

# A Prior Predictive Loss Function for the Evaluation of Inequality Constrained Hypotheses

Rens van de Schoot <sup>\*†</sup>, Herbert Hoijtink  
Department of Methods and Statistics, Utrecht University

Daniel Brugman  
Department of Developmental Psychology, Utrecht University,

Jan-Willem Romeijn  
Department of Philosophy, Groningen University

July 8, 2011

## Abstract

In many types of statistical modeling inequality constraints are imposed between the parameters of interest. As we will show in this paper, the DIC (i.e., posterior Deviance Information Criterion as proposed as a Bayesian model selection tool by Spiegelhalter et al., 2002) fails when comparing inequality constrained hypotheses. In this paper we will derive the prior DIC and show that it also fails when comparing inequality constrained hypotheses. However, it will be shown that a modification of the prior predictive loss function that is minimized by the prior DIC renders a criterion that does have the properties needed in order to be able to compare inequality constrained hypotheses. This new criterion will be called the Prior Information Criterion (PIC) and will be illustrated and evaluated using simulated data and examples. The PIC has a close connection with the marginal likelihood in combination with the encompassing prior approach and both methods will be compared. All in all, the main message of the current paper is: (1) do not use the classical DIC when evaluating inequality constrained hypotheses, better use the PIC; and (2) the PIC is considered a proper

---

<sup>\*</sup>Correspondence should be addressed to Rens van de Schoot: Department of Methods and Statistics, Utrecht University, P.O. Box 80.140, 3508TC, Utrecht, The Netherlands; Tel.: +31 302534468; Fax: +31 2535797; E-mail address: a.g.j.vandeschoot@uu.nl.

<sup>†</sup>This research is financed by The Netherlands Organization for Scientific Research (NWO-VICI-453-05-002).

model selection tool in the context of evaluating inequality constrained hypotheses.

*keywords:* Bayesian Model Selection, Inequality Constrained Hypothesis, Deviance Information Criterion, DIC.

1 A Prior Predictive Loss Function for the Evaluation  
2 of Inequality Constrained Hypotheses

3

4

July 8, 2011

5

**Abstract**

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

In many types of statistical modeling inequality constraints are imposed between the parameters of interest. As we will show in this paper, the DIC (i.e., posterior Deviance Information Criterion as proposed as a Bayesian model selection tool by Spiegelhalter, Best, Carlin, & Van Der Linde, 2002) fails when comparing inequality constrained hypotheses. In this paper we will derive the prior DIC and show that it also fails when comparing inequality constrained hypotheses. However, it will be shown that a modification of the prior predictive loss function that is minimized by the prior DIC renders a criterion that does have the properties needed in order to be able to compare inequality constrained hypotheses. This new criterion will be called the Prior Information Criterion (PIC) and will be illustrated and evaluated using simulated data and examples. The PIC has a close connection with the marginal likelihood in combination with the encompassing prior approach and both methods will be compared. All in all, the main message of the current paper is: (1) do not use the classical DIC when evaluating inequality constrained hypotheses, better use the PIC; and (2) the PIC is considered a proper model selection tool in the context of evaluating inequality constrained hypotheses.

*keywords:* Bayesian Model Selection, Inequality Constrained Hypothesis, Deviance Information Criterion, DIC.

27

## 1 Introduction

28

29

30

31

32

In many types of statistical modeling inequality constraints are imposed between the parameters of interest (Barlow, Bartholomew, Bremner, & Brunk, 1972; Hoijsink, Klugkist, & Boelen, 2008; Robertson, Wright, & Dykstra, 1988; Silvapulle & Sen, 2004; Van de Schoot, Hoijsink, & Deković, 2010). For an overview of literature about inequality constrained hypotheses see

33 Van de Schoot, Romeijn, and Hoijtink (2011). More specifically, the current  
34 paper considers model parameters such as means or regression coefficients  
35 that can be constrained to being greater or smaller than either a fixed value  
36 or other means or regression coefficients. Phrases like “The mean outcome  
37 in both experimental groups is expected to be larger than in the control  
38 group” and “women score higher than men in each condition” can be found  
39 in many applied papers. These specific expectations may be derived from  
40 theories, or empirical evidence, and are translated into statistical hypotheses  
41 formulated with inequality constraints. For applications, see, for example  
42 Kammers, Mulder, De Vignemont, and Dijkerman (2009); Meeus, Van de  
43 Schoot, Keijsers, Schwartz, and Branje (2010); Van de Schoot and Wong  
44 (2010); Van Well, Kolk, and Klugkist (2009). Evaluating such inequality  
45 constrained hypotheses can be done using model selection procedures. For  
46 an overview of literature about inequality constrained hypotheses see Van  
47 de Schoot et al. (2011). There is a variety of such model selection tools  
48 commonly used in practical applications, most notably Akaike’s Informa-  
49 tion Criterium (AIC; Akaike, 1973), the Bayesian Information Criterium  
50 (BIC; Schwarz, 1978), minimal description length (MDL, see, for exam-  
51 ple Grnwald, Myung, & Pitt, 2005), Bayes factors (BF; see, e.g., Kass &  
52 Raftery, 1995) and the recently developed Deviance Information Criterium  
53 (DIC; Spiegelhalter et al., 2002).

54 However, all these tools are not equipped to properly deal with inequality  
55 constrained hypotheses. Klugkist, Laudy, and Hoijtink (2005) showed that  
56 the Bayes factor can only be used in combination with an encompassing prior  
57 approach (see also, Mulder, Hoijtink, & Klugkist, 2009). Both the AIC and  
58 BIC fail when evaluating inequality constrained hypotheses because these  
59 criteria are not equipped to deal with inequality constraints between the  
60 parameters of a model. Alternatives are the order restricted information  
61 criterion (ORIC; Anraku, 1999; Kuiper & Hoijtink, 2010) which is limited  
62 to analysis of variance, and the prior-adapted-BIC (Romeijn, Van de Schoot,  
63 & Hoijtink, 2011), respectively. The MDL in relation to a reduction of the  
64 parameter space is discussed in Balasubramanian (2005). The DIC is up till  
65 now not discussed in relation to its behavior in the context of evaluating  
66 inequality constraints and this is exactly what we do in the current paper.

67 The DIC has an important role in statistical model comparison, see  
68 for example its availability in software like WinBUGS (Lunn, Thomas,  
69 Best, & Spiegelhalter, 2000), MlwiN (Rasbash, Charlton, Browne, Healy,  
70 & Cameron, 2009) or Mplus (Muthen & Muthen, 2010). However, as we  
71 will show, the DIC fails when evaluating inequality constraint hypotheses.  
72 The plan of this paper is as follow. After introducing some examples in Sec-

73 tion 2, we introduce in Section 3 the original DIC and we show that, it can  
74 not be used to choose between a set of inequality constrained hypotheses. In  
75 Section 4 we provide an alternative for the classical DIC, namely the prior  
76 Deviance Information Criterium (*prior* DIC). Unfortunately, also the *prior*  
77 DIC does not work well in the context of inequality constrained hypotheses.  
78 To accommodate for this, we propose a new loss function in Section 5, which  
79 is minimized by the Prior Information Criterion (PIC). The PIC can be used  
80 to evaluate a set of inequality constrained hypothesis. We evaluate its per-  
81 formance, see Section 6, and we show that it is connected to the marginal  
82 likelihood and thus to the Bayes factor approach of, for example, Klugkist  
83 et al. (2005).

## 84 2 Examples

85 In this section we provide three different situations where inequality con-  
86 strained hypotheses can be of interest and we describe two real-life examples  
87 where the hypotheses of interest are specified using inequality constraints.  
88 We will use Example 1 as a case study throughout the paper to investigate  
89 the performance of the *posterior* DIC, *prior* DIC and the PIC. In Section 6  
90 we briefly reconsider all other examples. Note that the scope of our proposed  
91 method is limited to the multivariate normal linear model.

### 92 2.1 Example 1

93 First, consider an example of a univariate model with where persons from  
94 two groups receive a score on one dependent variable,  $y_i$  ( $i = 1, \dots, N$ ):

$$y_i = \mu_1 d_{i1} + \mu_2 d_{i2} + \epsilon_i, \quad (2.1)$$

95 where  $\mu_1$  and  $\mu_2$  denote the mean score on  $y$  for group 1 and 2 respectively  
96 and where the residuals  $\epsilon_i$  are assumed to be normally distributed  $N(0, \sigma^2)$ .  
97 The group membership of a person is denoted by  $d_{ig} \in \{0, 1\}$ , where 1 and  
98 0 denote that a person is either a member or not a member of group  $g$ .  
99 Suppose we want to evaluate two hypotheses:  $H_0 : \mu_1, \mu_2$  and  $H_1 : \mu_1 < \mu_2$ .

100 This example has its counterparts in applied papers, see, for example,  
101 Van Well et al. (2009) about the relationship between sex, gender role iden-  
102 tification, and the gender relevance of a stressor. The authors examined  
103 mean scores for eight groups on the dependent variable stress responses,  
104 to investigate sex and gender (mis)match effects. They formulated several  
105 hypotheses by imposing inequality constraints upon group means (i.e., one

106 or more group means are expected to be larger or smaller than one or more  
107 other group means).

## 108 **2.2 Example 2**

109 Next, consider a second example of a multivariate model with two dependent  
110 variables (denoted by  $y_{1i}$  and  $y_{2i}$  for  $i = 1, \dots, N$ ),

$$\begin{aligned} y_{1i} &= \mu_1 + \epsilon_{i1} \\ y_{2i} &= \mu_2 + \epsilon_{i2} , \end{aligned} \tag{2.2}$$

111 where the residuals are assumed to be normally distributed

$$\begin{bmatrix} \epsilon_{i1} \\ \epsilon_{i2} \end{bmatrix} \sim N(0, \Sigma) , \Sigma = \begin{bmatrix} \sigma_{y_1}^2 & \rho\sigma_{y_1}\sigma_{y_2} \\ \rho\sigma_{y_1}\sigma_{y_2} & \sigma_{y_2}^2 \end{bmatrix} . \tag{2.3}$$

112 Suppose we want to evaluate two hypotheses:  $H_0 : \mu_1, \mu_2$  and  $H_1 : \mu_1 >$   
113  $0; \mu_2 > 0$ .

114 Also this multivariate model has its counterparts in applied papers, see,  
115 for example Kammers et al. (2009) about the number of body representa-  
116 tions in the brain. The authors examined the main problems that are en-  
117 countered when trying to dissociate multiple body representations in healthy  
118 individuals with the use of bodily illusions. Several models were specified  
119 within a multivariate normal model using (in)equality constraints between  
120 five repeated measurements.

## 121 **2.3 Example 3**

122 Finally, consider an example of a non-linear regression model with one de-  
123 pendent variable,  $y_i$  ( $i = 1, \dots, N$ ) with a linear, i.e.  $X_i$ , and a non-linear  
124 predictor, i.e.  $X_i^2$ :

$$y_i = \beta_0 + \beta_1 X_i + \beta_2 X_i^2 + \epsilon_i , \tag{2.4}$$

125 where  $\beta_0$  denote the intercept and where the residuals  $\epsilon_i$  are assumed to be  
126 normally distributed with  $N(0, \sigma^2)$ .

127 Suppose we want to evaluate two hypotheses:  $H_0 : \beta_2$  and  $H_1 : \beta_2 > 0$ .  
128 Such expectations are of interest in, for example, the research of retention  
129 memory (see, e.g., Myung, 2003). See Section 2.5 for a real life example.

## 130 **2.4 Real-life Example 1: Moral Judgment Competence**

131 Leenders and Brugman (2005) investigated whether moral judgment com-  
132 petence and attitude towards delinquent behavior create a domain shift in

133 young adolescents. That is, a certain behavior which in society as a whole  
134 is considered to be not moral (e.g. aggression, violence), might be a group  
135 convention in certain adolescent groups. In total 135 pupils of intermedi-  
136 ate secondary schools in the Netherlands were asked to report whether the  
137 respondent had committed such behaviour (never, once, more than once).  
138 They were also asked to judge aggressive acts and vandalistic acts in hypo-  
139 theoretical situations on how moral they thought the behavior was. For each  
140 hypothetical situation, questions were asked (on a 4-point scale) about the  
141 acceptability (Is it wrong or right to do such a thing?), the seriousness (How  
142 bad is it to do such a thing?), the generalizability (If everybody were doing  
143 such things, would they then be wrong or right?) and the rule/authority  
144 contingency (If nobody saw it, would it then be wrong or right?) of the  
145 transgression. Just like in the original article, for each category the sum  
146 scores were computed in a way that a high criterion score indicated a more  
147 non-moral (conventional/personal) judgment. The researchers had specific  
148 ideas about differences in the level of morality in these hypothetical situa-  
149 tions between pupils that did or did not report to conduct aggressive acts  
150 themselves.

151 The model under investigation is given by

$$\begin{aligned} y_{1i} &= \mu_{11}d_{ig1} + \mu_{12}d_{ig2} + \epsilon_{1i} \\ y_{2i} &= \mu_{21}d_{ig1} + \mu_{22}d_{ig2} + \epsilon_{2i} \end{aligned} \quad (2.5)$$

152 where  $\mu_1$ . and  $\mu_2$ . denote the mean score on the hypothetical construct van-  
153 dalism (denoted by  $y_1$ ) and the hypothetical construct aggression (denoted  
154 by  $y_2$ ) and where  $\mu_{.1}$  and  $\mu_{.2}$  denote the mean for the group reported not  
155 to conduct aggressive acts and the group that did report to conduct aggres-  
156 sive acts, respectively. Again, group membership of a person is denoted by  
157  $d_{ig} \in 0, 1$  and the residuals are assumed to be normally distributed with

$$\begin{bmatrix} \epsilon_{i1} \\ \epsilon_{i2} \end{bmatrix} \sim N(0, \Sigma), \Sigma = \begin{bmatrix} \sigma_{y_1}^2 & \rho\sigma_{y_1}\sigma_{y_2} \\ \rho\sigma_{y_1}\sigma_{y_2} & \sigma_{y_2}^2 \end{bmatrix}. \quad (2.6)$$

158 Note that this example is a combination of (2.2) and (2.4).

159 There are three hypotheses of interest:

$$\begin{aligned} H_0 &: \mu_{12}, \mu_{11} \text{ and } \mu_{22}, \mu_{21} \\ H_1 &: \mu_{12} > \mu_{11} \text{ and } \mu_{22} > \mu_{21} \\ H_2 &: \mu_{12} = \mu_{11} \text{ and } \mu_{22} > \mu_{21} \end{aligned} \quad (2.7)$$

160 The first hypothesis, is an unconstrained hypothesis ( $H_0$ ). A second hy-  
161 pothesis,  $H_1$ , postulates that the aggressive group ( $\mu_{.2}$ ) also judge the same

162 behavior in all hypothetical situations to be more conventional and as such  
163 morally more appropriate than their peers who do not report such behavior  
164 ( $\mu_1$ ). The third hypothesis,  $H_2$ , is that there is a domain shift in the judge-  
165 ment about hypothetical situations. That is, for pupils that reported to have  
166 conducted some delinquent behavior (i.e. aggression), in the same hypothet-  
167 ical situation, they will judge it to be more morally accepted compared to  
168 adolescents that did not report to conduct the same behavior. However, in  
169 hypothetical situations concerning other delinquent behavior that was not  
170 reported by these same adolescents (i.e. vandalism), they will judge the hy-  
171 pothetical situation to be equally morally condemnable as adolescents that  
172 did not report any antisocial behavior. In Section 6.3 the data of Leenders  
173 and Brugman (2005) will be used to re-evaluate these hypotheses.

## 174 2.5 Real-life Example 2: Ph.D. delays

175 Sonneveld, Yerkes, and Van de Schoot (2009) report on Ph.D. trajectories  
176 and employment outcomes of recent Dutch Ph.D. recipients at four uni-  
177 versities in the Netherlands. The report provides detailed information on  
178 the background of Ph.D. candidates, their Ph.D. trajectory, including su-  
179 pervision and the performance of Ph.D. candidates, as well as their initial  
180 employment after obtaining their Ph.D.

181 In the Netherlands it is possible to differentiate between three different  
182 types of Ph.D. status, including: (a) a Ph.D. candidate that is employed by  
183 the university, (b) scholarship recipients and (c) external and/or dual Ph.D.  
184 candidates. Full employment contracts for Ph.D. candidates are the excep-  
185 tion and not the rule throughout Europe. Only the Netherlands, Finland  
186 and Turkey have doctoral educational structures in which different types of  
187 Ph.D. status exist simultaneously. The majority of respondents surveyed  
188 (71.1%) reported that their main formal status was 'employee'. In the cur-  
189 rent paper we will only focus on employees ( $n = 304$ ).

190 Among many other questions, the researchers asked the Ph.D. recipients  
191 how long it took them to finish their Ph.D thesis. It appeared that Ph.D.  
192 recipients took an average of 59.8 months (five years and four months) to  
193 complete their Ph.D. trajectory. In the current paper we will answer the  
194 question why some Ph.D. recipients took longer than other by using age as a  
195 predictor ( $M = 30.7$ ,  $SD = 4.48$ , min-max = 26-69 ). The relation between  
196 completion time and age is expected to be non-linear. This might be due  
197 to the fact that at a certain point in your life (i.e., mid-thirties), family life  
198 takes up more of your time than when you are in your twenties or when you  
199 are older.



200 However, we expect that if you are in your mid-thirties and you are  
 201 doing a Ph.D. you also take this extra time into account. The researchers  
 202 asked to Ph.D. candidates about their planned graduation day according  
 203 to the original contract and their actual graduation day. The average gap  
 204 between the two data appeared to be 9.6 months (SD = 14.4, min-max =  
 205 -3 - 69). We expect that the lag between planned and actual time spent  
 206 on the trajectory is less prone to non-linearity compared to actual project  
 207 time.

208 If  $y_{1i}$  denotes actual project time and  $y_{2i}$  denotes the lag between actual  
 209 and planned project time, the model under investigation is given by

$$\begin{aligned} y_{1i} &= \beta_{01} + \beta_1 Age_i + \beta_2 Age_i^2 + \epsilon_{1i} , \\ y_{2i} &= \beta_{02} + \beta_3 Age_i + \beta_4 Age_i^2 + \epsilon_{2i} . \end{aligned} \quad (2.8)$$

210 To avoid multicollinearity  $Age$  will be centered. The residuals are assumed  
 211 to be normally distributed

$$\begin{bmatrix} \epsilon_{i1} \\ \epsilon_{i2} \end{bmatrix} \sim N\left(0, \Sigma\right), \Sigma = \begin{bmatrix} \sigma_{y_1}^2 & \rho\sigma_{y_1}\sigma_{y_2} \\ \rho\sigma_{y_1}\sigma_{y_2} & \sigma_{y_2}^2 \end{bmatrix} . \quad (2.9)$$

212 The following hypotheses are of interest:

$$\begin{aligned} H_0 &: \beta_2, \beta_4 \\ H_1 &: \beta_2 > 0 \text{ and } \beta_4 > 0 \\ H_2 &: \beta_2 > \beta_4 > 0 . \end{aligned} \quad (2.10)$$

213 In Section 6.4 the data of Sommeveld et al. (2009) will be used to evaluate  
 214 these hypotheses.

### 215 3 Posterior DIC

216 One way of evaluating hypotheses, is to use a model selection approach.  
 217 This is not a test of the model in the sense of hypothesis testing, rather it  
 218 is an evaluation between models using a trade-off of model fit and model  
 219 complexity. The likelihood of an hypothesis is a measure of model fit, and the  
 220 number of parameters involved in the hypothesis is a measure of complexity.  
 221 The greater the number of parameters, the larger the compensation for  
 222 model complexity becomes. So, adding a parameter should be accompanied  
 223 by an increase in model fit to accommodate for the increase in complexity.  
 224 Several competing statistical models may be ranked according to their value  
 225 on the model selection tool used and the one with the best trade-off is the  
 226 winner of the model selection competition.

227 The Deviance Information Criterion (DIC), is proposed in Spiegelhalter  
 228 et al. (2002) as a Bayesian criterion for minimizing the posterior predic-  
 229 tive loss. In this section we briefly introduce the DIC, thereafter we show  
 230 with our running example that the DIC fails when comparing inequality  
 231 constrained hypotheses.

232 Note that from now on, we will use the *posterior* DIC whenever we refer  
 233 to the DIC of Spiegelhalter et al. (2002) and we will use *prior* DIC whenever  
 234 we refer to our adjustment of the DIC, that will be introduced in Section 4.

### 235 3.1 Definition

236 The *posterior* DIC minimizes the posterior expectation of the expected loss  
 237 (Gelman, Carlin, Stern, & Rubin, 2004). It is defined as the error that is  
 238 expected when a statistical model estimated by the observed data set  $\mathbf{y}$  is  
 239 applied to a future data set  $\mathbf{x}$ . Let  $f(\cdot)$  denote the likelihood, then the  
 240 expected loss is given by

$$E_{f(\mathbf{x}|\boldsymbol{\theta}^*)}[-2 \log f(\mathbf{x} | \bar{\boldsymbol{\theta}}_y)] , \quad (3.1)$$

241 where  $-2 \log f(\cdot)$  is the loss function of a future data set  $\mathbf{x}$  in which  $\bar{\boldsymbol{\theta}}_y$  is  
 242 the expected a-posteriori estimate of the model parameters  $\boldsymbol{\theta}$  based on the  
 243 observed data set  $\mathbf{y}$ . If we would know the true parameter value  $\boldsymbol{\theta}^*$ , the  
 244 expectation in (3.1) could be computed. However, since these are unknown,  
 245 the *posterior* DIC takes the posterior expectation of (3.1). Let  $E_{g(\boldsymbol{\theta}|\mathbf{y})}$  de-  
 246 notes the expectation with respect to the posterior distribution  $g(\boldsymbol{\theta} | \mathbf{y})$ ,  
 247 then

$$E_{g(\boldsymbol{\theta}|\mathbf{y})} \left\{ E_{f(\mathbf{x}|\boldsymbol{\theta})} [-2 \log f(\mathbf{x} | \bar{\boldsymbol{\theta}}_y)] \right\} \approx \\
 -2 \log f(\mathbf{y} | \bar{\boldsymbol{\theta}}_y) + 2 \left[ -2 \overline{\log f(\mathbf{y} | \boldsymbol{\theta})} + 2 \log f(\mathbf{y} | \bar{\boldsymbol{\theta}}_y) \right] , \quad (3.2)$$

248 where (3.2) is the definition of the *posterior* DIC. The term  $-2 \log f(\mathbf{y} | \bar{\boldsymbol{\theta}}_y)$   
 249 in (3.2) is often interpreted as model (mis)fit and the term  $[-2 \overline{\log f(\mathbf{y} | \boldsymbol{\theta})} +$   
 250  $2 \log f(\mathbf{y} | \bar{\boldsymbol{\theta}}_y)]$  in (3.2) is often interpreted as the effective number of pa-  
 251 rameters and is considered a penalty term.

### 252 3.2 Estimation

253 The *posterior* DIC can be computed using Monte Carlo simulation and is  
 254 available in several software packages, for example, WinBUGS (Lunn et al.,  
 255 2000), MLwiN (Rasbash et al., 2009) and Mplus (Muthen & Muthen, 2010).

256 Let  $\boldsymbol{\theta}^1 \dots \boldsymbol{\theta}^L$  be  $L$  draws from the posterior distribution  $g(\boldsymbol{\theta} \mid \mathbf{y})$ , then  
 257  $-2\log f(\mathbf{y} \mid \bar{\boldsymbol{\theta}})$  can be estimated by

$$\sum_{l=1}^L \frac{-2\log f(\mathbf{y} \mid \boldsymbol{\theta}^l)}{L}, \quad (3.3)$$

258 and  $-2\log f(\mathbf{y} \mid \bar{\boldsymbol{\theta}}_y)$  can be estimated by

$$-2\log f(\mathbf{y} \mid \sum_{l=1}^L \frac{\boldsymbol{\theta}_1^l}{L}, \dots, \sum_{l=1}^L \frac{\boldsymbol{\theta}_k^l}{L}), \quad (3.4)$$

259 where  $k$  is an index for the parameters in  $\boldsymbol{\theta}$  ( $k = 1, \dots, K$ ).

260 An important issue when computing the *posterior* DIC is the specifica-  
 261 tion of the prior distribution. A default approach is to specify a vague or  
 262 low-informational prior distribution. In that case, the computation of the  
 263 *posterior* DIC is independent of the specified prior because the posterior  
 264 distribution,  $g(\boldsymbol{\theta} \mid \mathbf{y})$ , is dominated by the data.

### 265 3.3 Behavior of the Posterior DIC in Constrained Model Se- 266 lection

267 To inspect the behavior of the *posterior* DIC in the context of evaluating  
 268 inequality constraint hypotheses, we consider Example 1 with  $H_0 : \mu_1, \mu_2$   
 269 and  $H_1 : \mu_1 < \mu_2$ . According to Mulder, Hoijsink, and Klugkist (2009), the  
 270 prior distribution for Example 1 is given by

$$h_0(\mu_1, \mu_2, \sigma^2) = N(\mu_1 \mid \mu_0, \tau_0^2) \times N(\mu_2 \mid \mu_0, \tau_0^2) \times Inv\chi^2(\sigma^2 \mid \nu_0, \sigma_0^2), \quad (3.5)$$

271 where  $\mu_0$  is the prior mean and  $\tau_0^2$  is the prior variance. Hypothesis  $H_1$  is  
 272 nested in  $H_0$ , therefore  $h_1(\cdot)$  is proportional to  $h_0(\cdot)$ , with

$$h_1(\cdot) : \begin{cases} c^{-1}h_0(\cdot) & \text{if } (\mu_1, \mu_2) \in H_1 \\ 0 & \text{otherwise} \end{cases}, \quad (3.6)$$

273 where  $c$  is a normalization constant given by

$$c = \int_{(\mu_1, \mu_2) \in H_1} h_0(\mu_1, \mu_2) d(\mu_1, \mu_2). \quad (3.7)$$

274 Using this encompassing prior approach only the prior distribution for  $H_0$   
 275 needs to be specified. Note that the encompassing prior approach has been

276 used in computing Bayes factors, which will not be considered in the current  
 277 paper, but see for more information Mulder, Hoijtink, and Klugkist (2009).

278 Now let  $g_0(\cdot)$  denote the posterior distribution of the unconstrained hy-  
 279 potheses and  $g_1(\cdot)$  the posterior distribution of  $H_1$ , then  $g_0(\cdot) \propto f(\cdot) \times h_0(\cdot)$   
 280 and  $g_1(\cdot) \propto f(\cdot) \times h_1(\cdot)$ . Then,  $g_1(\cdot) = d^{-1}g_0(\cdot)$  where

$$d = \int_{(\mu_1, \mu_2) \in H_1} g_0(\mu_1, \mu_2) \partial(\mu_1, \mu_2) . \quad (3.8)$$

281 For  $\mu_2 - \mu_1 \rightarrow \infty$ ,  $g_0(\mu_1, \mu_2, \sigma^2 | \mathbf{y}) - g_1(\mu_1, \mu_2, \sigma^2 | \mathbf{y}) \rightarrow 0$ . That is, if  
 282 the population from which the data are generated is strongly in agreement  
 283 with  $H_1$ , the difference between the posterior distributions for  $H_0$  and  $H_1$   
 284 goes to zero. Since the *posterior* DIC is computed using samples of  $\mu_1, \mu_2$   
 285 and  $\sigma^2$  obtained from the posterior distribution, see Equations (3.3) and  
 286 (3.4), for  $\mu_2 - \mu_1 \rightarrow \infty$ , samples obtained under  $H_0$  and  $H_1$  are exchange-  
 287 able. Consequently,  $\text{DIC}_{H_0}$  and  $\text{DIC}_{H_1}$  have the same values. This result  
 288 is counterintuitive and unwanted because  $H_1$  is more parsimonious than  $H_0$   
 289 and hence it contains more information (cf. Sober, 2006), so it should be  
 290 preferred by the DIC.

291 A simulation study was performed to illustrate the failure of the *poste-*  
 292 *rior* DIC. You can also derive analytic expressions for the behavior of the  
 293 posterior probability distribution, on which the behavior of the *posterior*  
 294 DIC hinges. For more details see Romeijn et al. (2011). Here we discuss  
 295 this behavior merely for the purpose of illustration. Seven data sets from  
 296 seven populations were considered. Data were constructed in such a way  
 297 that the sample means and variance are exactly equal to the population pa-  
 298 rameters (with  $\sigma^2 = 1$  and  $n = 20$  for each group). The population means  
 299 for the seven data sets are displayed on the x-axis in Figure 1. Note that  
 300 the first four data sets are in agreement with the constraints of  $H_1$ , whereas  
 301 the last three data sets are constructed in such a way that they violate the  
 302 constraints of  $H_1$ . The difference between the seven data sets is that the size  
 303 of the difference between the two group means varies from small to large.  
 304 We also considered an equality constrained hypothesis,  $H_2 : \mu_1 = \mu_2$ , to in-  
 305 vestigate the performance of the *posterior* DIC. For this hypothesis,  $\mu_1$  and  
 306  $\mu_2$  can be replaced by  $\mu$ . For each data set we used WinBUGS to compute  
 307 the *posterior* DIC.

308 Next, the hypotheses of interest were evaluated for all seven data sets  
 309 with the *posterior* DIC. The results are presented in Figure 1. When looking  
 310 at populations 1-5 in Figure 1, it can be seen that the values of the *posterior*  
 311 DIC for  $H_0$  and  $H_1$  are equal. Hence, the *posterior* DIC can not distinguish

312  $H_0$  and  $H_1$ . This is counterintuitive because the population values satisfy  
313 the constraints of  $H_1$  and  $H_1$  is more parsimonious than  $H_0$ . For population  
314 4 and 5, the two data sets with the smallest difference in sample means, the  
315 value of the *posterior* DIC for  $H_2$  is lowest. This result is in line with what  
316 would be expected because the means are approximately equal. When the  
317 population means do not fit the constraints imposed by  $H_1$  (i.e. populations  
318 6 and 7) the values for the *posterior* DIC for  $H_0$ ,  $H_1$  and  $H_2$  are in line  
319 with what would be expected: the lowest value for  $H_0$  followed by  $H_2$  and  
320  $H_1$ , respectively. In sum, the *posterior* DIC fails to distinguish between  
321 hypotheses  $H_0$  and  $H_1$  when the data are strongly in agreement with the  
322 most constrained hypothesis,  $H_1$ .

## 323 4 Prior DIC

324 Within the Bayesian framework, there are two perspectives on model selec-  
325 tion: a prior predictive approach (e.g. Box, 1980; Kass & Raftery, 1995) and  
326 a posterior predictive approach (e.g. Gelman et al., 2004; Gelman, Meng,  
327 & Stern, 1996). Spiegelhalter, Best, Carlin, and Van Der Linde (2002)  
328 derived the posterior Deviance Information Criterion (*posterior* DIC) to  
329 choose between a set of competing hypotheses. As we have seen in the pre-  
330 vious section, the *posterior* DIC failed to choose between a set of inequality  
331 constrained hypotheses. In this section we will derive the prior Deviance  
332 Information Criterion (*prior* DIC).

### 333 4.1 Definition

334 The point of departure for the *prior* DIC is the same as for the *posterior*  
335 DIC, namely the expected loss given in (3.1). However, to deal with the  
336 unknown parameters  $\boldsymbol{\theta}^*$ , for the *prior* DIC, we take the expectation of the  
337 expected loss with respect to the prior distribution,  $h(\boldsymbol{\theta})$ , instead of the  
338 posterior distribution,  $g(\boldsymbol{\theta} \mid \mathbf{y})$ , as was the case for the *posterior* DIC:

$$E_{h(\boldsymbol{\theta})} \left\{ E_{f(\mathbf{x} \mid \boldsymbol{\theta})} [-2 \log f(\mathbf{x} \mid \bar{\boldsymbol{\theta}}_y)] \right\} \quad (4.1)$$

339 The major difference between (3.2) and (4.1) is that  $g(\boldsymbol{\theta} \mid \mathbf{y})$  is replaced by  
340  $h(\boldsymbol{\theta})$ . As will be shown, using  $h(\boldsymbol{\theta})$  instead of  $g(\boldsymbol{\theta} \mid \mathbf{y})$  is a final step towards  
341 an IC that does not suffer from the drawbacks discussed in the previous  
342 section.

343 The main problem now, is to find an expression for  $E_{h(\boldsymbol{\theta})} [c(\mathbf{y}, \boldsymbol{\theta}, \bar{\boldsymbol{\theta}}_y)]$   
344 and this is what we do in Appendix A resulting in the definition of the *prior*

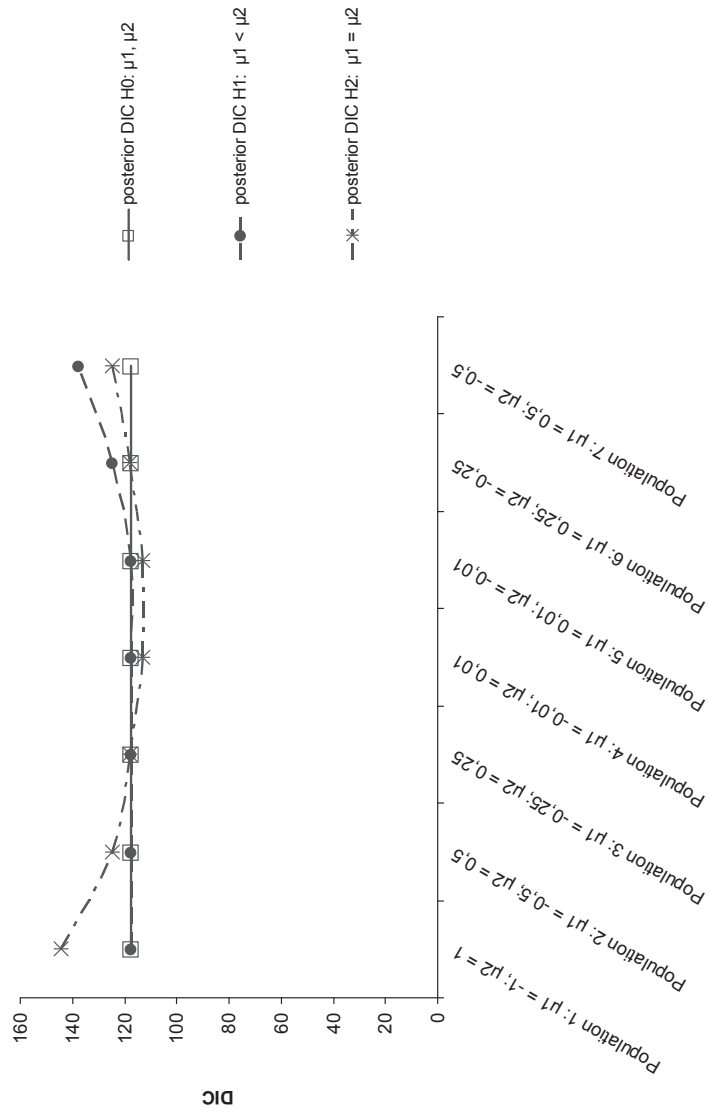


Figure 1: Values of the Posterior DIC for  $H_0$ ,  $H_1$  and  $H_2$ ; and for populations 1-7 of Example 1.

345 DIC:

$$\begin{aligned} E_{h(\boldsymbol{\theta})} \left\{ E_{f(\mathbf{x}|\boldsymbol{\theta})} [-2 \log f(\mathbf{x} | \bar{\boldsymbol{\theta}}_y)] \right\} \approx \\ C + 2 \log f(\mathbf{y} | \bar{\boldsymbol{\theta}}_y) + E_{h(\boldsymbol{\theta})} [-2 \log f(\mathbf{y} | \boldsymbol{\theta})] , \end{aligned} \quad (4.2)$$

346 where  $C = E_{h(\boldsymbol{\theta})} \left\{ E_{f(\mathbf{x}|\boldsymbol{\theta})} [-2 \log f(\mathbf{x} | \boldsymbol{\theta})] \right\}$  is constant when comparing  
347 inequality constrained hypotheses, see Appendix B, and consequently can  
348 be ignored.

349 Note the two major differences between the *prior* and *posterior* DIC:  
350 the first term of (4.2) (i.e.  $C$ ) does not have a corresponding part in the  
351 definition of the *posterior* DIC, see (3.2) and the third term on the right  
352 hand side of (4.2) (i.e.  $E_{h(\boldsymbol{\theta})} [-2 \log f(\mathbf{y} | \boldsymbol{\theta})]$ ) is the expectation with respect  
353 to the prior distribution whereas the corresponding term in (3.2) is the  
354 expectation with respect to the posterior distribution.

## 355 4.2 Estimation

356 The *prior* DIC can be computed using Monte Carlo simulation, for example  
357 using R (R Development Core Team, 2006). Let  $\boldsymbol{\theta}^1 \dots \boldsymbol{\theta}^L$  be  $L$  draws from  
358 the posterior distribution, then  $2 \log f(\mathbf{y} | \bar{\boldsymbol{\theta}}_y)$ , in Equation (4.2) can be  
359 estimated by

$$2 \log f(\mathbf{y} | \frac{1}{L} \sum_{l=1}^L \boldsymbol{\theta}_1^l, \dots, \frac{1}{L} \sum_{l=1}^L \boldsymbol{\theta}_k^l). \quad (4.3)$$

360 Furthermore, let  $\boldsymbol{\theta}^1 \dots \boldsymbol{\theta}^K$  be  $K$  draws from the prior distribution, then  
361  $E_{h(\boldsymbol{\theta})} [-2 \log f(\mathbf{y} | \boldsymbol{\theta})]$  in Equation (4.2) can be estimated by

$$\frac{1}{K} \sum_{k=1}^K -2 \log f(\mathbf{y} | \boldsymbol{\theta}^k). \quad (4.4)$$

362 Just like for the *posterior* DIC, the specification of the prior distribution  
363 is of importance. For the *prior* DIC it is even essential that the prior distri-  
364 bution is specified correctly because only then background knowledge in the  
365 form of inequality constraints between the parameters of interested can be  
366 incorporated. In order to incorporate the constraints in the prior distribu-  
367 tion, we use the encompassing prior approach as was discussed before. The  
368 prior is given in Equation (3.5) and we assume the same prior distribution for  
369 each parameter that is subjected to constraints,  $h_0(\mu_1) = h_0(\mu_2)$ . Specifying  
370 the parameters of the prior distribution in constrained hypotheses selection

371 is further explained in Mulder, Hoijtink, and Klugkist (2009) and Mulder,  
 372 Klugkist, et al. (2009). The actual computation of the second term of the  
 373 *prior* DIC for  $H_0$  and  $H_1$  can be done using samples from  $g_0(\mu_1, \mu_2, \sigma^2 | \mathbf{y})$   
 374 and  $g_1(\mu_1, \mu_2, \sigma^2 | \mathbf{y})$ , respectively. These samples can be obtained using the  
 375 Gibbs sampler for  $g_0(\cdot)$  (see, Gelman et al., 2004) and the constrained Gibbs  
 376 sampler for  $g_1(\cdot)$  (see, Klugkist et al., 2005). The third term of the *prior*  
 377 DIC can be computed using a sample from the prior distribution of the  
 378 hypotheses under investigation.

### 379 4.3 Behavior of The prior DIC in Constrained Model Selec- 380 tion

381 To show that the *prior* DIC can be used to choose between a set of con-  
 382 strained hypotheses if the population from which the data are generated  
 383 is fully in agreement with the most constrained hypothesis, whereas the  
 384 *posterior* DIC fails to do so, we reconsider Example 1.

385 If we would compare  $H_0$  and  $H_1$  with the *prior* DIC, the first term  
 386 of the *prior* DIC given in Equation (4.2) is constant (see Appendix B).  
 387 Now, consider the same situation as in the beginning of Section 3.3 where  
 388 the population from which the data was generated is strongly in agreement  
 389 with  $H_1$ . In this case, the second term in Equation (4.2) does also not  
 390 differ between  $H_0$  and  $H_1$ , because for  $\mu_1 - \mu_2 \rightarrow \infty$ ,  $\bar{\mu}_1 | H_0 \rightarrow \bar{\mu}_1 | H_1$   
 391 and  $\bar{\mu}_2 | H_0 \rightarrow \bar{\mu}_2 | H_1$ . So, the third term,  $E_{h(\mu_1, \mu_2, \sigma^2)}[\cdot]$ , should make the  
 392 difference between  $H_0$  and  $H_1$ .

393 Since samples of  $\mu_1$  and  $\mu_2$  are taken from the prior distribution  $h_0(\mu_1, \mu_2, \sigma^2)$   
 394 and since  $h_0(\mu_1, \mu_2, \sigma^2) \neq h_1(\mu_1, \mu_2, \sigma^2)$  because of the normalization of  
 395 the prior distribution according to Equation (3.6), samples from the prior  
 396 distribution are different for  $H_0$  and  $H_1$ . For  $\mu_1 - \mu_2 \rightarrow \infty$ , the third  
 397 term of (4.2) when computed for  $H_0$  is based on more large values of  
 398  $-2 \log f(\mathbf{y} | \mu_1, \mu_2, \sigma^2)$  than when it is computed for  $H_1$ . Consequently,  
 399 the third term of (4.2) for  $H_1$  is smaller than the third term of (4.2) for  $H_0$ .

400 Again, a simulation study was performed where data sets from the seven  
 401 populations of Section 3.3 were considered. The exact specification of the  
 402 parameters of the prior distribution for population 1 with  $\mu_1 = -1$  and  
 403  $\mu_2 = 1$ , are  $\mu_0 = 0$ ,  $\tau_0^2 = 0.97$ ,  $v_0 = 2$  and  $\sigma_0^2 = 1.95$ .

404 In contrast to the *posterior* DIC, the *prior* DIC is able to correctly distin-  
 405 guish between  $H_0$  and  $H_1$  when the data are in agreement of the constraints  
 406 of  $H_1$ , see populations 1-3 in Figure 2, where the *prior* DIC is lowest for  
 407  $H_1$ . For the data with the smallest differences in sample means (population  
 408 4 and 5), the *prior* DIC is lowest for  $H_2$ . When the constraints are not sup-



409 ported by the data, populations 6-7, the value for  $H_0$  should be the lowest  
 410 value, but as can be seen in Figure 2, this is not the case! So, when the  
 411 data are fully in agreement with  $H_1$  the *prior* DIC outperforms the *posterior*  
 412 DIC, but when the data do not support  $H_1$ , the *prior* DIC fails to correctly  
 413 distinguish the three hypotheses.

414 What goes wrong? Consider the prior expectation of the expected loss  
 415 given in (3.1), which is approximated by the *prior* DIC as was shown in  
 416 Appendix A:

$$E_{h(\theta)} \left\{ E_{f(\mathbf{x}|\theta)} [-2 \log f(\mathbf{x} | \bar{\theta}_y)] \right\} \quad (4.5)$$

$$\approx 2 \log f(\mathbf{y} | \bar{\theta}_y) + E_{h(\theta)} [-2 \log f(\mathbf{y} | \theta)] . \quad (4.6)$$

417 The loss function in Equation (4.5) captures how well replicated data fit  
 418 a certain hypothesis, that is, how good  $\bar{\theta}_y$  is a summary of  $\mathbf{x}$ . However,  
 419 this loss function does not accommodate ‘bad’ fitting hypotheses, that is, if  
 420 for a hypothesis  $\bar{\theta}_y$  is not a good summary of  $\mathbf{y}$ , this will not be detected  
 421 by the loss function in (4.5). Note that it might appear the correction for  
 422 ‘bad’ fitting hypotheses is done by the first term of the approximation of  
 423 the loss function, see Equation (4.6). However, the second term cancels the  
 424 influence of the first term because the second term can be written as a Taylor  
 425 expansion around the first term, see Appendix A and Equation (A.7).

426 Let us return to the loss function in Equation (4.6) and consider the  
 427 situation of Example 1. Suppose that a population is not in agreement with  
 428 the inequality constrained hypothesis,  $H_1 : \mu_1 < \mu_2$ , for example Population  
 429 7 with population means  $\mu_1 = 0.5; \mu_2 = -0.5$ . In this situation the *prior*  
 430 DIC chooses  $H_1$  as the best hypothesis, see Figure 2. This result is unwanted  
 431 because the means in the data satisfy  $\mu_1 > \mu_2$ .

432 Under the assumption  $\mu_1 < \mu_2$  in the data were  $\mu_1 = 0.5; \mu_2 = -0.5$ ,  
 433 the prior mean that fits these constraints will have a mean of zero because  
 434 it is set at the boundary of the admissible parameter space. For the com-  
 435 putation of (4.5), data are replicated based on  $\theta$  from a prior distribution  
 436 with  $\mu_0 = 0$ . These replicated data are adequately summarized by  $\bar{\mu}_1$  and  
 437  $\bar{\mu}_2$ . However, what is not accounted for in (4.5) is that the observed data  
 438  $\mathbf{y}$  are not adequately summarized by  $\bar{\mu}_1$  and  $\bar{\mu}_2$ . This leads to situations  
 439 where the loss function in (4.5) has a preference for ‘bad’ fitting inequality  
 440 constrained hypotheses.

441 In conclusion, neither the *prior* DIC, nor the *posterior* DIC are proper  
 442 model selection tools for the evaluation of inequality constrained hypotheses.  
 443 In the next section the prior predictive loss function will be adjusted such

444 that its estimate, the PIC, can be used to select the best of a set of equality  
 445 and inequality constrained hypotheses.

## 446 **5 A New Loss Function for the Evaluation of In-** 447 **equality Constrained Hypotheses**

448 The solution of the aforementioned problem (i.e. that neither the *prior* DIC,  
 449 nor the *posterior* DIC are proper model selection tools for the evaluation of  
 450 inequality constrained hypotheses) is to adjust the loss function that is used  
 451 to select the best hypothesis such that it also accounts for the agreement  
 452 between  $\bar{\boldsymbol{\theta}}_y$  and  $\mathbf{y}$ . The loss function in (4.5) can be rewritten as

$$-2 \operatorname{E}_{h(\boldsymbol{\theta})} \left\{ \operatorname{E}_{f(\mathbf{x}|\boldsymbol{\theta})} [\log f(\mathbf{x} | \bar{\boldsymbol{\theta}}_y)] \right\} + \log f(\mathbf{y} | \bar{\boldsymbol{\theta}}_y) \quad (5.1)$$

$$\approx \operatorname{E}_{h(\boldsymbol{\theta})} [-2 \log f(\mathbf{y} | \boldsymbol{\theta})]. \quad (5.2)$$

453 The new loss function determines not only how well replicated data fit with a  
 454 certain hypothesis (the term between accolades in 5.1), but it also determines  
 455 how well a hypothesis fits the data (the second term between accolades in  
 456 5.1). It is approximated by the third term of the *prior* DIC and is our final  
 457 model selection tool, to be called Prior Information Criterium (PIC) given  
 458 by (5.2).

### 459 **5.1 Behavior of the PIC in Constrained Model Selection**

460 In Figure 3 the PIC values for populations 1-7 of Example 1 are shown. As  
 461 can be seen, the PIC chooses for  $H_1$  as the best hypothesis in situations  
 462 where this hypothesis is true in the population, see populations 1-3. The  
 463 PIC chooses for  $H_2$  as the best hypothesis where this hypothesis is strongly  
 464 supported by the population values, see populations 4 and 5. Finally, the  
 465 PIC chooses for the unconstrained hypothesis,  $H_0$ , where the (in)equality  
 466 constraints for both  $H_1$  and  $H_2$  are not supported by the data, see popula-  
 467 tions 6 and 7. These results makes the PIC outperform both the *posterior*  
 468 and *prior* DIC in all situations.

### 469 **5.2 Influence of Prior Specification**

470 Since the specification of the prior has an impact on the results, we evaluated  
 471 the influence of the prior specifications on the PIC. To do so, we performed  
 472 a simulation study where  $\mu_0$ ,  $\tau_0^2$ ,  $v_0$  and  $\sigma_0$  were varied across populations.

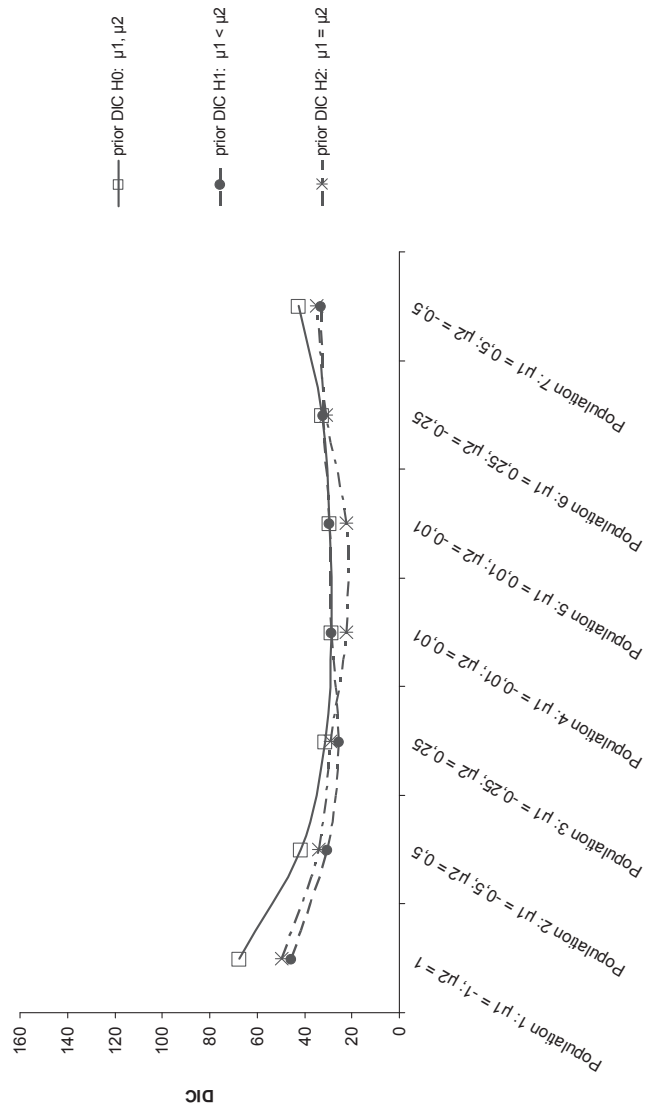


Figure 2: Values of the Prior DIC for  $H_0$ ,  $H_1$  and  $H_2$ ; and for populations 1-7 of Example 1.

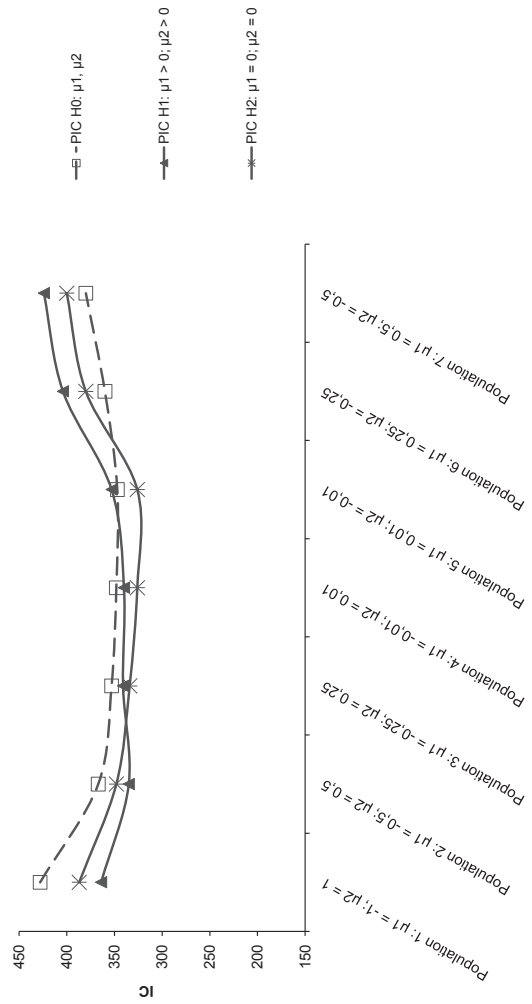


Figure 3: The PIC for populations 1-7 of Example 1.

473 We evaluated  $H_0$ ,  $H_1$  and  $H_2$  for populations 1, 4, and 7 with: (1)  $\mu_0 - 1$ ,  
474  $\mu_0 + 0$  and  $\mu_0 + 1$ ; (2)  $\tau_0 \times .5$ ,  $\tau_0 \times 1$ , and  $\tau_0 \times 5$ ; (3)  $\nu_0 = 2$  and  $\nu_0 = 5$ ; (4)  
475  $\sigma_0 \times .5$ ,  $\sigma_0 \times 1$ , and  $\sigma_0 \times 5$ .

476 The results are presented in Table 1 with in bold the correct conclusions.  
477 As can be seen, the specification of the prior influences the results. However,  
478 as can be seen for different prior specifications the influence is mainly on the  
479 height of PIC and not the relative ordering of  $\text{PIC}_{H_0}$ ,  $\text{PIC}_{H_1}$ , and  $\text{PIC}_{H_2}$ .

### 480 5.3 PIC versus Marginal Likelihood

481 The PIC is related to the marginal likelihood (ML) which is given by

$$\text{ML} \approx -2 \log E_{h_t(\theta)} [f(\mathbf{y} \mid \theta)] \quad (5.3)$$

482 The difference between (5.2) and (5.3) is the position of the log: inside  
483 (PIC) or outside (ML) the expectation. If within the constrained model  
484  $h_1(\cdot) = c \times h_0(\cdot)$ , see Equation (3.7), and under the further assumptions made  
485 about encompassing and constrained priors made in this paper then the  
486 relation between the PIC and ML shows a monotone relation. To exemplify  
487 this relation, we performed a small simulation study. In Figure 4  $\text{PIC}_1 - \text{PIC}_2$   
488 and  $\text{ML}_1 - \text{ML}_2$  are displayed for populations 1-7 of Example 1. As can be  
489 seen, there is a monotone relation between both selection tools. So, the  
490 PIC is related to the marginal likelihood approach, which is often used for  
491 inequality constrained model selection (see for example, Klugkist et al., 2005;  
492 Mulder, Hoijsink, & Klugkist, 2009).

## 493 6 Examples Reconsidered

494 After we have evaluated the performance of the *posterior* DIC, the *prior*  
495 DIC, the PIC and the ML for Example 1, it is now time to reconsider the  
496 other examples. For Examples 2 and 3 we only consider two populations:  
497 one population in agreement with the inequality constrained hypothesis and  
498 one population not in agreement with the constraints.

### 499 6.1 Example 2 continued

500 Let us return to Example 2 with  $H_0 : \mu_1, \mu_2$  and  $H_1 : \mu_2 > 0, \mu_1 > 0$ . To  
501 evaluate  $H_0$  and  $H_1$ , we performed a small simulation study where data  
502 sets from two different populations were considered. Population 1 satisfy  
503 the constraints of  $H_1$  and population 2 is not in agreement with  $H_1$ . The

		$\mu_0 - 1$		$\mu_0$		$\mu_0 + 1$			
	$H_0$	$H_1$	$H_2$	$H_0$	$H_1$	$H_2$	$H_0$	$H_1$	$H_2$
Population 1	203	<b>180</b>	203	182	<b>164</b>	183	202	<b>181</b>	202
Population 4	186	185	<b>178</b>	144	187	<b>136</b>	187	187	<b>178</b>
Population 7	<b>189</b>	200	228	<b>155</b>	167	187	<b>188</b>	199	228
		$\tau_0 \times .5:$		$\tau_0 \times 1$		$\tau_0 \times 5$			
	$H_0$	$H_1$	$H_2$	$H_0$	$H_1$	$H_2$	$H_0$	$H_1$	$H_2$
Population 1	185	<b>160</b>	186	182	<b>164</b>	183	292	<b>270</b>	195
Population 4	150	151	<b>143</b>	144	143	<b>136</b>	215	214	<b>207</b>
Population 7	<b>162</b>	174	175	<b>155</b>	167	162	<b>236</b>	257	238
		$v_0 = 2$		$v_0 = 5$					
	$H_0$	$H_1$	$H_2$	$H_0$	$H_1$	$H_2$			
Population 1	182	<b>160</b>	186	168	<b>145</b>	186			
Population 4	144	143	<b>136</b>	130	129	<b>127</b>			
Population 7	<b>155</b>	167	162	<b>141</b>	152	146			
		$\sigma_0^2 \times .5$		$\sigma_0^2 \times 1$		$\sigma_0^2 \times 5$			
	$H_0$	$H_1$	$H_2$	$H_0$	$H_1$	$H_2$	$H_0$	$H_1$	$H_2$
Population 1	213	<b>167</b>	246	183	<b>161</b>	186	200	<b>196</b>	177
Population 4	167	165	<b>156</b>	144	143	<b>136</b>	169	169	<b>159</b>
Population 7	<b>180</b>	203	218	<b>177</b>	181	179	<b>177</b>	181	179

Table 1: PIC values for different prior specifications. The bold numbers represent the hypothesis that is preferred by the PIC.

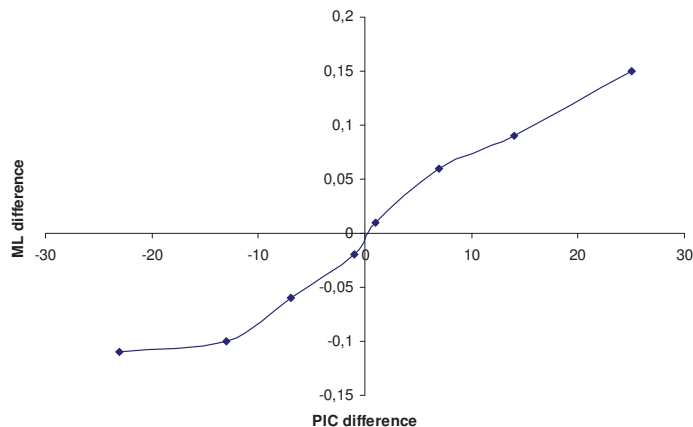


Figure 4: The differences between  $H_1$  and  $H_2$  are displayed for both the PIC and the ML for populations of Example 1.

504 two data sets were constructed in such a way that the sample means and  
 505 variance-covariance matrix are exactly equal to the population parameters  
 506 ( $\rho = .4; \sigma_1^2 = 1; \sigma_2^2 = 1; n = 40$ ). For each of these data sets, we computed  
 507 the *posterior* DIC, the *prior* DIC, and the PIC for  $H_0$  and  $H_1$ .

508 According to Mulder, Hoijtink, and Klugkist (2009) the prior distri-  
 509 bution,  $\theta_n = h_0(\mu_{11}, \mu_{12}, \mu_{21}, \mu_{22}) h_0(\Sigma)$ , can be given by a multivariate  
 510 normal distribution for the means and an inverse Wishart distribution for  
 511 the variance-covariance matrix

$$h_0(\mu_{11}, \mu_{12}, \mu_{21}, \mu_{22}, \Sigma) = MVN(\boldsymbol{\mu} | \boldsymbol{\mu}_0, \boldsymbol{\tau}_0^2) \times W^{-1}(\Sigma | \nu_0, \boldsymbol{\Sigma}_0), \quad (6.1)$$

512 where  $\boldsymbol{\mu} = \{\mu_{11}, \mu_{12}, \mu_{21}, \mu_{22}\}$  and  $\boldsymbol{\mu}_0 = \{\mu_0, \mu_0, \mu_0, \mu_0\}$ . For the Inverse  
 513 Wishart, we used  $\nu_0 = 3$  and for  $\boldsymbol{\Sigma}_0$ , which is the scale matrix, we used

$$\begin{bmatrix} \sigma_0^2 & 0 \\ 0 & \sigma_0^2 \end{bmatrix}, \quad (6.2)$$

514 For population 1 with  $\mu_1 = 1$  and  $\mu_2 = 1$ , the priors are  $\mu_0 = 0$ ,  $\tau_0 =$   
 515  $0.98$ ,  $\nu_0 = 3$  and  $\sigma_0^2 = 3.95$ . The results are shown in Table 2. As is  
 516 illustrated in Table 2, the situation for this example is analogous to Example  
 517 1. Analogously to Example 1, the *prior* DIC does not correctly distinguish  
 518  $H_0$  and  $H_1$  because the loss function does not take ‘bad’ fitting hypotheses  
 519 into account.

			<i>post.</i> DIC	<i>prior</i> DIC	PIC
Example 2	Population 1: $\mu_1 = 1, \mu_2 = 1$	$H_0$	234	489	428
		$H_1$	234	437	364
	Population 2: $\mu_1 = -.5, \mu_2 = -.5$	$H_0$	360	296	370
		$H_1$	404	294	406
Example 3	Population 1: $\beta_2 = 0.2$	$H_0$	275	919	820
		$H_1$	275	916	754
	Population 2: $\beta_2 = -0.2$	$H_0$	275	922	819
		$H_1$	411	924	960

Table 2: Results for Example 2 and 3.

## 520 6.2 Example 3 continued

521 For Example 3 we compared,  $H_0 : \beta_2$  and  $H_1 : \beta_2 > 0$ . Analogously to Ex-  
522 ample 1 and 2, the *posterior* and *prior* DIC do not correctly distinguish  $H_0$   
523 and  $H_1$  for Example 3. We also performed a small simulation study to eval-  
524 uate  $H_0$  and  $H_1$ . Data sets from two different populations were considered,  
525 see Table 2, where population 1 satisfy the constraints of  $H_1$  and population  
526 2 is not in agreement with  $H_1$ . The two data sets were constructed in such  
527 a way that the sample means and variance-covariance matrix are exactly  
528 equal to the population parameters ( $\beta_0 = 1.0; \beta_1 = 0.5; \beta_2 = 0.2; n = 50$ ).  
529 The prior parameters we used are  $\beta_1 = 0, \beta_2 = 0$ , and  $\sigma_0^2 = 1.95$ . For each  
530 of these data sets, we computed the *posterior* DIC, the *prior* DIC, and the  
531 PIC for  $H_0$  and  $H_1$ . The results are shown in Table 2 and it can be seen  
532 that the PIC outperforms the *posterior* and *prior* DIC.

## 533 6.3 Real-life Example 1

534 We re-evaluated the hypotheses given in (2.7). In Table 3 group means and  
535 standard deviations (SD) are provided. We computed the *posterior* DIC,  
536 the *prior* DIC, and the PIC for  $H_0, H_1$  and  $H_2$ . The results of the model  
537 selection procedure are presented in Table 4. As can be seen in this table  
538 the *posterior* DIC is indifferent for all hypotheses, whereas both the *prior*  
539 DIC and the PIC choose for  $H_2$ . This result can be confirmed when looking  
540 at the group means in Table 3 where  $\mu_{22}$  is larger than  $\mu_{21}$  and  $\mu_{11}$  is close  
541 to  $\mu_{12}$ .

542 The theoretical conclusion is that there is support for a domain shift  
543 in the judgement about hypothetical situations. That is, for pupils that



Table 3: Descriptive Statistics for real-life example 1 ( $n_1 = 38; n_2 = 97; \rho = .52$ )

	Mean	SD
$\mu_{11}$	5.37	1.23
$\mu_{12}$	5.68	1.62
$\mu_{21}$	5.27	1.27
$\mu_{22}$	6.71	2.14

Table 4: Model Selection Results for the Real-life data 1.

Hypothesis	<i>post.</i> DIC	<i>prior</i> DIC	PIC
$H_0$	935	1976	1044
$H_1$	935	1952	1023
$H_2$	935	1803	872

544 reported to have conducted some delinquent behavior (i.e. aggression), in  
 545 the same hypothetical situation, they will judge it to be more morally ac-  
 546 cepted compared to adolescents that did not report to conduct the same  
 547 behavior. However, in hypothetical situations concerning other delinquent  
 548 behavior that was not reported by these same adolescents (i.e. vandalism),  
 549 they will judge the hypothetical situation to be equally morally condemnable  
 550 as adolescents that did not report any antisocial behavior.

## 551 6.4 Real-life Example 2

552 We evaluated the hypothesis given in (2.10) using the *posterior* DIC, the  
 553 *prior* DIC, and the PIC. The results are shown in Table 5. As can be seen in  
 554 this table the *posterior* DIC fails to correctly distinguish the hypotheses of  
 555 interest, whereas both the *prior* DIC and the PIC choose for  $H_2$  as the best  
 556 hypothesis. This result can be confirmed when looking at the group means  
 557 in Table 6 where both  $\beta_2 Age^2$  as well as  $\beta_4 Age^2$  are both smaller than zero  
 558 and  $\beta_2 Age^2$  is smaller than  $\beta_4 Age^2$ .

559 The theoretical conclusion is that the relation between on the one hand  
 560 age, and on the other hand either time to complete a Ph.D. trajectory  
 561 or the gap between planned and actual project time are both non-linear.  
 562 Moreover, this non-linear effect is stronger for time to complete a Ph.D.  
 563 trajectory compared to the gap. This might be due to the fact that Ph.D  
 564 candidates in the middle thirties take more time to finish their Ph.D thesis,  
 565 but they also plan extra time.

Table 5: Model Selection Results for the Real-life data 2.

Hypothesis	<i>post.</i> DIC	<i>prior</i> DIC	PIC
$H_0$	4869	1862	1896
$H_1$	4872	1839	1855
$H_2$	4884	1834	1840

Table 6: Descriptive Statistics real-life example 2.

	Mean	SD
$\beta_{01}$	61.58	0.79
$\beta_1 Age$	2.88	0.27
$\beta_2 Age^2$	-0.96	0.01
$\beta_{02}$	10.54	0.82
$\beta_3 Age$	1.43	0.29
$\beta_2 Age^2$	-0.04	0.01

## 566 7 Conclusion

567 The main message of the current paper is: (1) although the DIC (Spiegelhalter  
568 et al., 2002) is often used in model selection, do not use it when evaluating  
569 inequality constrained hypotheses, better use the PIC which is derived in  
570 the current paper; and (2) the PIC is related to the marginal likelihood ap-  
571 proach, which is often used for inequality constrained model selection (see  
572 for example, Klugkist et al., 2005; Mulder, Hoijsink, & Klugkist, 2009). We  
573 showed how to obtain the *prior* DIC based on the derivation of the *poste-*  
574 *rior* DIC presented in Spiegelhalter et al. (2002). The point of departure  
575 for the *prior* DIC is the same as for the *posterior* DIC, namely the expected  
576 loss. The derivation of the *prior* DIC is provided and the choice for the  
577 prior distribution, which is based on training data is motivated (see also  
578 Mulder, Hoijsink, & Klugkist, 2009). Its performance is illustrated using  
579 examples and we showed that the *prior* DIC can be used to choose between  
580 a set of constrained hypotheses if the population from which the data are  
581 generated is fully in agreement with the most constrained hypothesis, where  
582 the *posterior* DIC failed to do so. However, the *prior* DIC fails to choose  
583 between a set of inequality constrained hypotheses if the population is *not*  
584 in agreement with the constrained hypothesis.

585 In conclusion, neither the *prior* DIC, nor the *posterior* DIC are proper  
586 model selection tools for the evaluation of inequality constrained hypotheses.  
587 To accommodate for this, the loss function that is minimized by the *prior*

588 DIC was adjusted. The proposed loss function determines not only how well  
589 replicated data fit with a certain hypothesis, but it also determines how well  
590 a hypothesis fits the data. It is approximated by a new model selection tool,  
591 the Prior Information Criterium (PIC). We demonstrated with examples  
592 that the PIC is able to select the best of a set of (in)equality constrained  
593 hypotheses. More research is needed to evaluate under what conditions the  
594 PIC is expected to work well and under what other conditions is it expected  
595 to fail. However, since we showed that the marginal likelihood is highly  
596 related to the PIC, we expect that the PIC behaves similar as the marginal  
597 likelihood approach. The current paper adds to the growing body of evidence  
598 that classical model selection tools, like AIC, BIC, MDL and now also the  
599 DIC, are not equipped to deal with inequality constraints and offers a viable  
600 alternative.

## 601 References

- 602 Akaike, H. (1973). Information theory as an extension of the maximum  
603 likelihood principle. In B. N. Petrov & F. Csaki (Eds.), *Second inter-*  
604 *national symposium on information theory* (p. 267 - 281). Budapest:  
605 Akademiai Kiado.
- 606 Anraku, K. (1999). An information criterion for parameters under a simple  
607 order restriction. *Journal of the Royal Statistical Society, series B*,  
608 *86*, 141-152.
- 609 Balasubramanian, V. (2005). Minimum description length. theory and ap-  
610 plications. In P. J. Grunwald, I. J. Myung, & M. A. Pitt (Eds.),  
611 (p. 81-99). MIT Press: Boston.
- 612 Barlow, R. E., Bartholomew, D. J., Bremner, H. M., & Brunk, H. D. (1972).  
613 *Statistical inference under order restrictions*. New York: Wiley.
- 614 Box, G. E. P. (1980). Sampling and bayes inference in scientific modelling  
615 and robustness. *Journal of the Royal Statistical Society, series A*, *143*,  
616 383 - 430.
- 617 Gelman, A., Carlin, J. B., Stern, H. S., & Rubin, D. B. (2004). *Bayesian*  
618 *data analysis* (2nd ed.). London: Chapman&HallCRC.
- 619 Gelman, A., Meng, X., & Stern, H. (1996). Posterior predictive assessment  
620 of model fitness via realized discrepancies (with discussion). *Statistica*  
621 *Sinica*, *6*, 733 -807.
- 622 Grunwald, P. D., Myung, I. J., & Pitt, M. a. (2005). *Advances in min-*  
623 *imum description length. theory and applications*. The MIT Press:  
624 Cambridge.

- 625 Hoijtink, H., Klugkist, I., & Boelen, P. A. (2008). *Bayesian evaluation of*  
626 *informative hypotheses*. New-York: Springer.
- 627 Kammers, M., Mulder, J., De Vignemont, F., & Dijkerman, H. (2009). The  
628 weight of representing the body: Addressing the potentially indefinite  
629 number of body representations in healthy individuals. *Experimental*  
630 *Brain Research, Published on-line, 22 sept. 2009*.
- 631 Kass, R. E., & Raftery, A. E. (1995). Bayes factors. *Journal of the American*  
632 *Statistical Association, 90*, 773-795.
- 633 Klugkist, I., Laudy, O., & Hoijtink, H. (2005). Inequality constrained anal-  
634 ysis of variance: A bayesian approach. *Psychological Methods, 10*, 477  
635 - 493.
- 636 Kuiper, R. M., & Hoijtink, H. (2010). Comparisons of means using ex-  
637 ploratory and confirmatory approaches. *Psychological Methods, 15*,  
638 69-86.
- 639 Leenders, I., & Brugman, D. (2005). Moral/non-moral domain shift in  
640 young adolescents in relation to delinquent behaviour. *British Journal*  
641 *of Developmental Psychology, 23*, 65 - 79.
- 642 Lunn, D., Thomas, A., Best, N., & Spiegelhalter, D. (2000). Winbugs - a  
643 bayesian modelling framework: concepts, structure, and extensibility.  
644 *Statistics and Computing, 10*, 325 - 337.
- 645 Meeus, W., Van de Schoot, R., Keijsers, L., Schwartz, S. J., & Branje, S.  
646 (2010). On the progression and stability of adolescent identity forma-  
647 tion. a five-wave longitudinal study in early-to-middle and middle-to-  
648 late adolescence. *Child Development*.
- 649 Mulder, J., Hoijtink, H., & Klugkist, I. (2009). Equality and inequality  
650 constrained multivariate linear models: Objective model selection us-  
651 ing constrained posterior priors. *Journal of Statistical Planning and*  
652 *Inference, 140*, 887-906.
- 653 Mulder, J., Klugkist, I., Van de Schoot, R., Meeus, W., Selfhout, M., & Hoi-  
654 jtink, H. (2009). Bayesian model selection of informative hypotheses  
655 for repeated measurements. *Journal of Mathematical Psychology, 53*,  
656 530-546.
- 657 Muthen, L. K., & Muthen, B. (2010). *Mplus users guide. sixth edition*. Los  
658 Angeles: Muthen & Muthen.
- 659 Myung, J. (2003). Tutorial on maximum likelihood estimation. *Journal of*  
660 *Mathematical Psychology 00, 47*, 90-100.
- 661 Press, S. J. (2005). *Applied multivariate analysis: Using bayesian and*  
662 *frequentist methods of inference (2nd ed)*. Malabar, FL: Krieger.
- 663 R Development Core Team. (2006). R:a language and environment for statisti-  
664 cal computing (no. isbn 3-900051-07-0) [Computer software manual].

- 665 Rasbash, J., Charlton, C., Browne, W., Healy, M., & Cameron, B. (2009).  
666 *Mlwin version 2.1*. Centre for Multilevel Modelling, University of Bris-  
667 toln.
- 668 Robertson, T., Wright, F. T., & Dykstra, R. L. (1988). *Order restricted*  
669 *statistical inference*. New York : Wiley.
- 670 Romeijn, J.-W., Van de Schoot, R., & Hoijsink, H. (2011). One size does  
671 not fit all: Derivation of an adapted bic. In D. Dieks et al. (Eds.),  
672 *Probabilities, laws, and structures* (p. xxx-xxx). Berlin: Springer.
- 673 Schwarz. (1978). Estimating the dimension of a model. *Annals of Statistics*,  
674 *6*, 461-464.
- 675 Silvapulle, M. J., & Sen, P. K. (2004). *Constrained statistical inference:*  
676 *Order, inequality, and shape constraints*. London: John Wiley Sons.
- 677 Sober, E. (2006). Parsimony. In J. Pfeifer & S. Sarkar (Eds.), *The phi-*  
678 *losophy of science: An encyclopedia* (Vol. 2, p. 530-541). New York:  
679 Routledge.
- 680 Sonneveld, H., Yerkes, M., & Van de Schoot, R. (2009). *Ph.d. trajecto-*  
681 *ries and labor market mobility: A survey of doctoral graduates in the*  
682 *netherland*. (Tech. Rep.). Report for Netherlands Centre for Graduate  
683 and Research Schools in the Netherlands and was subsidized by the  
684 Dutch Ministry of Education, Culture and Science.
- 685 Spiegelhalter, D. J., Best, N. G., Carlin, B. P., & Van Der Linde, A. (2002).  
686 Bayesian measures of model complexity and fit. *Journal of Royal*  
687 *Statistical Society, series B*, *64*, 583-639.
- 688 Van de Schoot, R., Hoijsink, H., & Deković, M. (2010). Testing inequality  
689 constrained hypotheses in sem models. *Structural Equation Modeling*,  
690 *17*, 443-463.
- 691 Van de Schoot, R., Romeijn, J.-W., & Hoijsink, H. (2011). Moving beyond  
692 traditional null hypothesis testing: Evaluating expectations directly.  
693 *Frontiers in Quantitative Psychology and Measurement*.
- 694 Van de Schoot, R., & Wong, T. (2010). Do antisocial young adults have a  
695 high or a low level of self-concept? *Self and Identity*.
- 696 Van Well, S., Kolk, A. M., & Klugkist, I. (2009). The relationship between  
697 sex, gender role identification, and the gender relevance of a stres-  
698 sor on physiological and subjective stress responses: Sex and gender  
699 (mis)match effects. *International Journal of Psychophysiology*, *32*,  
700 427-449.

## 701 A Derivation of Prior Predictive DIC

702 In this appendix we show how to obtain the *prior* DIC based on the deriva-  
 703 tion of the *posterior* DIC presented in Spiegelhalter et al. (2002). The point  
 704 of departure for the *prior* DIC is the same as for the *posterior* DIC, namely  
 705 the expected loss given in (3.1). However, to deal with the unknown pa-  
 706 rameters  $\boldsymbol{\theta}^*$ , we take the expectation with respect to the *prior* distribution,  
 707  $h(\boldsymbol{\theta})$ , instead of the *posterior* expectation of the expected loss:

$$\begin{aligned} \mathbb{E}_{h(\boldsymbol{\theta})} \left\{ \mathbb{E}_{f(\mathbf{x}|\boldsymbol{\theta})} [-2 \log f(\mathbf{x} | \bar{\boldsymbol{\theta}}_y)] \right\} = \\ -2 \log f(\mathbf{y} | \bar{\boldsymbol{\theta}}_y) + \mathbb{E}_{h(\boldsymbol{\theta})} [c(\mathbf{y}, \boldsymbol{\theta}, \bar{\boldsymbol{\theta}}_y)] . \end{aligned} \quad (\text{A.1})$$

708 The main problem now, is to find an expression for the second term on the  
 709 right hand side in (A.1). Using  $D(\mathbf{a}, \mathbf{b}) = -2 \log f(\mathbf{a} | \mathbf{b})$ ,  $c(\mathbf{y}, \boldsymbol{\theta}, \bar{\boldsymbol{\theta}}_y)$  in  
 710 (A.1) can be rewritten to

$$\begin{aligned} c(\mathbf{y}, \boldsymbol{\theta}, \bar{\boldsymbol{\theta}}_y) &= \mathbb{E}_{f(\mathbf{x}|\boldsymbol{\theta})} [D(\mathbf{x}, \bar{\boldsymbol{\theta}}_y) - D(\mathbf{y}, \bar{\boldsymbol{\theta}}_y)] \\ &= \mathbb{E}_{f(\mathbf{x}|\boldsymbol{\theta})} [D(\mathbf{x}, \bar{\boldsymbol{\theta}}_y) - D(\mathbf{x}, \boldsymbol{\theta})] \\ &+ \mathbb{E}_{f(\mathbf{x}|\boldsymbol{\theta})} [D(\mathbf{x}, \boldsymbol{\theta}) - D(\mathbf{y}, \boldsymbol{\theta})] \\ &+ D(\mathbf{y}, \boldsymbol{\theta}) - D(\mathbf{y}, \bar{\boldsymbol{\theta}}_y) . \end{aligned} \quad (\text{A.2})$$

711 Now,  $D(\mathbf{x}, \bar{\boldsymbol{\theta}}_y)$  in (A.2) can be approximated by taking a second order Taylor  
 712 expansion about  $\boldsymbol{\theta}$ ,

$$\begin{aligned} D(\mathbf{x}, \bar{\boldsymbol{\theta}}_y) \approx -2 \log f(\mathbf{x} | \boldsymbol{\theta}) - 2 \left\{ \frac{\partial \log f(\mathbf{x}|\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right\}^T (\bar{\boldsymbol{\theta}}_y - \boldsymbol{\theta}) - \\ - (\bar{\boldsymbol{\theta}}_y - \boldsymbol{\theta})^T \left\{ \frac{\partial^2 \log f(\mathbf{x}|\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \right\} (\bar{\boldsymbol{\theta}}_y - \boldsymbol{\theta}) . \end{aligned} \quad (\text{A.3})$$

713 Since  $-2 \log f(\mathbf{x} | \boldsymbol{\theta})$  is equal to  $D(\mathbf{x}, \boldsymbol{\theta})$  and the expectation of the second  
 714 term on the right hand side of (A.3) with respect to  $f(\mathbf{x} | \boldsymbol{\theta})$  is zero (p. 604  
 715 Spiegelhalter et al., 2002),

$$\begin{aligned} \mathbb{E}_{f(\mathbf{x}|\boldsymbol{\theta})} [D(\mathbf{x}, \bar{\boldsymbol{\theta}}_y) - D(\mathbf{x}, \boldsymbol{\theta})] \approx \\ \mathbb{E}_{f(\mathbf{x}|\boldsymbol{\theta})} \left[ -(\bar{\boldsymbol{\theta}}_y - \boldsymbol{\theta})^T \left\{ \frac{\partial^2 \log f(\mathbf{x}|\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \right\} (\bar{\boldsymbol{\theta}}_y - \boldsymbol{\theta}) \right] . \end{aligned} \quad (\text{A.4})$$

716 The expression on the right hand side of (A.4) can be rewritten as  $\text{tr} \{ \mathbf{I}(\boldsymbol{\theta}) (\bar{\boldsymbol{\theta}}_y -$   
 717  $\boldsymbol{\theta}) (\bar{\boldsymbol{\theta}}_y - \boldsymbol{\theta})^T \}$  and since  $\mathbf{x}$  and  $\mathbf{y}$  stem from the same data generating mech-  
 718 anism, the Fisher information matrix  $\mathbf{I}(\boldsymbol{\theta})$  can be approximated by the ob-  
 719 served Fisher information matrix,  $\mathbf{I}(\bar{\boldsymbol{\theta}}_y)$  (p. 604 Spiegelhalter et al., 2002),

720 where  $\mathbf{I}(\bar{\boldsymbol{\theta}}_y) = -\partial^2 \log f(\mathbf{y} \mid \bar{\boldsymbol{\theta}}_y) / \partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T$ . Using  $E\{\text{tr}(\cdot)\} = \text{tr}\{E(\cdot)\}$ , the  
 721 prior expectation of  $c(\mathbf{y}, \boldsymbol{\theta}, \bar{\boldsymbol{\theta}}_y)$  can now be approximated by:

$$\begin{aligned} E_{h(\boldsymbol{\theta})}[c(\mathbf{y}, \boldsymbol{\theta}, \bar{\boldsymbol{\theta}}_y)] &\approx \text{tr}\{\mathbf{I}(\bar{\boldsymbol{\theta}}_y)\boldsymbol{\Lambda}\} + \\ &+ E_{h(\boldsymbol{\theta})}\left\{E_{f(\mathbf{x}|\boldsymbol{\theta})}[D(\mathbf{x}, \boldsymbol{\theta}) - D(\mathbf{y}, \boldsymbol{\theta})]\right\} + d, \end{aligned} \quad (\text{A.5})$$

722 where  $\boldsymbol{\Lambda} = E_{h(\boldsymbol{\theta})}[(\bar{\boldsymbol{\theta}}_y - \boldsymbol{\theta})(\bar{\boldsymbol{\theta}}_y - \boldsymbol{\theta})^T]$  denotes the variation in the prior distri-  
 723 bution around  $\bar{\boldsymbol{\theta}}_y$ . The last term on the right hand side of (A.5) is defined  
 724 as

$$\begin{aligned} d &= E_{h(\boldsymbol{\theta})}[D(\mathbf{y}, \boldsymbol{\theta})] - E_{h(\boldsymbol{\theta})}[D(\mathbf{y}, \bar{\boldsymbol{\theta}}_y)] \\ &= E_{h(\boldsymbol{\theta})}[D(\mathbf{y}, \boldsymbol{\theta})] - D(\mathbf{y}, \bar{\boldsymbol{\theta}}_y). \end{aligned} \quad (\text{A.6})$$

725 To show that  $\text{tr}\{\mathbf{I}(\bar{\boldsymbol{\theta}}_y)\boldsymbol{\Lambda}\}$  is approximately equal to  $d$ , we use a second order  
 726 Taylor expansion about  $\bar{\boldsymbol{\theta}}_y$ :

$$\begin{aligned} E_{h(\boldsymbol{\theta})}[D(\mathbf{y}, \boldsymbol{\theta})] &\approx D(\mathbf{y}, \bar{\boldsymbol{\theta}}_y) + E_{h(\boldsymbol{\theta})}\left[-2\left\{\frac{\partial \log f(\mathbf{y} \mid \bar{\boldsymbol{\theta}}_y)}{\partial \boldsymbol{\theta}}\right\}^T (\boldsymbol{\theta} - \bar{\boldsymbol{\theta}}_y) - \right. \\ &\quad \left. - (\boldsymbol{\theta} - \bar{\boldsymbol{\theta}}_y)^T \left\{\frac{\partial^2 \log f(\mathbf{y} \mid \bar{\boldsymbol{\theta}}_y)}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T}\right\} (\boldsymbol{\theta} - \bar{\boldsymbol{\theta}}_y)\right]. \end{aligned} \quad (\text{A.7})$$

727 Since,  $\bar{\boldsymbol{\theta}}_y \rightarrow \bar{\boldsymbol{\theta}}_{ML}$  for  $n \rightarrow \infty$ ,  $-2\left\{\frac{\partial \log f(\mathbf{y}|\bar{\boldsymbol{\theta}}_y)}{\partial \boldsymbol{\theta}}\right\}^T$  is asymptotically zero  
 728 (Gelman et al., 2004). This way,  $E_{h(\boldsymbol{\theta})}[D(\mathbf{y}, \boldsymbol{\theta})]$  can now be approximated  
 729 by

$$\begin{aligned} E_{h(\boldsymbol{\theta})}[D(\mathbf{y}, \boldsymbol{\theta})] &\approx D(\mathbf{y}, \bar{\boldsymbol{\theta}}_y) + E_{h(\boldsymbol{\theta})}\left[\text{tr}\left\{-\frac{\partial^2 \log f(\mathbf{y} \mid \bar{\boldsymbol{\theta}}_y)}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} (\boldsymbol{\theta} - \bar{\boldsymbol{\theta}}_y)(\boldsymbol{\theta} - \bar{\boldsymbol{\theta}}_y)^T\right\}\right] \\ &\approx D(\mathbf{y}, \bar{\boldsymbol{\theta}}_y) + \text{tr}\{\mathbf{I}(\bar{\boldsymbol{\theta}}_y)\boldsymbol{\Lambda}\} \end{aligned} \quad (\text{A.8})$$

730 To show that  $\text{tr}\{\mathbf{I}(\bar{\boldsymbol{\theta}}_y)\boldsymbol{\Lambda}\}$  is approximately equal to  $d$ ,  $D(\mathbf{y}, \bar{\boldsymbol{\theta}}_y)$  is subtracted  
 731 from both sides of (A.8)

$$\text{tr}\{\mathbf{I}(\bar{\boldsymbol{\theta}}_y)\boldsymbol{\Lambda}\} \approx E_{h(\boldsymbol{\theta})}[D(\mathbf{y}, \boldsymbol{\theta})] - D(\mathbf{y}, \bar{\boldsymbol{\theta}}_y) = d. \quad (\text{A.9})$$

732 Equation (A.5) then becomes

$$\begin{aligned} E_{h(\boldsymbol{\theta})}[c(\mathbf{y}, \boldsymbol{\theta}, \bar{\boldsymbol{\theta}}_y)] &\approx E_{h(\boldsymbol{\theta})}\left\{E_{f(\mathbf{x}|\boldsymbol{\theta})}[D(\mathbf{x}, \boldsymbol{\theta}) - D(\mathbf{y}, \boldsymbol{\theta})]\right\} + \\ &+ 2\left\{E_{h(\boldsymbol{\theta})}[D(\mathbf{y}, \boldsymbol{\theta})] - D(\mathbf{y}, \bar{\boldsymbol{\theta}}_y)\right\}. \end{aligned} \quad (\text{A.10})$$

733 The *prior* DIC can now be written as

$$\begin{aligned} & \mathbb{E}_{h(\boldsymbol{\theta})} \left\{ \mathbb{E}_{f(\mathbf{x}|\boldsymbol{\theta})} [-2 \log f(\mathbf{x} | \bar{\boldsymbol{\theta}}_y)] \right\} \approx \\ & \mathbb{E}_{h(\boldsymbol{\theta})} \left\{ \mathbb{E}_{f(\mathbf{x}|\boldsymbol{\theta})} [D(\mathbf{x}, \boldsymbol{\theta})] \right\} - D(\mathbf{y}, \bar{\boldsymbol{\theta}}_y) + \mathbb{E}_{h(\boldsymbol{\theta})} [D(\mathbf{y}, \boldsymbol{\theta})] \end{aligned} \quad (\text{A.11})$$

734 whereas, using the same notation, the *posterior* DIC can be written as

$$\begin{aligned} & \mathbb{E}_{h(\boldsymbol{\theta})} \left\{ \mathbb{E}_{f(\mathbf{x}|\boldsymbol{\theta})} [-2 \log f(\mathbf{x} | \bar{\boldsymbol{\theta}}_y)] \right\} \approx \\ & D(\mathbf{y}, \bar{\boldsymbol{\theta}}_y) + 2 \left\{ \mathbb{E}_{g(\boldsymbol{\theta}|\mathbf{y})} [D(\mathbf{y}, \boldsymbol{\theta})] - D(\mathbf{y}, \bar{\boldsymbol{\theta}}_y) \right\}. \end{aligned} \quad (\text{A.12})$$

## 735 B Simplifying the prior DIC for constrained hy- 736 potheses

737 Let  $H_t (t = 1, \dots, T)$  denote a hypothesis specified using constraints and let  
738  $H_0$  denote an unconstrained hypothesis. All hypotheses  $H_t$  are nested in  
739  $H_0$ . As we will prove in this section,  $\mathbb{E}_{h_t(\boldsymbol{\theta})} \left\{ \mathbb{E}_{f(\mathbf{x}|\boldsymbol{\theta})} [D(\mathbf{x}, \boldsymbol{\theta})] \right\}$  in (A.11)  
740 is constant between constrained hypotheses. In this context the *prior* DIC  
741 reduces to

$$\text{prior DIC} = C + 2 \log f(\mathbf{y} | \bar{\boldsymbol{\theta}}_y) + \mathbb{E}_{h_t(\boldsymbol{\theta})} [-2 \log f(\mathbf{y} | \boldsymbol{\theta})], \quad (\text{B.1})$$

742 where  $C = \mathbb{E}_{h_t(\boldsymbol{\theta})} \left\{ \mathbb{E}_{f(\mathbf{x}|\boldsymbol{\theta})} [-2 \log f(\mathbf{x} | \boldsymbol{\theta})] \right\}$  and can be ignored for all  $H_t$ .

### 743 B.1 Example 1 Continued

744 For Example 1,  $h_t(\boldsymbol{\theta}_c)h_t(\boldsymbol{\theta}_u) = h_t(\mu_1, \mu_2)h_t(\sigma^2)$  where  $h_t(\sigma^2)$  is the same,  
745 but  $h_t(\mu_1, \mu_2)$  differs across hypotheses because of the normalization of  
746 the prior distribution in Equation (3.6). In the remainder of this subsection  
747 we drop the subscript  $t$  to simplify the notation. We will prove that  
748  $\mathbb{E}_{h(\sigma^2)h(\mu_1, \mu_2)} \left\{ \mathbb{E}_{f(\mathbf{x}|\mu_1, \mu_2, \sigma^2)} [-2 \log f(\mathbf{x} | \mu_1, \mu_2, \sigma^2)] \right\}$  is constant over all  
749 hypotheses under consideration. When comparing constrained hypotheses  
750 we have to prove that the term within accolades is independent of  $\mu_1, \mu_2$ ,  
751 and  $\sigma^2$ . First using

$$f(\mathbf{x} | \mu_1, \mu_2, \sigma^2) = \left( \frac{1}{\sqrt{2\pi\sigma^2}} \right)^N \exp \left[ -\frac{1}{2} \frac{\sum_{i=1}^N (x_i - \mu_1 d_1 - \mu_2 d_2)^2}{\sigma^2} \right], \quad (\text{B.2})$$

752 the term being constant can be written as



$$\int_{\sigma^2} \int_{\mu_1, \mu_2} \int_{\mathbf{x}} 2N \log \sqrt{2\pi\sigma^2} \partial f(\mathbf{x} | \mu_1, \mu_2, \sigma^2) \partial h(\mu_1, \mu_2) \partial h(\sigma^2) + \quad (\text{B.3})$$

$$+ \int_{\sigma^2} \int_{\mu_1, \mu_2} \int_{\mathbf{x}} \sum_{i=1}^N \frac{(x_i - \mu_1 d_1 - \mu_2 d_2)^2}{\sigma^2} \partial f(\mathbf{x} | \mu_1, \mu_2, \sigma^2) \partial h(\mu_1, \mu_2) \partial h(\sigma^2) .$$

753 The first term of (B.3) is independent of  $\mu_1, \mu_2$ , and since  $h(\sigma^2)$  is the same  
754 for each hypothesis, the second term integrated over  $\sigma^2$  in (B.3) should be  
755 constant for every value for  $\sigma^2$  to render (B.3) constant. Let  $\mathbf{x} = \{\mathbf{x}_1, \mathbf{x}_2\}$   
756 denote subgroups with sample sizes  $N_1$  and  $N_2$  for  $\mathbf{x}_1$  and  $\mathbf{x}_2$ , respectively.  
757 Omitting the integral over  $\sigma^2$ , we can now rewrite the second term in (B.3)  
758 to

$$\int_{\mu_1} \int_{\mathbf{x}_1} \sum_{i=1}^{N_1} \frac{(x_i - \mu_1)^2}{\sigma^2} \partial f(\mathbf{x}_1 | \mu_1, \sigma^2) \partial h(\mu_1) + \quad (\text{B.4})$$

$$+ \int_{\mu_2} \int_{\mathbf{x}_2} \sum_{i=1}^{N_2} \frac{(x_i - \mu_2)^2}{\sigma^2} \partial f(\mathbf{x}_2 | \mu_2, \sigma^2) \partial h(\mu_2) .$$

759 Note, that for the first group in (B.4)  $x_i \sim N(\mu_1, \sigma^2)$  and for the second  
760 group  $x_i \sim N(\mu_2, \sigma^2)$ . Using  $x_i^* = \frac{x_i - \mu_1}{\sigma^2}$  with  $x_i^* \sim N(0, 1)$  in the first  
761 group, and  $x_i^* = \frac{x_i - \mu_2}{\sigma^2}$  with  $x_i^* \sim N(0, 1)$  in the second group, the integral  
762 over  $\mu_1$  and  $\mu_2$  drop out of (B.4):

$$\int_{\mathbf{x}_1^*} \sum_{i=1}^N (x_i^*)^2 \partial f(\mathbf{x}_i^* | 0, 1) + \int_{\mathbf{x}_2^*} \sum_{i=1}^N (x_i^*)^2 \partial f(\mathbf{x}_i^* | 0, 1) . \quad (\text{B.5})$$

763 Consequently, for every value of  $\sigma^2$ , (B.4) is independent of  $\mu_1, \mu_2$ . That  
764 is, for this example,  $E_{h(\sigma^2)h(\mu_1, \mu_2)} \left\{ E_{f(\cdot)} [-2 \log f(\cdot)] \right\}$  is constant over con-  
765 strained hypotheses.

## 766 B.2 Example 2 Continued

767 For Example 2,  $h_t(\boldsymbol{\theta}_c)h_t(\boldsymbol{\theta}_u) = h_t(\mu_1, \mu_2)h_t(\Sigma)$  where  $h_t(\Sigma)$  is the same, but  
768  $h_t(\mu_1, \mu_2)$  differs across hypotheses because of the normalization of the prior  
769 distribution in Equation (3.6). In the remainder of this subsection we drop  
770 the subscript  $t$  to simplify the notation. We now have to prove that the  
771 term between accolades in

$$E_{h(\mu_1, \mu_2)h(\sigma_{x_1}, \sigma_{x_2}, \rho)} \left\{ E_{f(\cdot)} [-2 \log f(\mathbf{x}_1, \mathbf{x}_2 | \mu_1, \mu_2, \sigma_{x_1}, \sigma_{x_2}, \rho)] \right\} \quad (\text{B.6})$$

772 is constant over hypotheses for  $\mu_1, \mu_2$ , and  $\Sigma$ . Using

$$f(\mathbf{x}_1, \mathbf{x}_2 \mid \mu_1, \mu_2, \sigma_{x_1}, \sigma_{x_2}, \rho) = \left( \frac{1}{2\pi\sigma_{x_1}\sigma_{x_2}\sqrt{1-\rho^2}} \right)^N \exp \left[ -\frac{1}{2(1-\rho^2)} \left\{ \frac{\sum_{i=1}^N (x_{1i} - \mu_1)^2}{\sigma_{x_1}^2} + \frac{\sum_{i=1}^N (x_{2i} - \mu_2)^2}{\sigma_{x_2}^2} - \frac{2\rho \sum_{i=1}^N (x_{1i} - \mu_1)(x_{2i} - \mu_2)}{\sigma_{x_1}\sigma_{x_2}} \right\} \right], \quad (\text{B.7})$$

773 (B.6) can be written as the sum of

$$\int_{\sigma_{x_1}, \sigma_{x_2}, \rho} \int_{\mu_1, \mu_2} \int_{\mathbf{x}_1, \mathbf{x}_2} 2N \log 2\pi\sigma_{x_1}\sigma_{x_2}\sqrt{1-\rho^2} \partial f(\mathbf{x}_1, \mathbf{x}_2 \mid \mu_1, \mu_2, \sigma_{x_1}, \sigma_{x_2}, \rho) \partial h(\mu_1, \mu_2) \partial h(\sigma_{x_1}, \sigma_{x_2}, \rho), \quad (\text{B.8})$$

774 and

$$\int_{\sigma_{x_1}, \sigma_{x_2}, \rho} \int_{\mu_1, \mu_2} \int_{\mathbf{x}_1, \mathbf{x}_2} \frac{1}{(1-\rho^2)} \left\{ \frac{\sum_{i=1}^N (x_{1i} - \mu_1)^2}{\sigma_{x_1}^2} + \frac{\sum_{i=1}^N (x_{2i} - \mu_2)^2}{\sigma_{x_2}^2} - \frac{2\rho \sum_{i=1}^N (x_{1i} - \mu_1)(x_{2i} - \mu_2)}{\sigma_{x_1}\sigma_{x_2}} \right\} \partial f(\mathbf{x}_1, \mathbf{x}_2 \mid \mu_1, \mu_2, \sigma_{x_1}, \sigma_{x_2}, \rho) \partial h(\mu_1, \mu_2) \partial h(\sigma_{x_1}, \sigma_{x_2}, \rho). \quad (\text{B.9})$$

775 Since  $h(\Sigma)$  is the same for each hypothesis, the integrals in (B.9) integrated  
776 over  $\sigma_{x_1}, \sigma_{x_2}, \rho$  should be constant for every value of  $h(\Sigma)$  to render (B.9)  
777 constant. Also, in this situation (B.8) is constant over constrained hypothe-  
778 ses. Using  $x_{1i}^* = \frac{x_{1i} - \mu_1}{\sigma}$  and  $x_{2i}^* = \frac{x_{2i} - \mu_2}{\sigma}$ , (B.9) can be rewritten into

$$\int_{\rho} \int_{\mathbf{x}_1^*, \mathbf{x}_2^*} \sum_{i=1}^N \frac{1}{(1-\rho^2)} \left\{ (x_{1i}^*)^2 + (x_{2i}^*)^2 - 2\rho^2 x_{1i}^* x_{2i}^* \right\} \partial f(\mathbf{x}_1^*, \mathbf{x}_2^* \mid 0, 0, 1, 1, \rho) \partial(\rho). \quad (\text{B.10})$$

779 Consequently, for every  $\Sigma$ , (B.9) is independent of  $\mu_1$  and  $\mu_2$ . That is, for  
780 this example,  $E_{h(\mu_1, \mu_2)h(\sigma_{x_1}, \sigma_{x_2}, \rho)} \left\{ E_{f(\cdot)} [-2 \log f(\cdot)] \right\}$  is constant over con-  
781 strained hypotheses.

### 782 B.3 Multivariate Models

783 Finally, consider a multivariate example with two groups with mean scores  
784 on two dependent variables:

$$\begin{aligned} y_{1i} &= \mu_{11}d_{ig1} + \mu_{12}d_{ig2} + \epsilon_{1i} \\ y_{2i} &= \mu_{21}d_{ig1} + \mu_{22}d_{ig2} + \epsilon_{2i}, \end{aligned} \quad (\text{B.11})$$

785 where  $\mu_{1\cdot}$  and  $\mu_{2\cdot}$  denote the mean score on  $y_1$  and  $y_2$  respectively and  
 786 where  $\mu_{\cdot 1}$  and  $\mu_{\cdot 2}$  denote the mean for group 1 and 2 respectively. Again,  
 787 group membership of a person is denoted by  $d_{ig} \in \{0, 1\}$  and the residuals are  
 788 assumed to be normally distributed with

$$\begin{bmatrix} \epsilon_{i1} \\ \epsilon_{i2} \end{bmatrix} \sim N(0, \Sigma), \Sigma = \begin{bmatrix} \sigma_{y_1}^2 & \rho\sigma_{y_1}\sigma_{y_2} \\ \rho\sigma_{y_1}\sigma_{y_2} & \sigma_{y_2}^2 \end{bmatrix}. \quad (\text{B.12})$$

789 Note that this example is a combination of (2.4) and (2.2). Also for con-  
 790 strained hypotheses in this multivariate example it can be proved that  
 791  $E_{h_t(\mu_{11}, \mu_{12}, \mu_{21}, \mu_{22})h_t(\Sigma)} \left\{ E_{f(\cdot)} [-2 \log f(\mathbf{y}_1, \mathbf{y}_2 \mid \mu_{11}, \mu_{12}, \mu_{21}, \mu_{22}, \Sigma)] \right\}$  is con-  
 792 stant over constrained hypotheses. Even so, using the same steps as pre-  
 793 sented in Section B.1 and B.2, it can be proved for the general multivari-  
 794 ate normal linear model (Press (2005), pp. 252-257) that  $E_{h_t(\boldsymbol{\theta})} \left\{ E_{f(\mathbf{x}|\boldsymbol{\theta})} \right.$   
 795  $\left. [-2 \log f(\mathbf{x} \mid \boldsymbol{\theta})] \right\}$  is constant over constrained hypotheses.