns lists where the columns show the values of each saved state for the parameters. \ Empty means the parameter does not apply to the model. The last column is the DIC value \ that was generated. The list is defined in the following order: e, rho, z, d, c, DIC.

## Examples

```
geneticmodels_function()
```

## C Convergence Plots



**Figure 15:** Trace plots (15a-15f) and density plots (15g-15i) from four MH-MCMC chains with 2000.000 iterations, with a thinning factor of 100 resulting in 20.000 states given in the plots. In each trace plots, two chains are visible. All four chains are plotted simultaneously in the density plot.

## D Results: Parameter Estimates

**Table 14:** The parameter estimates (posterior means) with the standard deviation between brackets, based on the results of the simulations performed on data set 1, defined in Table 5.

Parameter	Simulated values	Model 1	Model 2	Model 3	Model 4	Model 5	Model 6
$e$	0.6	0.650 (0.010)	0.65 (0.023)	0.640 (0.021)	-0.549 (0.014)	-0.637 (0.022)	-0.637 (0.021)
$\rho$	-	-	0.013 (0.018)	0.014 (0.013)	0.040 (0.023)	-	-
$z$	-	-	-	-0.024 (0.039)	0.317 (0.013)	-	-
$d$	-	-	-	-	0.057 (0.043)	-	0.127 (0.080)
$c$	-	-	-	-	-	0.068 (0.050)	0.085 (0.058)

**Table 15:** The parameter estimates (posterior means) with the standard deviation between brackets, based on the results of the simulations performed on data set 2, defined in Table 5.

Parameter	Simulated values	Model 1	Model 2	Model 3	Model 4	Model 5	Model 6
$e$	0.76	-0.795 (0.024)	0.75 (0.014)	-0.650 (0.012)	0.758 (0.021)	-0.754 (0.022)	-0.650 (0.032)
$\rho$	0.33	-	0.36 (0.0019)	0.400 (0.029)	0.353 (0.032)	-	-
$z$	-	-	-	0.300 (0.008)	0.034 (0.041)	-	-
$d$	-	-	-	-	0.185 (0.122)	-	0.454 (0.056)
$c$	-	-	-	-	-	0.420 (0.023)	0.451 (0.031)

**Table 16:** The parameter estimates (posterior means) with the standard deviation between brackets, based on the results of the simulations performed on data set 3, defined in Table 5.

Parameter	Simulated values	Model 1	Model 2	Model 3	Model 4	Model 5	Model 6
$e$	0.63	0.621 (0.037)	-0.624 (0.023)	-0.535 (0.018)	0.621 (0.032)	0.578 (0.017)	-0.500 (0.044)
$\rho$	0.33	-	0.311 (0.018)	0.354 (0.031)	0.264 (0.032)	-	-
$z$	0.05	-	-	0.239 (0.012)	0.087 (0.051)	-	-
$d$	-	-	-	-	0.103 (0.070)	-	0.283 (0.092)
$c$	-	-	-	-	-	0.384 (0.029)	0.410 (0.034)

**Table 17:** The parameter estimates (posterior means) with the standard deviation between brackets, based on the results of the simulations performed on data set 4, defined in Table 5.

Parameter	Simulated values	Model 1	Model 2	Model 3	Model 4	Model 5	Model 6
$e$	0.31	0.757 (0.039)	-0.782 (0.012)	0.737 (0.013)	-0.415 (0.157)	0.748 (0.014)	-0.002 (0.071)
$\rho$	0.37	-	0.35 (0.010)	0.312 (0.037)	0.356 (0.033)	-	-
$z$	-0.03	-	-	-0.382 (0.016)	0.329 (0.150)	-	-
$d$	0.83	-	-	-	0.639 (0.221)	-	0.934 (0.017)
$c$	-	-	-	-	-	0.041 (0.031)	0.343 (0.022)

**Table 18:** The parameter estimates (posterior means) with the standard deviation between brackets, based on the results of the simulations performed on data set 5, defined in Table 5.

Parameter	Simulated values	Model 1	Model 2	Model 3	Model 4	Model 5	Model 6
$e$	0.1	-0.5 (0.032)	-0.46 (0.0169)	0.953 (0.033)	0.960 (0.030)	-0.016 (0.092)	-0.001 (0.085)
$\rho$	-	-	0.71 (0.011)	0.642 (0.023)	0.356 (0.033)	-	-
$z$	-	-	-	0.462 (0.008)	0.460 (0.008)	-	-
$d$	-	-	-	-	0.025 (0.020)	-	0.067 (0.041)
$c$	0.79	-	-	-	-	0.784 (0.010)	0.787 (0.011)

**Table 19:** The parameter estimates (posterior means) with the standard deviation between brackets, based on the results of the simulations performed on data set 6, defined in Table 5.

Parameter	Simulated values	Model 1	Model 2	Model 3	Model 4	Model 5	Model 6
$e$	0.05	0.68 (0.026)	0.68 (0.038)	0.641 (0.015)	0.200 (0.034)	-0.673 (0.026)	0.009 (0.121)
$\rho$	-	-	0.024 (0.022)	0.010 (0.010)	0.017 (0.014)	-	-
$z$	-	-	-	-0.400 (0.015)	0.244 (0.172)	-	-
$d$	0.85	-	-	-	0.864 (0.042)	-	0.843 (0.023)
$c$	0.21	-	-	-	-	0.027 (0.021)	0.089 (0.051)