

# Regularized Bayesian estimation of generalized threshold regression models

Friederike Greb <sup>1, 2</sup>    Tatyana Krivobokova <sup>2, 3</sup>    Axel Munk <sup>3, 4</sup>  
Stephan von Cramon-Taubadel <sup>1</sup>  
Georg-August-Universität Göttingen

July 7, 2013

## Abstract

In this article we discuss estimation of generalized threshold regression models in settings when the threshold parameter lacks identifiability. In particular, if estimation of the regression coefficients is associated with high uncertainty and/or the difference between regimes is small, estimators of the threshold and, hence, of the whole model can be strongly affected. A new regularized Bayesian estimator for generalized threshold regression models is proposed. We derive conditions for superiority of the new estimator over the standard likelihood one in terms of mean squared error. Simulations confirm excellent finite sample properties of the suggested estimator, especially in the critical settings. The practical relevance of our approach is illustrated by two real-data examples already analyzed in the literature. *Key words and phrases: empirical Bayes, regularization, threshold identification.*

## 1 Introduction

1 Modeling a response variable as a linear combination of some covariates with regression  
2 coefficients that vary between (possibly several) regimes is known as threshold regression.  
3 The choice of regime is determined by a transition function, which depends on a transition  
4 variable as well as a threshold parameter. Transition functions can be either smooth  
5 (Van Dijk et al., 2002, provide a comprehensive overview) or step functions. In the

---

<sup>1</sup>Department of Agricultural Economics and Rural Development, Georg-August-Universität Göttingen, Platz der Göttinger Sieben 5, 37073 Göttingen, Germany

<sup>2</sup>Courant Research Centre "Poverty, Equity and Growth in Developing Countries", Georg-August-Universität Göttingen, Wilhelm-Weber-Str. 2, 37073 Göttingen, Germany

<sup>3</sup>Institute for Mathematical Stochastics, Georg-August-Universität Göttingen, Goldschmidstrasse 7, 37077 Göttingen, Germany

<sup>4</sup>Max Planck Institute for Biophysical Chemistry, Am Faßberg 11, 37077 Göttingen, Germany

6 following, we restrict attention to the latter. In principle, the response variable can  
7 follow any distribution from the exponential family. However, such generalized threshold  
8 regression models have only recently been formally introduced by Samia and Chan (2011),  
9 and most of the literature on threshold regression deals with models with a piecewise  
10 linear mean. In this article we concentrate on generalized regression models with regimes  
11 controlled by a step transition function and refer to such models as generalized threshold  
12 regression models.

13 Generalized threshold regression models are employed in a wide range of different fields  
14 of application. Hansen (2011) provides an overview of the extensive use of generalized  
15 threshold regression models in economic applications including e.g. models of output  
16 growth, forecasting, and the term structure of interest rates or stock returns. Samia et al.  
17 (2007) employ a generalized threshold regression model to analyze plague outbreaks, and  
18 Lee et al. (2011) complement these applications with examples in finance, sociology, and  
19 biostatistics among others.

20 Obviously, a good threshold estimator is crucial for the entire threshold regression model  
21 estimation. In this paper we discuss settings in which threshold identification becomes  
22 difficult. Typically, threshold parameters are estimated by the maximization of the cor-  
23 responding profile likelihood using a grid search, as the likelihood function is not differ-  
24 entiable with respect to the threshold parameter. This estimation procedure itself has an  
25 intrinsic problem: the profile likelihood is not defined for thresholds that leave fewer ob-  
26 servations in one of the regimes than are necessary to estimate the regression coefficients.  
27 Hence, in practice it is unavoidable to restrict the domain of the threshold parameters  
28 depending on the dimension of the regression coefficients. The literature offers arbitrary  
29 constraints including one observation per dimension of the regression coefficient (Samia  
30 and Chan, 2011) or 15% of the observations (Andrews, 1993) to give just two examples.  
31 This restriction can be problematic in small samples, especially if the true threshold is

32 close to the boundary of its domain.

33 Another problem occurs, if the threshold parameter itself lacks identifiability. In particu-  
34 lar, if differences between regimes are small and/or the regression coefficients' estimators  
35 are highly variable, the uncertainty of the threshold estimator increases. Note that the  
36 large variance of the regression coefficients' estimator is likely to be found in small sam-  
37 ples, for the true threshold at the boundary of its domain and also if the signal-to-noise  
38 ratio is low. We are not aware of any work that points out these deficiencies of the  
39 common threshold estimator even though the problematic settings frequently occur in  
40 empirical applications. Macro-economic data are often only available for a small sam-  
41 ple, e.g. if observations correspond to different countries. Spatial arbitrage modeling is  
42 another example (Greb et al., 2013).

43 Bayesian methods are also popular to estimate threshold regression models. In the litera-  
44 ture Bayesian estimation is typically based on non-informative priors, leading to what we  
45 refer to as the non-informative Bayesian estimator. For the threshold estimator in case of  
46 a threshold regression model with piecewise linear mean, Yu (2012) shows that, regardless  
47 of the choice of priors, Bayesian threshold estimators are asymptotically efficient among  
48 all estimators in the locally asymptotically minimax sense. However, in the critical small  
49 sample settings described above, the non-informative Bayesian estimator shares all the  
50 drawbacks of the standard likelihood estimator and can completely fail in certain cases,  
51 as we discuss in Section 3.2.

52 In this article, we suggest an alternative estimator, which we call the regularized Bayesian  
53 estimator. Contrary to previous work on estimation in threshold regression (Samia and  
54 Chan, 2011; Yu, 2012), we focus on the estimator's performance in critical small sample  
55 situations. Simulations confirm that it yields good results even in settings in which likeli-  
56 hood and non-informative Bayesian estimator are highly susceptible to faults. Given the  
57 threshold parameter's crucial function within the model, our idea is to improve estimation

58 of the whole model by improving estimation of this essential parameter.

59 To summarize the intuition for the new threshold estimator: If regression coefficients  
60 were known, none of the problems in threshold estimation outlined above would exist.  
61 This suggests that stabilizing their estimates might help to prevent them from distorting  
62 the threshold estimates. In addition, regularization of regression coefficient estimates  
63 allows us to obtain a posterior density that is well-defined on the entire domain of the  
64 threshold parameters. We achieve regularization by a particular specification of priors.  
65 While it proves to be beneficial in the critical small sample situations, the choice of priors  
66 does not have an impact asymptotically (as Yu, 2012, shows for a threshold regression  
67 model with piecewise linear mean and independent observations). We further derive  
68 an explicit (approximate) expression of the posterior density, which allows us to utilize  
69 existing functions for mixed models in standard software to easily compute the threshold  
70 estimator and simultaneously obtain estimates for the remaining model parameters.

71 The rest of this article is organized as follows. We specify the generalized threshold  
72 regression model in the second section. In the third section, we review existing estimators  
73 for threshold regression models and point out their deficiencies. Here, we concentrate  
74 on estimators for the crucial threshold parameter. The regularized Bayesian estimator is  
75 introduced in the fourth section. In the fifth section, we derive conditions under which the  
76 regularized Bayesian estimates fare better than their likelihood counterparts. Simulation  
77 results are presented in the sixth section. We use the last section to discuss two empirical  
78 applications. The appendix contains some technical details.

## 79 **2 Model**

80 Observations  $(y_i, \mathbf{X}_i^T, q_i) \in \mathbb{R} \times \mathbb{R}^p \times \mathbb{R}$ ,  $i = 1, \dots, n$ , are assumed to be realizations  
81 of random variables that follow a generalized threshold regression model with threshold  
82 parameter  $\psi \in \mathbb{R}$ , regression coefficients  $\beta_1, \beta_2 \in \mathbb{R}^p$  and scale (or dispersion) parameter

83  $\phi \in \mathbb{R}^+$ , that is

$$\mu_i = \mathbb{E}(y_i | \mathbf{X}_i^T, q_i) = h(\eta_i) \quad (1)$$

84 where  $h$  is a known one-to-one function, the inverse of the link function  $g = h^{-1}$ , and

$$\eta_i = I(q_i \leq \psi) \mathbf{X}_i^T \boldsymbol{\beta}_1 + I(q_i > \psi) \mathbf{X}_i^T \boldsymbol{\beta}_2, \quad (2)$$

85 with  $I(\cdot)$  as the indicator function. Moreover, conditional on the design vector  $\mathbf{X}_i^T$  and  
 86 the transition variable  $q_i$ , the response variables  $y_i$  are independently drawn from an  
 87 exponential family distribution with density

$$f(y_i | \psi, \phi, \boldsymbol{\beta}_1, \boldsymbol{\beta}_2) = \exp \left\{ \frac{y_i \theta_i - b(\theta_i)}{\phi} + c(y_i, \phi) \right\}, \quad (3)$$

88 characterized by known functions  $b$  and  $c$  together with the natural parameter  $\theta_i = \theta(\mu_i)$ .

89 Above and in the following, the same symbol denotes both a random variable and its  
 90 realization; the context should eliminate ambiguities. To use matrix notation, we define

91 vectors  $\boldsymbol{\mu}$ ,  $\boldsymbol{\eta}$ ,  $\mathbf{y}$ ,  $\mathbf{q}$ ,  $\mathbf{I}(\mathbf{q} \leq \psi)$  and  $\mathbf{I}(\mathbf{q} > \psi)$  by stacking  $\mu_i$ ,  $\eta_i$ ,  $y_i$ ,  $q_i$ ,  $I(q_i \leq \psi)$  and

92  $I(q_i > \psi)$ , respectively, and create an  $n \times p$  matrix  $\mathbf{X}$  with rows  $\mathbf{X}_i^T$ ,  $i = 1, \dots, n$ . With

93  $\text{diag}\{\mathbf{I}(\cdot)\}$  the diagonal matrix with entries  $\mathbf{I}(\cdot)$  along the diagonal and  $\boldsymbol{\beta} = (\boldsymbol{\beta}_1^T, \boldsymbol{\beta}_2^T)^T$ ,

94 we can write

$$\boldsymbol{\eta} = \text{diag}\{\mathbf{I}(\mathbf{q} \leq \psi)\} \mathbf{X} \boldsymbol{\beta}_1 + \text{diag}\{\mathbf{I}(\mathbf{q} > \psi)\} \mathbf{X} \boldsymbol{\beta}_2 = \mathbf{X}_1 \boldsymbol{\beta}_1 + \mathbf{X}_2 \boldsymbol{\beta}_2 = \mathbf{X}_\psi \boldsymbol{\beta}.$$

95 We consider generalized threshold regression models with one threshold to keep the expo-  
 96 sition simple; extension to generalized threshold regression models with more thresholds  
 97 is straightforward (see e.g. Greb et al., 2013).

98 Naturally, our model covers  $y_i = I(q_i \leq \psi) \mathbf{X}_i^T \boldsymbol{\beta}_1 + I(q_i > \psi) \mathbf{X}_i^T \boldsymbol{\beta}_2 + \varepsilon_i$ ,  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$

99 and  $i = 1, \dots, n$ . This is by far the most frequently encountered generalized threshold

100 regression model in the literature. It is broad enough to comprise the popular threshold

101 autoregressive model in which the transition variable  $q_i$  is an element of  $\mathbf{X}_i$  (Tong and

102 Lim, 1980; Tong, 2011, for a review of the development of the model).

103 Depending on the assumptions on the data generating process, inferences (or estimators)  
104 for model (1) – (3) can take on different asymptotic behavior. A first differentiation re-  
105 gards the transition variable  $q_i$ . Change point models are characterized by deterministic  
106  $q_i = i$ , while for threshold models  $q_i$  is a random variable which follows any continu-  
107 ous distribution. This is reflected in distinct limit likelihood ratio processes and, hence,  
108 asymptotic behavior of the maximum likelihood estimators for  $\psi$  in the two models. The  
109 limiting likelihood ratio process involves a functional of random walks for change point  
110 models and of compound Poisson processes for threshold models. Check Bai (1997) for  
111 more details on the asymptotic properties in the former, and Samia and Chan (2011) for  
112 the limiting behavior of the profile log-likelihood and the asymptotic distribution of the  
113 profile likelihood threshold estimator in the latter case.

114 If the transition variable coincides with one of the covariates and the regression function is  
115 continuous at the threshold, least squares estimates are known to be normally distributed  
116 (for threshold models, see Chan and Tsay, 1998; Feder, 1975, treats change-point models),  
117 which simplifies inference. Clearly, once the data is sampled, the estimation procedure in  
118 both change point and threshold models is the same. Referring to a threshold regression  
119 model with piecewise linear mean, Hansen (2000) points out that “if the observed values  
120 of  $q_i$  are distinct, the parameters can be estimated by sorting the data based on  $q_i$ , and  
121 then applying known methods for change point problems”.

122 As the focus of this article is on estimation problems that arise in small samples, we do  
123 not further differentiate between models. In the real-data examples, we concentrate on  
124 discontinuous threshold models since they are frequently encountered in applications and  
125 have not been studied as extensively as change point models due to their more intricate  
126 limiting behavior.

### 127 3 Estimation of threshold regression models

#### 128 3.1 The likelihood estimator

129 As noted in the introduction, the prevalent estimator of threshold regression models is  
 130 the likelihood estimator, see e.g. Samia and Chan (2011) or Hansen (2000). Thereby, the  
 131 threshold parameter is estimated from the corresponding profile likelihood  $\mathcal{L}_p$ , which is  
 132 constructed from the likelihood function  $\mathcal{L}$ , by replacing nuisance parameters  $\boldsymbol{\beta}^T \in \mathbb{R}^{2p}$   
 133 and  $\phi \in \mathbb{R}$  with their maximum likelihood estimates at given values of  $\psi$  (which are  
 134 just standard (weighted) least squares estimators). More specifically, we work with the  
 135 conditional profile likelihood function given  $\mathbf{X}$  and  $\mathbf{q}$ ,

$$\mathcal{L}_p(\psi) = \prod_{i=1}^n f(y_i | \psi, \hat{\phi}_\psi, \hat{\boldsymbol{\beta}}_\psi) = \exp \left[ \sum_{i=1}^n \left\{ \frac{y_i \hat{\theta}_i - b(\hat{\theta}_i)}{\hat{\phi}_\psi} + c(y_i, \hat{\phi}_\psi) \right\} \right],$$

136 where  $\hat{\theta}_i = \theta \circ h(\hat{\eta}_i) = \theta \circ h \left\{ I(q_i \leq \psi) \mathbf{X}_i^T \hat{\boldsymbol{\beta}}_{1_\psi} + I(q_i > \psi) \mathbf{X}_i^T \hat{\boldsymbol{\beta}}_{2_\psi} \right\}$  and  $\hat{\boldsymbol{\beta}}_\psi$  and  $\hat{\phi}_\psi$  are  
 137 maximum likelihood estimators at a fixed  $\psi$ . In the following, we assume a canonical link,  
 138 that is,  $\theta_i = \eta_i$ . All developments still hold approximately if this assumption does not  
 139 hold. We denote the profile log-likelihood with  $\ell_p(\psi) = \log \mathcal{L}_p(\psi)$ .

140 In generalized threshold regression models, the domain of the threshold parameter  $\psi$   
 141 is restricted to a random set  $\Psi = \{\psi \in \mathbb{R} | q_{(1)} \leq \psi \leq q_{(n)}\} \subseteq \mathbb{R}$ , where  $q_{(i)}$  denotes the  
 142  $i$ th order statistic. To measure the proximity of a threshold  $\psi$  to the boundary of its  
 143 domain  $\Psi$ , we introduce  $d(\psi) = \min(j, n - j)/p$  with  $j$  such that  $q_{(j)} \leq \psi < q_{(j+1)}$ . The  
 144 quantity  $d(\psi)$  is the distance between  $\psi$  and  $\Psi$ 's boundary in terms of the number of  
 145 observations between them relative to the dimension of the regression coefficients,  $p =$   
 146  $\dim(\boldsymbol{\beta}_k)$ ,  $k = 1, 2$ . When  $d(\psi) = 1$ ,  $\psi$  assigns at least  $p$  observations to each of the  
 147 regimes. The allocation of 5% of the observations into one of the regimes can be expressed  
 148 as  $d(\psi) = 0.05 n/p$ .

149 Clearly,  $\mathcal{L}_p(\psi)$  is not defined for  $d(\psi) < 1$ , since in this case  $\psi$  does not leave enough  
 150 observations for the estimation of  $\boldsymbol{\beta}_k$  in one of the regimes. Hence, in practice it is

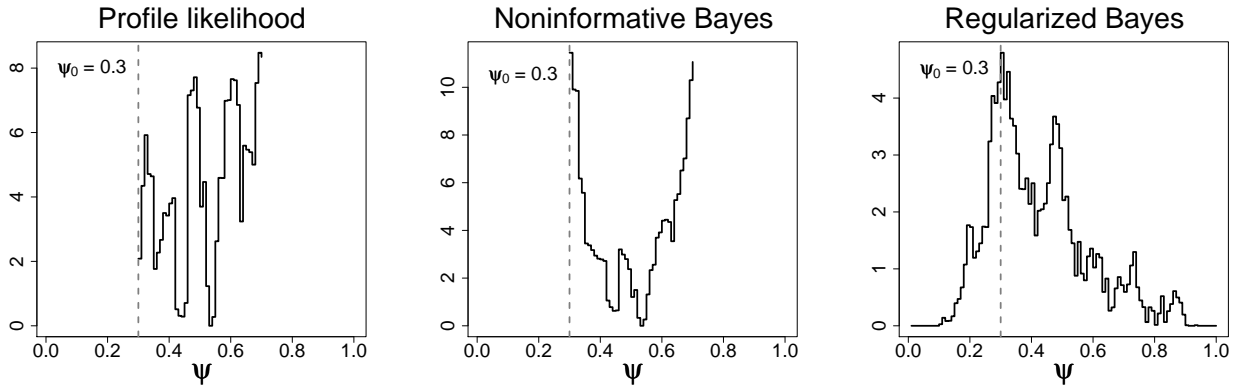


Figure 1: For a sample run corresponding to setting 1C of Section 6,  $\ell_p(\psi)$  is shown on the left,  $\log p_{nB}(\psi|\mathbf{y}, \mathbf{X}, \mathbf{q})$  in the middle and  $\log p_{rB}(\psi|\mathbf{y}, \mathbf{X}, \mathbf{q})$  on the right.

151 inevitable to restrict  $\Psi$  to  $\Psi^*(c) = \{\Psi | d(\psi) > c\}$  for some  $c \geq 1$ . In the literature different  
 152 heuristic suggestions for the choice of  $c$  have been proposed. For example, Hansen and  
 153 Seo (2002) propose  $c = 0.05 n/p$ , we find  $c = 0.15 n/p$  in Andrews (1993) and Samia and  
 154 Chan (2011) even use  $c = 0.25 n/p$  for their application.

155 The profile likelihood threshold estimator is then given by

$$\hat{\psi}_{pL} = \operatorname{argmax}_{\psi \in \Psi^*(c)} \mathcal{L}_p(\psi).$$

156 This definition based on the restricted domain  $\Psi^*(c)$  immediately suggests that in settings  
 157 in which  $d(\psi_0) < c$  for a true threshold  $\psi_0$ ,  $\hat{\psi}_{pL}$  is inconsistent. The left panel of Fig. 1  
 158 illustrates this showing the profile log-likelihood for a sample run of a generalized threshold  
 159 regression model corresponding to the simulation setting 1C detailed in Section 6. If  
 160  $\Psi^*(1) = [0.3, 0.7]$  would be restricted any further, e.g. to be  $[0.31, 0.69]$ , then the true  
 161 threshold  $\psi_0 = 0.3$  would be excluded from the threshold domain and  $\hat{\psi}_{pL}$  would move to  
 162 the next extremum. For small  $n$ , large  $p$  and  $\psi_0$  close to the boundary of  $\Psi$ ,  $d(\psi_0) < c$   
 163 is likely to be the case. Altogether, subjective restriction of the threshold domain is an  
 164 undesirable property of threshold estimation based on the profile likelihood.

165 The same plot in Fig. 1 also exemplifies that in certain small-sample settings the pro-  
 166 file (log-)likelihood can be jagged and have multiple extrema, leading to an estimated



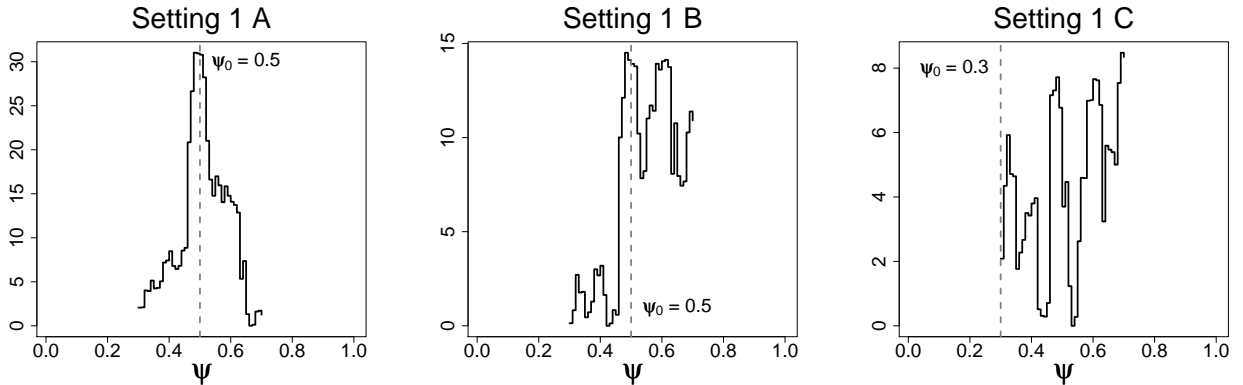


Figure 2: Sample (log) profile likelihood functions  $\ell_p(\psi)$  for different settings.

167 threshold that is very sensitive to the initialization of the search. Large variance of  $\hat{\beta}_\psi$   
168 and/or small differences between regimes compared to the noise level can have a strong  
169 distorting effect on the profile (log-)likelihood and are associated with settings charac-  
170 terized by small  $n$  relative to  $p$ , but can also be due to low signal-to-noise ratio, model  
171 misspecifications (e.g. overdispersion), or a threshold that is close to the boundary of  
172 its domain. This is exposed in the left as compared with the middle plot of Fig. 2; the  
173 log-likelihoods depicted in these plots belong to models which only differ in one aspect: in  
174 the plot on the left-hand side, the residual standard deviation is 0.75, while in the middle  
175 plot it is 1.5, increasing the signal-to-noise ratio and  $\text{var}(\hat{\beta}_\psi)$ . Clearly, the log-likelihood  
176 in the middle plot is highly distorted over the whole range of  $\Psi$ , triggering multiple ex-  
177 tremas and a highly variable estimator for  $\psi$ . Moving the true threshold closer to the  
178 boundary, as shown in the right plot of Fig. 2, leads to an even stronger deformation of  
179 the log-likelihood.

180 In summary, in small samples and particular settings exemplified above, the profile like-  
181 lihood threshold estimator can perform poorly, being very sensitive to inappropriate esti-  
182 mates of the nuisance parameters and relying on a subjective restriction of its domain.

## 183 3.2 The Bayesian estimator

184 For threshold regression models with piecewise linear mean, there is a long tradition of  
185 using Bayesian techniques in applied work beginning with Bacon and Watts (1971) and  
186 including Geweke and Terui (1993) among many others. This popularity can be at least  
187 partially attributed to practical advantages, since the Bayesian approach offers a natural  
188 framework for inference and accounts for the uncertainty of the nuisance parameters. The  
189 Bayesian regression coefficients estimators coincide with the maximum likelihood ones for  
190 non-informative priors. The theoretical properties of Bayesian threshold estimators in  
191 certain generalized threshold regression models have been investigated by Yu (2012). He  
192 shows that for independently and identically distributed observations Bayesian threshold  
193 estimators are asymptotically efficient among all estimators in the locally asymptotically  
194 minimax sense and strictly more efficient than the maximum likelihood estimator. In a  
195 related paper, Chan and Kutoyants (2012) examine asymptotic properties of Bayesian  
196 estimators in threshold autoregression models. They note that in the limit, the variance  
197 of the Bayesian estimator is smaller than that of the maximum likelihood estimator.

198 Without any prior knowledge of possible parameter values, it is natural to assume a  
199 uniform prior for the threshold parameter and non-informative priors for the regression  
200 coefficients; these choices are (almost) omnipresent in the Bayesian literature on gener-  
201 alized threshold regression models with piecewise linear mean. While the priors do not  
202 have an impact asymptotically, it turns out that they do affect the performance of the  
203 Bayesian threshold estimator in finite samples. We show that non-informative priors can  
204 distort estimates, especially in small samples.

205 It is straightforward to obtain an approximation of a generalized threshold regression  
206 model's posterior density  $p_{n_B}(\psi|\phi, \mathbf{y}, \mathbf{X}, \mathbf{q})$  associated with non-informative (improper)  
207 priors  $p(\boldsymbol{\beta}) \propto 1$  and  $p(\psi|\mathbf{q}) \propto I(\psi \in \Psi)$  based on a Laplace approximation (Shun and

208 McCullagh, 1995; Severini, 2000) of the integral for fixed  $p \ll n$

$$\int_{\mathbb{R}^{2p}} p(y|\psi, \phi, \boldsymbol{\beta}, \mathbf{X}, \mathbf{q}) d\boldsymbol{\beta} = \mathcal{L}_p(\psi)(2\pi)^p \left| -\frac{\partial^2 \ell}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^T}(\psi, \phi, \hat{\boldsymbol{\beta}}_\psi) \right|^{-1/2} + \mathcal{O}(n^{-1}),$$

209 with  $\ell(\psi, \phi, \boldsymbol{\beta}) = \log \mathcal{L}(\psi, \phi, \boldsymbol{\beta})$ . As  $\left| -\partial^2 \ell / \partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^T(\psi, \phi, \hat{\boldsymbol{\beta}}_\psi) \right| = |\mathbf{X}_\psi^T \mathbf{W} \mathbf{X}_\psi|$ , we get

$$p_{nB}(\psi|\phi, \mathbf{y}, \mathbf{X}, \mathbf{q}) = \mathcal{L}_p(\psi)(2\pi)^p |\mathbf{X}_\psi^T \mathbf{W} \mathbf{X}_\psi|^{-1/2} I(\psi \in \Psi) / p(\mathbf{y}) + \mathcal{O}(n^{-1}).$$

210 With this, the prevalent Bayesian threshold estimator in the literature is the posterior  
 211 mean  $\hat{\psi}_{nB} = \int_{\Psi^*} \psi p_{nB}(\psi|\phi, \mathbf{y}, \mathbf{X}, \mathbf{q}) d\psi$ . Comparing  $p_{nB}(\psi|\phi, \mathbf{y}, \mathbf{X}, \mathbf{q})$  with  $\mathcal{L}_p(\psi)$ , we note  
 212 that they differ by a term proportional to  $|\mathbf{X}_\psi^T \mathbf{W} \mathbf{X}_\psi|^{-1/2}$ . In the case of Gaussian  
 213 observations,  $\mathbf{W} = \mathbf{I}_n / \sigma^2$ . Since  $|\mathbf{X}_\psi^T \mathbf{W} \mathbf{X}_\psi| = |\mathbf{X}_1^T \mathbf{W} \mathbf{X}_1| \cdot |\mathbf{X}_2^T \mathbf{W} \mathbf{X}_2| \rightarrow 0$  for  $d(\psi) \rightarrow$   
 214  $0$ ,  $p_{nB}(\psi|\phi, \mathbf{y}, \mathbf{X}, \mathbf{q})$  becomes very large for  $\psi$  close to the boundary of  $\Psi$ . Moreover, as the  
 215 profile likelihood function requires  $d(\psi) \geq 1$  to be well-defined, so does the calculation of  
 216 the posterior density. Again, the only solution in the literature is to restrict the parameter  
 217 space  $\Psi$  (which in our Bayesian framework is equivalent to working with a uniform prior  
 218  $\psi \sim U[\Psi^*]$  instead of  $\psi \sim U[\Psi]$ ). In this case, however,  $p_{nB}(\psi|\phi, \mathbf{y}, \mathbf{X}, \mathbf{q})$  becomes largest  
 219 exactly for values of  $\psi$  which are arbitrarily included or excluded from  $\Psi^*$  by varying  $c$ .  
 220 Consequently, expanding or reducing  $\Psi^*$  critically affects the Bayesian threshold estimate,  
 221 whether it is calculated as the posterior mode, mean or median. The middle plot in Fig. 1  
 222 illustrates this problem.

## 223 4 The regularized Bayesian estimator

224 When rethinking the threshold regression estimation, there are good arguments for con-  
 225 tinuing to pursue Bayesian options. In general, Bayesian estimators naturally incorporate  
 226 the uncertainty of nuisance parameters and there are reasons to expect the threshold  
 227 estimators to be (at least asymptotically) the most efficient estimators, as discussed in  
 228 Section 3.2.

229 Our idea now is to exploit understanding of when reliable estimation becomes particularly  
 230 difficult in order to regularize the posterior density. First, we define

$$\boldsymbol{\eta} = \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\beta}_2 = (\mathbf{X}_1 + \mathbf{X}_2)\boldsymbol{\beta}_1 + \mathbf{X}_2(\boldsymbol{\beta}_2 - \boldsymbol{\beta}_1) = \mathbf{X}\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\delta}. \quad (4)$$

231 Here,  $\mathbf{X}$  is independent of  $\psi$ , while  $\mathbf{X}_2 = \mathbf{X}_2(\psi) = \text{diag}\{\mathbf{I}(\mathbf{q} > \psi)\}\mathbf{X}$ . Hence, if  $\boldsymbol{\delta}$  is  
 232 small and/or its estimators are highly variable, it becomes hard to identify the threshold  
 233  $\psi$ . We, therefore, suggest to regularize the estimator for  $\boldsymbol{\delta}$ . In a Bayesian framework the  
 234 natural approach is to assume  $\boldsymbol{\delta} \sim \mathcal{N}(0, \sigma_\delta^2 \mathbf{I}_p)$ . When  $\sigma_\delta^2$  tends towards infinity, this prior  
 235 becomes non-informative. However, for small values  $\sigma_\delta^2$ , we introduce prior knowledge  
 236 suggesting that  $\boldsymbol{\delta}$  takes values close to zero, that is there is no threshold in the model.  
 237 The most important characteristic of this new choice of priors is that it regularizes the  
 238 posterior density for  $\psi$  close to the boundary of  $\Psi$ . Putting priors on  $\sigma_\delta^2$  (e.g. an inverse  
 239 Gamma distribution) and  $\psi$  specifies a fully Bayesian model and allows for estimation  
 240 with Markov chain Monte Carlo techniques.

241 Alternatively, we suggest to use a Laplace approximation to get the approximate pos-  
 242 terior  $p(\psi|\phi, \sigma_\delta^2, \mathbf{y}, \mathbf{X}, \mathbf{q})$ . This accelerates estimation and enables us to illustrate the  
 243 regularizing effect. To evaluate the posterior density

$$p(\psi|\phi, \sigma_\delta^2, \mathbf{y}, \mathbf{X}, \mathbf{q}) = \frac{p(\psi|\mathbf{q})}{p(\mathbf{y}|\phi, \sigma_\delta^2, \mathbf{X}, \mathbf{q})} \int_{\mathbb{R}^p} \int_{\mathbb{R}^p} p(\mathbf{y}|\boldsymbol{\beta}_1, \boldsymbol{\delta}, \psi, \phi, \sigma_\delta^2, \mathbf{X}, \mathbf{q}) p(\boldsymbol{\delta}|\sigma_\delta^2) d\boldsymbol{\delta} d\boldsymbol{\beta}_1,$$

244 we use a Laplace approximation and follow a line of reasoning closely resembling Breslow  
 245 and Clayton (1993) to obtain

$$\begin{aligned} & \int_{\mathbb{R}^p} \int_{\mathbb{R}^p} p(\mathbf{y}|\boldsymbol{\beta}_1, \boldsymbol{\delta}, \psi, \phi, \sigma_\delta^2, \mathbf{X}, \mathbf{q}) p(\boldsymbol{\delta}|\sigma_\delta^2) d\boldsymbol{\delta} d\boldsymbol{\beta}_1 \\ &= (2\pi)^{p/2} \exp \left\{ -\frac{1}{2} (\tilde{\mathbf{z}} - \mathbf{X}\hat{\boldsymbol{\beta}}_1)^T \mathbf{V}^{-1} (\tilde{\mathbf{z}} - \mathbf{X}\hat{\boldsymbol{\beta}}_1) + \sum_{i=1}^n c(y_i, \phi) \right\} \\ & \quad \cdot |\sigma_\delta^2 \mathbf{X}_2^T \mathbf{W} \mathbf{X}_2 + \mathbf{I}_p|^{-1/2} |\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X}|^{-1/2} + \mathcal{O}(n^{-1}), \end{aligned} \quad (5)$$

246 with the working variable  $\tilde{\mathbf{z}}$  defined as  $\tilde{\mathbf{z}} = \mathbf{X}\hat{\boldsymbol{\beta}}_1 + \mathbf{X}_2\hat{\boldsymbol{\delta}} + \mathbf{G}(\mathbf{y} - \boldsymbol{\mu})$ ,  
 247  $\mathbf{G} = \text{diag}\{g'(\mu_i)\}$ , and  $\mathbf{V} = \mathbf{W}^{-1} + \sigma_\delta^2 \mathbf{X}_2 \mathbf{X}_2^T$  for  $\mathbf{W}^{-1} = \text{diag}\{\phi b''(\theta_i) g'(\mu_i)^2\}$ .

248 Here,  $\boldsymbol{\mu}$ ,  $\mathbf{G}$ ,  $\mathbf{W}$  and  $\mathbf{V}$  are evaluated at the (approximate) pos-  
 249 terior mode  $(\hat{\boldsymbol{\beta}}_1, \hat{\boldsymbol{\delta}}) = \arg \max_{(\boldsymbol{\beta}_1, \boldsymbol{\delta}) \in \mathbb{R}^{2p}} p(\boldsymbol{\beta}_1, \boldsymbol{\delta} | \psi, \phi, \sigma_\delta^2, \mathbf{y}, \mathbf{X}, \mathbf{q})$ , that is,  
 250  $\hat{\boldsymbol{\beta}}_1 = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \tilde{\mathbf{z}}$  and  $\hat{\boldsymbol{\delta}} = \sigma_\delta^2 \mathbf{X}_2^T \mathbf{V}^{-1} (\tilde{\mathbf{z}} - \mathbf{X} \hat{\boldsymbol{\beta}}_1)$ . Note that these re-  
 251 gression parameter estimators are regularized and are different from usual likelihood  
 252 estimators. Details on the derivation of (5) are provided in the appendix.

253 In contrast to the posterior based on non-informative priors, the term  $|\mathbf{X}_\psi^T \mathbf{W} \mathbf{X}_\psi|$   
 254 disappears, and with it the deteriorations near the boundary of  $\Psi$  observed for  
 255  $p_{nB}(\psi | \phi, \mathbf{y}, \mathbf{X}, \mathbf{q})$ . Moreover,  $p(\psi | \phi, \sigma_\delta^2, \mathbf{y}, \mathbf{X}, \mathbf{q})$  is well-defined for all  $\psi \in \Psi$ , independent  
 256 of  $d(\psi)$ . It is easy to see that  $\hat{\boldsymbol{\delta}} \rightarrow 0$  and  $\hat{\boldsymbol{\beta}}_1 \rightarrow (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \tilde{\mathbf{z}}$  at the boundary of  
 257  $\Psi$ , for  $\mathbf{X}_2 = 0$  or  $\mathbf{X}_2 = \mathbf{X}$ . We do not encounter the ill-posed problem of estimating  $p$   
 258 nuisance parameters from  $m < p$  observations, or calculating  $\hat{\boldsymbol{\beta}}_\psi$  when  $d(\psi) < 1$ , as in  
 259 profile likelihood or non-informative Bayesian estimation. Consequently, there is no need  
 260 to subjectively restrict the parameter space.

261 Considering

$$\hat{\boldsymbol{\delta}} = \sigma_\delta^2 \mathbf{X}_2^T \mathbf{V}^{-1} (\tilde{\mathbf{z}} - \mathbf{X} \hat{\boldsymbol{\beta}}_1) = \arg \min_{\boldsymbol{\delta} \in \mathbb{R}^p} (\tilde{\mathbf{z}} - \mathbf{X} \hat{\boldsymbol{\beta}}_1 - \mathbf{X}_2 \boldsymbol{\delta})^T \mathbf{W} (\tilde{\mathbf{z}} - \mathbf{X} \hat{\boldsymbol{\beta}}_1 - \mathbf{X}_2 \boldsymbol{\delta}) + \frac{1}{\sigma_\delta^2} \boldsymbol{\delta}^T \boldsymbol{\delta}, \quad (6)$$

262 it becomes evident that the proposed prior leads to the strategy of turning an ill-posed  
 263 into a well-posed problem tracing back to Tikhonov et al. (1977). For small values of the  
 264 regularization parameter  $1/\sigma_\delta^2$ , the first term of the functional to be minimized in (6) will  
 265 drive the resulting  $\hat{\boldsymbol{\delta}}$ , for large values it is the latter. For the nuisance parameter estimates  
 266  $\hat{\boldsymbol{\beta}}_1$  and  $\hat{\boldsymbol{\beta}}_2 = \hat{\boldsymbol{\beta}}_1 + \hat{\boldsymbol{\delta}}$ , basic matrix algebra reveals that  $\hat{\boldsymbol{\beta}}_1 \rightarrow (\mathbf{X}_1^T \mathbf{W} \mathbf{X}_1)^{-1} \mathbf{X}_1^T \mathbf{W} \tilde{\mathbf{z}}$  and  
 267  $\hat{\boldsymbol{\beta}}_2 \rightarrow (\mathbf{X}_2^T \mathbf{W} \mathbf{X}_2)^{-1} \mathbf{X}_2^T \mathbf{W} \tilde{\mathbf{z}}$  for  $\sigma_\delta^2 \rightarrow \infty$ , while for  $\sigma_\delta^2 \rightarrow 0$ , both  $\hat{\boldsymbol{\beta}}_1$  and  $\hat{\boldsymbol{\beta}}_2$  converge to  
 268  $(\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \tilde{\mathbf{z}}$ .

269 Clearly, the choice of the regularization parameter  $\sigma_\delta^2$  is essential to any estimate based on  
 270  $p(\psi | \phi, \sigma_\delta^2, \mathbf{y}, \mathbf{X}, \mathbf{q})$ . It can naturally be estimated in the fully Bayesian framework. How-  
 271 ever, pursuing our approximate approach further we prefer to make use of the empirical

272 Bayes paradigm. In general, the empirical Bayes approach to modeling observations  $\mathbf{y}$   
 273 differs from the usual Bayesian setup in that the hyperparameters for the highest level in  
 274 the model's hierarchy are replaced by their maximum likelihood estimates. In our case,  
 275 we obtain  $\hat{\sigma}_\delta^2$  for fixed  $\mathbf{X}$ ,  $\mathbf{q}$  and  $\psi$  by maximizing

$$p(\mathbf{y}|\psi, \phi, \sigma_\delta^2, \mathbf{X}, \mathbf{q}) = \int_{\mathbb{R}^p} \int_{\mathbb{R}^p} p(\mathbf{y}|\boldsymbol{\beta}_1, \boldsymbol{\delta}, \psi, \phi, \sigma_\delta^2, \mathbf{X}, \mathbf{q}) p(\boldsymbol{\delta}|\sigma_\delta^2) d\boldsymbol{\delta} d\boldsymbol{\beta}_1,$$

276 so as to base threshold estimation on

$$p_{rB}(\psi|\mathbf{y}, \mathbf{X}, \mathbf{q}) = p(\psi|\mathbf{y}, \mathbf{X}, \mathbf{q}, \hat{\phi}_\psi, \hat{\sigma}_\delta^2) \propto \left| \hat{\sigma}_\delta^2 \mathbf{X}_2^T \mathbf{W} \mathbf{X}_2 + \mathbf{I}_p \right|^{-1/2} \left| \mathbf{X}^T \hat{\mathbf{V}}^{-1} \mathbf{X} \right|^{-1/2} \\ \cdot \exp \left\{ -\frac{1}{2} (\tilde{\mathbf{z}} - \mathbf{X} \hat{\boldsymbol{\beta}}_1)^T \hat{\mathbf{V}}^{-1} (\tilde{\mathbf{z}} - \mathbf{X} \hat{\boldsymbol{\beta}}_1) + \sum_{i=1}^n c(y_i, \hat{\phi}_\psi) \right\} I(\psi \in \Psi),$$

277 with  $\hat{\mathbf{V}}$  evaluated at  $\hat{\sigma}_\delta^2$ . The right plot in Fig. 1 shows the log of this posterior density for  
 278 a sample run corresponding to simulation setting 1 C of Section 6. It is clearly well-defined  
 279 over the whole domain of the threshold and its values are regularized at the boundary  
 280 regions, making the extremum more pronounced.

281 Once the posterior density is obtained, one can calculate  $\hat{\psi}_{rB}$ . We observed that in critical  
 282 small-sample settings the posterior density is often characterized by multiple modes. Thus,  
 283 obtaining an estimate based on numerical maximization (the posterior mode) is likely to be  
 284 challenging. The posterior mean presents a more robust alternative. However, when the  
 285 true threshold is located close to the boundary of  $\Psi$ , the posterior distribution is skewed  
 286 towards this boundary. As a result, the posterior mean tends to be drawn towards the  
 287 middle of  $\Psi$  (Doodson, 1917; Kendall, 1943, page 35). Hence, we opt for the posterior  
 288 median as a compromise between the latter two. Accordingly, we suggest calculating a  
 289 regularized Bayesian threshold estimator  $\hat{\psi}_{rB}$  as

$$\int_{q(1)}^{\hat{\psi}_{rB}} p_{rB}(\psi|\mathbf{y}, \mathbf{X}, \mathbf{q}, \phi) d\psi = 0.5$$

290 assuming a prior  $p(\psi|\mathbf{q}) \propto I(\psi \in \Psi)$  for  $\psi$ .

291 By definition, the restricted (or residual) likelihood function (Harville, 1977) of a gener-  
 292 alized linear mixed model is the approximate posterior (5). Hence, the function `g1mmPQL`  
 293 in the R-package `MASS` readily provides us with the desired estimate  $\hat{\sigma}_\delta^2$ . Moreover, the  
 294 function simultaneously produces an estimate  $\hat{\phi}_\psi$ . For the Gaussian case, we can employ  
 295 the function `lme` directly (with its parameter `method` left at the default value `REML`). It  
 296 is part of the R-package `nlme`. This possibility to take advantage of existing functions  
 297 implemented for mixed models greatly facilitates computation of our proposed estimator,  
 298 which can be performed in seconds.

299 Inference about all of the model parameters naturally follows in this Bayesian framework.  
 300 In particular, confidence regions for  $\psi$  are formed as credible sets; an equi-tailed credible  
 301 set  $C$  of level  $1 - 2\alpha$  is defined as

$$C = \int_{q_p(\alpha)}^{q_p(1-\alpha)} p(\psi|\mathbf{y}, \mathbf{X}, \mathbf{q}, \phi) d\psi, \quad q_p(\alpha) = \inf_{x \in \Psi} \left\{ x \mid \int_{\psi \leq x} p(\psi|\mathbf{y}, \mathbf{X}, \mathbf{q}, \phi) d\psi \geq \alpha \right\}.$$

302 These credible sets are valid for change-point and threshold models, both continuous and  
 303 discontinuous. By contrast, in the frequentist framework it is straightforward to obtain  
 304 confidence intervals for continuous models. For discontinuous models the asymptotic  
 305 distribution does not readily provide a feasible way to construct confidence intervals as it  
 306 depends on (a possibly large number of) nuisance parameters.

## 307 **5 Comparison of regularized Bayesian and maximum** 308 **likelihood estimation**

309 Our new estimation procedure results in new regularized regression coefficients estima-  
 310 tors, whose properties have not been investigated so far. In the following, we compare  
 311 regularized Bayesian and maximum likelihood approaches to estimation of threshold re-  
 312 gression models in terms of mean squared error under the frequentist model. Thereby, we  
 313 treat the threshold as fixed and known, but allow for any, not necessarily true threshold  
 314  $\psi$ .

315 A natural measure for comparing coefficient estimates is the mean squared error  
316  $M(\mathbf{X}_\psi \hat{\boldsymbol{\beta}}) = E(\mathbf{X}_\psi \hat{\boldsymbol{\beta}} - \mathbf{X}_\psi \boldsymbol{\beta})^T (\mathbf{X}_\psi \hat{\boldsymbol{\beta}} - \mathbf{X}_\psi \boldsymbol{\beta})$ , where E denotes the conditional expect-  
317 tation without averaging over the prior assumptions, i.e. expectation with respect to  
318 the distribution of  $\mathbf{Y}$  given  $\boldsymbol{\delta}$ , which corresponds to the usual frequentist framework.  
319 In the context of ridge regression, this approach has been criticized for indiscriminately  
320 putting together the mean squared errors of the components (Nelder, 1972; Theobald,  
321 1974). As an alternative, Theobald (1974) suggested to consider a weighted sum  
322  $M_A(\mathbf{X}_\psi \hat{\boldsymbol{\beta}}) = E(\mathbf{X}_\psi \hat{\boldsymbol{\beta}} - \mathbf{X}_\psi \boldsymbol{\beta})^T \mathbf{A} (\mathbf{X}_\psi \hat{\boldsymbol{\beta}} - \mathbf{X}_\psi \boldsymbol{\beta})$  for a non-negative definite matrix  $\mathbf{A}$ .  
323 Here,  $\psi$  is an arbitrary, fixed threshold. Of course, a comparison between  $M(\mathbf{X}_\psi \hat{\boldsymbol{\beta}})$  (or  
324  $M_A(\mathbf{X}_\psi \hat{\boldsymbol{\beta}})$ ) for different  $\hat{\boldsymbol{\beta}}$  is both interesting for such general  $\psi$  as well as the true  
325 threshold  $\psi_0$ . With this in mind, we state the following result.

326 **Theorem 1** For maximum likelihood estimates  $\hat{\boldsymbol{\beta}}_{ML} = (\mathbf{X}_\psi^T \mathbf{W} \mathbf{X}_\psi)^{-1} \mathbf{X}_\psi^T \mathbf{W} \mathbf{z}$  and regu-  
327 larized Bayesian estimates  $\hat{\boldsymbol{\beta}}_{rB} = (\mathbf{X}_\psi^T \mathbf{W} \mathbf{X}_\psi + \mathbf{H})^{-1} \mathbf{X}_\psi^T \mathbf{W} \mathbf{z}$  of  $\boldsymbol{\beta}$  based on a threshold  
328  $\psi \leq \psi_0$ ,  $\psi_0$  the true threshold,

$$(i) \quad M_A(\mathbf{X}_\psi \hat{\boldsymbol{\beta}}_{ML}) - M_A(\mathbf{X}_\psi \hat{\boldsymbol{\beta}}_{rB}) \geq 0 \text{ for all non-negative definite matrices } \mathbf{A}$$

$$\Leftrightarrow \mathbf{D} \{ (\mathbf{I} + \mathbf{C}) \mathbf{H} - (\mathbf{B} + \mathbf{H}) \boldsymbol{\beta} \boldsymbol{\beta}^T (\mathbf{B}^T + \mathbf{H}) + \mathbf{C} \mathbf{B} \boldsymbol{\beta} \boldsymbol{\beta}^T \mathbf{B}^T \mathbf{C}^T \} \mathbf{D}^T$$

is non-negative definite.

$$(ii) \quad M(\mathbf{X}_\psi \hat{\boldsymbol{\beta}}_{ML}) - M(\mathbf{X}_\psi \hat{\boldsymbol{\beta}}_{rB}) \geq 0$$

$$\Leftrightarrow \text{tr} \{ \mathbf{H} \mathbf{D}^T \mathbf{D} (\mathbf{I} + \mathbf{C}) \} - \boldsymbol{\beta}^T \{ (\mathbf{B}^T + \mathbf{H}) \mathbf{D}^T \mathbf{D} (\mathbf{B} + \mathbf{H}) + \mathbf{B}^T \mathbf{D}_0^T \mathbf{D}_0 \mathbf{B} \} \boldsymbol{\beta} \geq 0.$$

329 Here,  $\mathbf{W}^{-1} = \text{diag} \{ \phi b''(\theta_i) g'(\mu_i)^2 \}$ ,  $\mathbf{G} = \text{diag} \{ g'(\mu_i) \}$ , and  $\mathbf{z} = \mathbf{X}_\psi \boldsymbol{\beta} + \mathbf{G}(\mathbf{y} - \boldsymbol{\mu})$   
330 as before,  $\mathbf{H} = 1/\sigma_\delta^2 \begin{pmatrix} \mathbf{I}_p & -\mathbf{I}_p \\ -\mathbf{I}_p & \mathbf{I}_p \end{pmatrix}$ ,  $\mathbf{B} = \begin{pmatrix} 0 & 0 \\ -\mathbf{X}_{[\psi, \psi_0]}^T \mathbf{W} \mathbf{X}_{[\psi, \psi_0]} & \mathbf{X}_{[\psi, \psi_0]}^T \mathbf{W} \mathbf{X}_{[\psi, \psi_0]} \end{pmatrix}$   
331 with  $\mathbf{X}_{[\psi, \psi_0]} = \text{diag} \{ \mathbf{I}(\psi < \mathbf{q} \leq \psi_0) \} \mathbf{X}$ ,  $\mathbf{C} = \mathbf{I} + \mathbf{H} (\mathbf{X}_\psi^T \mathbf{W} \mathbf{X}_\psi)^{-1}$ ,  
332  $\mathbf{D} = \mathbf{X}_\psi (\mathbf{X}_\psi^T \mathbf{W} \mathbf{X}_\psi + \mathbf{H})^{-1}$ , and  $\mathbf{D}_0 = \mathbf{X}_\psi (\mathbf{X}_\psi^T \mathbf{W} \mathbf{X}_\psi)^{-1}$ .

**Remark 1** For the Gaussian model with  $\mathbf{W} = 1/\sigma^2 \mathbf{I}_n$  and at the true threshold  $\psi = \psi_0$ ,



equivalence (i) reduces to

$$\begin{aligned} & M_A \left( \mathbf{X}_\psi \hat{\boldsymbol{\beta}}_{ML} \right) - M_A \left( \mathbf{X}_\psi \hat{\boldsymbol{\beta}}_{rB} \right) \geq 0 \text{ for all non-negative definite matrices } \mathbf{A} \\ \Leftrightarrow & \boldsymbol{\delta}^T (2\sigma_\delta^2/\sigma^2 \mathbf{I} + \mathbf{Z})^{-1} \boldsymbol{\delta} \leq \sigma^2, \end{aligned} \quad (7)$$

where  $\mathbf{Z} = (\mathbf{X}_1^T \mathbf{X}_1)^{-1} + (\mathbf{X}_2^T \mathbf{X}_2)^{-1}$ , while equivalence (ii) reduces to

$$\begin{aligned} & M \left( \mathbf{X}_\psi \hat{\boldsymbol{\beta}}_{ML} \right) - M \left( \mathbf{X}_\psi \hat{\boldsymbol{\beta}}_{rB} \right) \geq 0 \\ \Leftrightarrow & \boldsymbol{\delta}^T \mathbf{Z} (\sigma^2/\sigma_\delta^2 \mathbf{I}_p + \mathbf{Z})^{-2} \boldsymbol{\delta} \leq \sigma^2 \left\{ p - \text{tr} \left( \mathbf{I}_p + \sigma^2/\sigma_\delta^2 \mathbf{Z} \right)^{-2} \right\} \end{aligned} \quad (8)$$

**Remark 2** Using a singular value decomposition  $\mathbf{Z} = \mathbf{U} \text{diag}(\eta_1, \dots, \eta_p) \mathbf{U}^T$  and writing  $\mathbf{U}^T \boldsymbol{\delta} = \boldsymbol{\alpha}$ , inequality (8) is equivalent to

$$\sum_{i=1}^p \frac{\eta_i (2\sigma_\delta^2/\sigma^2 + \eta_i - \alpha_i^2/\sigma^2)}{(\sigma_\delta^2/\sigma^2 + \eta_i)^2} \geq 0,$$

which holds in particular if

$$\frac{\alpha_{max}^2 - \eta_{min} \sigma^2}{2} \leq \sigma_\delta^2 \quad (9)$$

with  $\alpha_{max} = \max_{1 \leq i \leq p} \alpha_i$  and  $\eta_{min} = \min_{1 \leq i \leq p} \eta_i$ . Analogously, we obtain

$$\frac{p\alpha_{max}^2 - \eta_{min} \sigma^2}{2} \leq \sigma_\delta^2 \quad (10)$$

333 as a condition for inequality (7) to be satisfied.

334 **Remark 3** The left-hand side of inequalities (7) – (10) decreases when  $\delta_1, \dots, \delta_p$  dimin-  
 335 ish in magnitude, while the right-hand side increases with growing variance  $\sigma^2$ , that is,  
 336 when the signal-to-noise ratio becomes smaller. Hence, it is reasonable to expect regular-  
 337 ized Bayesian regression coefficient estimates to be particularly superior to their profile  
 338 likelihood counterparts in settings previously identified as problematic.

339 **Remark 4** The regularized Bayesian estimator for the regression coefficients  
 340  $\hat{\boldsymbol{\beta}}_{rB} = (\mathbf{X}_\psi^T \mathbf{W} \mathbf{X}_\psi + \mathbf{H})^{-1} \mathbf{X}_\psi^T \mathbf{W} \mathbf{z}$  closely resembles the ridge estimator. However, the  
 341 special form of the penalty matrix  $\mathbf{H} = \sigma_\delta^{-2} \begin{pmatrix} \mathbf{I}_p & -\mathbf{I}_p \\ -\mathbf{I}_p & \mathbf{I}_p \end{pmatrix}$  (instead of just  $\sigma_\delta^{-2} \mathbf{I}_{2p}$  in the  
 342 ridge regression) has considerable implications for the estimator.

## 343 6 Simulations

344 To assess the performance of the suggested approach and the estimator  $\hat{\psi}_{rB}$  in particular  
 345 we performed a simulation study. We report results for eight different settings covering  
 346 both situations in which common estimators produce reliable results and others in which  
 347 they are prone to be distorted.

348 The difference between setting 1 and setting 2 is in the conditional distribution of  $y_i$ : in  
 349 the first case,  $y_i | \mathbf{X}_i^T, q_i$  is normally distributed, in the second case it follows a Poisson  
 350 distribution. The design matrix  $\mathbf{X}$  is random, each entry  $x_{ij} \sim U[0, 1]$  for setting 1,  
 351  $x_{ij} \sim U[0, 0.01]$  for setting 2. The transition variable follows a uniform distribution  
 352  $q_i \sim U[0, 1]$ . As this implies  $P(d(\psi_0) < 1) \approx 0.46$  for setting C, we base our simulations  
 353 on a fixed sample of transition variables  $q_i = i/n, i = 1, \dots, n$ . This way, we ensure that  
 354  $d(\psi_0) = 1$ , hence, that  $\mathcal{L}_p(\psi_0)$  is always well-defined. While settings A and B differ from  
 355 setting C in the threshold ( $\psi_0 = 0.5$  for A and B;  $\psi_0 = 0.3$  for C), setting A is distinct  
 356 from settings B and C in the signal-to-noise ratio, which we control by the choice of

	Normal response (1)			
	A	B	C	D
$\psi_0$	0.5	0.5	0.3	0.3
$\delta$	$U[-0.5, 0.5]$	$U[-0.5, 0.5]$	$U[-0.5, 0.5]$	$U[-0.25, 0.25]$
$\text{var}(y_i)$	$0.75^2$	$1.5^2$	$1.5^2$	$0.25^2$
$x_{ij}$	$U[0, 1]$	$U[0, 1]$	$U[0, 1]$	$U[0, 1]$
$p$	30	30	30	10
	Poisson response (2)			
	A	B	C	D
$\psi_0$	0.5	0.5	0.3	0.3
$\delta$	$U[10, 20]$	$U[0, 10]$	$U[0, 10]$	$U[10, 20]$
$x_{ij}$	$U[0, 0.01]$	$U[0, 0.01]$	$U[0, 0.01]$	$U[0, 0.01]$
$p$	30	30	30	10

Table 1: Differences between simulation settings.

	MSE( $\hat{\psi}$ )			MSE( $\mathbf{X}_{\hat{\psi}}\hat{\beta}$ )	
	pL	nB	rB	pL	rB
1 A	0.006	0.035	0.002	0.00002	0.00001
1 B	0.040	0.093	0.024	0.00009	0.00005
1 C	0.272	0.264	0.089	0.00009	0.00005
1 D	0.401	0.738	0.191	0.00001	0.00001
2 A	0.000	0.003	0.000	0.05953	0.01947
2 B	0.013	0.115	0.004	0.07625	0.02916
2 C	0.083	0.116	0.014	0.57250	0.02266
2 D	0.146	0.358	0.036	0.72387	0.18669

Table 2: Simulation results.

357  $\delta = \beta_2 - \beta_1$  relative to the variance of the observations. For setting 1 A – C, the difference  
 358  $\delta \sim U[-0.5, 0.5]$  and random variables are simulated with variances  $\text{var}(y_i) = 0.75^2$   
 359 (setting A) and  $\text{var}(y_i) = 1.5^2$  (settings B and C). The effects of increasing the signal-to-  
 360 noise ratio and shifting  $\psi_0$  on  $\ell_p(\psi)$  are illustrated in Fig. 2. The mode of  $\ell_p(\psi)$  is less  
 361 pronounced in setting 1B than in 1A. Further, the number of local maxima rises and they  
 362 become more distinctive as we move to setting 1B and then to 1C. For setting 2 A the  
 363 difference  $\delta \sim U[10, 20]$ , whereas  $\delta \sim U[0, 10]$  for settings 2 B and C. Setting D features  
 364 less nuisance parameters than A – C;  $p = \dim(\beta_1) = \dim(\beta_2) = 10$  for D,  $p = 30$  for A –  
 365 C. The sample size is  $n = 100$ .

366 Table 1 sums up differences between settings. Regression coefficients  $\beta_1$  are drawn from  
 367 a Poisson distribution with mean 10. To be unambiguous, parameters  $\delta$  and  $\beta_1$  are fixed;  
 368 we randomly generate them once at the beginning of the simulation according to the  
 369 distributions specified. Our Monte Carlo sample contains  $R = 1000$  replications.

370 With regard to the threshold parameter, we summarize simulation results in Fig. 3, where  
 371 the boxplots of the threshold estimators are shown and in the left half of Table 2, where  
 372  $\text{MSE}(\hat{\psi}) = \frac{1}{R} \sum_{r=1}^R \left( \hat{\psi}^{(r)} / \psi - 1 \right)^2$  are reported. All three estimators  $\hat{\psi}_{pL}$ ,  $\hat{\psi}_{nB}$  and  $\hat{\psi}_{rB}$   
 373 perform well given a high signal-to-noise ratio and  $\psi_0$  in the middle of  $\Psi$  (setting A).

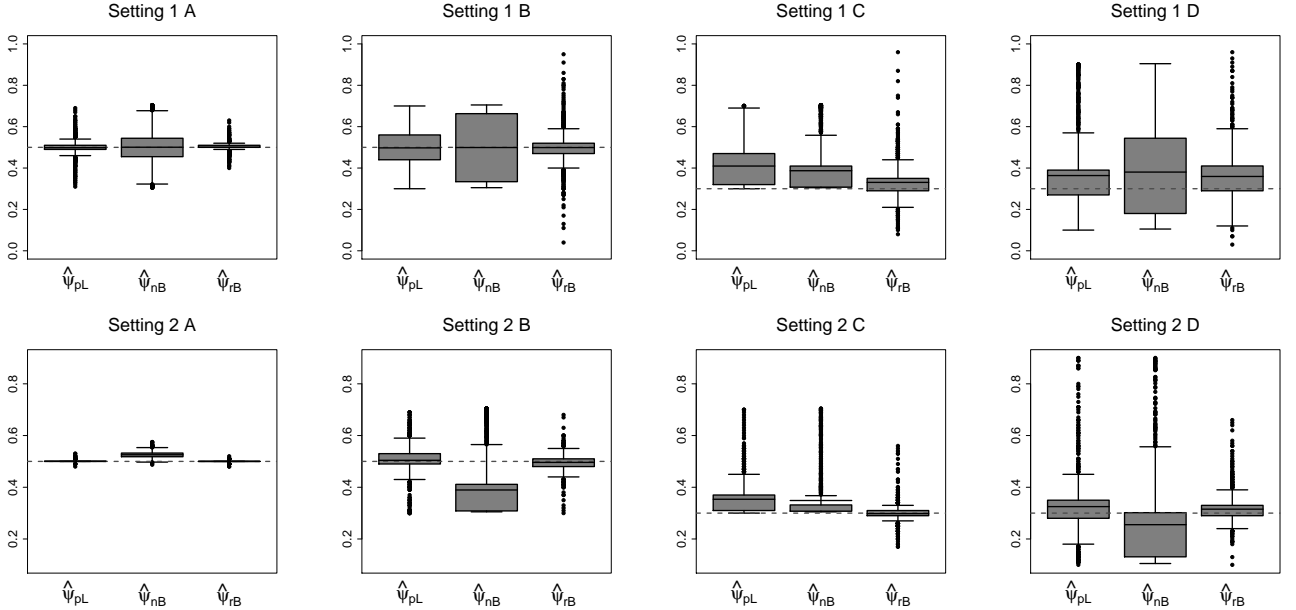


Figure 3: Boxplots for different threshold estimators and selected simulations. Dashed lines indicate the true threshold  $\psi_0$ , black lines in the boxes are sample means.

374 Lowering the signal-to-noise ratio (setting B) alters the results: we observe nearly unbiased  
 375 estimates  $\hat{\psi}_{pL}$ ,  $\hat{\psi}_{nB}$  and  $\hat{\psi}_{rB}$ , but due to its very small variance the latter stands out by its  
 376 small mean squared error. When we shift the true threshold towards the boundary of  $\Psi$   
 377 (setting C),  $\hat{\psi}_{rB}$  clearly outperforms both  $\hat{\psi}_{pL}$  and  $\hat{\psi}_{nB}$ . The differences in mean squared  
 378 error are more pronounced with a greater number of nuisance parameters  $p$ , but are still  
 379 visible in simulations with smaller ratio  $p/n$  (setting D).

380 To complement findings for the threshold estimators with results concerning estimation of  
 381 the model as a whole, in particular including the regression coefficients' estimator, we con-  
 382 sider the mean squared error for the entire model. The regularized Bayesian approach fares  
 383 better in general. While the mean squared error is much lower for simulations with normal  
 384 than with Poisson response, differences between the likelihood and regularized Bayesian  
 385 framework are more marked for the latter. The right half of Table 2 contains details. We  
 386 denote  $\text{MSE}(\mathbf{X}_{\hat{\psi}}\hat{\boldsymbol{\beta}}) = \frac{1}{R} \sum_{r=1}^R \frac{1}{n} \left( \mathbf{X}_{\hat{\psi}^{(r)}}\hat{\boldsymbol{\beta}}^{(r)} / \mathbf{X}_{\hat{\psi}}\boldsymbol{\beta} - \mathbf{1} \right)^T \left( \mathbf{X}_{\hat{\psi}^{(r)}}\hat{\boldsymbol{\beta}}^{(r)} / \mathbf{X}_{\hat{\psi}}\boldsymbol{\beta} - \mathbf{1} \right)$  with  
 387 the division  $\mathbf{X}_{\hat{\psi}^{(r)}}\hat{\boldsymbol{\beta}}^{(r)} / \mathbf{X}_{\hat{\psi}}\boldsymbol{\beta}$  defined elementwise and  $\mathbf{1} = (1, \dots, 1)^T \in \mathbb{R}^n$ . Note that

388 in settings 2 the Fisher scoring algorithm for the estimation of generalized regression  
389 models can be unstable for small sample sizes, sometimes leading to a false convergence.  
390 Therefore, we excluded such outliers (5% of the Monte Carlo sample) from the calculation  
391 of  $\text{MSE}(\mathbf{X}_{\hat{\psi}}\hat{\boldsymbol{\beta}})$  for settings 2 A – D.

## 392 **7 Applications**

393 This work is originally motivated by the application of threshold vector error correction  
394 models in price transmission analysis. Such models are rather involved, but one important  
395 characteristic in this context is that they contain a large number of parameters besides the  
396 threshold and available data series are typically short in relation to the complexity of the  
397 model. Greb et al. (2013) investigates the merits of the regularized Bayesian approach for  
398 this particular model; simulations demonstrate the superiority of the regularized Bayesian  
399 threshold estimator (see figure 1, figure 2, and table 1 in Greb et al., 2013) and two real  
400 data examples confirm its relevance in practice.

### 401 **7.1 Cross-country growth behavior**

402 As another application of the regularized Bayesian threshold estimator, we consider the  
403 case of economic growth modeling. Durlauf and Johnson (1995) estimate a standard  
404 growth model using cross-sectional data on a sample of 96 countries and investigate  
405 whether the coefficients of this model differ across sub-sets of countries depending on their  
406 initial conditions. Their analysis is based on the so-called regression tree methodology  
407 (Breiman et al., 1984), which suggests three thresholds based on two different transition  
408 variables for this application.

409 Hansen (2000) revisits their paper. Using the Durlauf and Johnson data he estimates a  
 410 regression

$$\begin{aligned} & \log(GDP)_{i,1985} - \log(GDP)_{i,1960} \\ &= \zeta + \beta \log(GDP)_{i,1960} + \pi_1 \log(INV)_i + \pi_2 \log(n_i + g + \delta) + \pi_3 \log(SCHOO)_{i,1960} + \varepsilon_i \end{aligned}$$

411 which explains real GDP growth between 1960 and 1985 in country  $i$ ,  
 412  $\log(GDP)_{i,1985} - \log(GDP)_{i,1960}$ , using real GDP in 1960  $GDP_{i,1960}$ , the invest-  
 413 ment to GDP ratio  $INV_i$ , the growth rate of the working-age population  $n_i$ , the rate of  
 414 technological change  $g$ , the rate of depreciation of physical and human capital stocks  $\delta$ ,  
 415 and the fraction of working-age population enrolled in secondary school  $(SCHOO)_{i,1960}$ .  
 416 With reference to Durlauf and Johnson (1995), he sets  $g + \delta = 0.05$ . He tests for a  
 417 threshold effect based on either one of transition variables they propose. He only finds  
 418 evidence based on the transition variable  $\log(GDP)_{i,1960}$  and calculates the profile  
 419 likelihood (or, equivalently, least squares) estimate as  $\hat{\psi}_{pL} = 6.76$  together with an  
 420 asymptotic 95% confidence interval [6.39, 7.49].

421 This corresponds to an estimate of \$863 per capita GDP in 1960 with an associated  
 422 confidence interval of [\$594, \$1794]. Hansen (2000) acknowledges that while the confidence  
 423 interval seems rather tight (given observations for  $GDP_{i,1960}$  ranging from \$383 to \$12362),  
 424 it effectively contains 40 of the 96 countries in the sample. This is in line with the number  
 425 of local maxima in the profile likelihood function which hints at the uncertainty inherent  
 426 in this method (Fig. 4). In addition, the fact that  $\hat{\psi}_{pL}$  leaves only 18 observations in  
 427 the first regime gives rise to concern that the threshold might be located close to the  
 428 boundary of  $\Psi$ . We know that the profile likelihood is typically distorted if this is the  
 429 case.

430 Hence, we reestimate the model with the regularized Bayesian estimator. The latter  
 431 depends on the parameterization of the transition variable. As  $\log(GDP)_{i,1960}$  is an  
 432 explanatory variable, we choose the parameterization  $q_i = \log(GDP)_{i,1960}$ . Figure 4

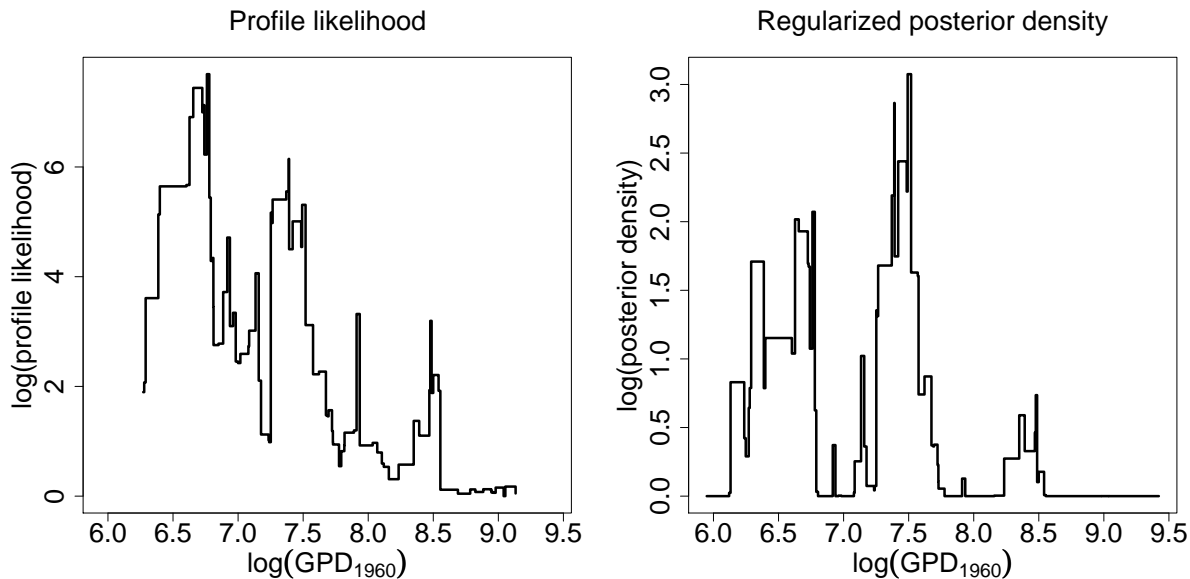


Figure 4: Profile likelihood and regularized posterior density for a threshold based on the transition variable  $q_i = \log(GDP)_{i,1960}$ .

433 shows that the resulting posterior density differs considerably from the profile likelihood  
 434 function and that the location of the maximum shifts. This is not surprising given the  
 435 deformations often observed for the profile likelihood function close to the boundary of  
 436 the threshold parameter space. The posterior median is located at  $\hat{\psi}_{rB} = 7.37$  compared  
 437 with Hansen's (2000)  $\hat{\psi}_{pL} = 6.76$ . It implies that, for the 43 poorest countries, coefficients  
 438 for the growth model are distinct from the rest, whereas the profile likelihood estimate  
 439 implicates that this is only the case for the poorest 18 countries.

440 While it is not possible to state conclusively that the regularized Bayesian estimate is more  
 441 appropriate from an economic perspective, the shapes of the likelihoods in Fig. 4 and the  
 442 fact that the profile likelihood estimate is near the boundary of its domain suggests that  
 443 the latter may be distorted by the weaknesses of the profile likelihood method discussed  
 444 above.

445 Comparing profile likelihood estimates for the regression coefficients with their regular-  
 446 ized Bayesian counterparts, we note that there is much less difference between regimes  
 447 (see table 7.1). Moreover, the difference between the two regimes as estimated within

448 the regularized Bayesian framework is negligible. This is in line with Hansen’s (2000)  
449 finding that the null hypothesis of no threshold is not rejected at the 5%-level (Hansen,  
450 2000, page 587). The example demonstrates the effect of using the suggested regularized  
451 Bayesian estimator instead of the profile likelihood estimator in small samples with a  
452 multi-modal profile likelihood and high uncertainty attached to the estimate  $\hat{\psi}_{pL}$  obtained  
453 by maximizing it.

	1st regime					2nd regime				
	$\hat{\zeta}$	$\hat{\beta}$	$\hat{\pi}_1$	$\hat{\pi}_2$	$\hat{\pi}_3$	$\hat{\zeta}$	$\hat{\beta}$	$\hat{\pi}_1$	$\hat{\pi}_2$	$\hat{\pi}_3$
pL	4.31 (3.21)	-0.66 (0.33)	0.23 (0.14)	-0.29 (0.92)	0.02 (0.11)	3.66 (0.85)	-0.32 (0.07)	0.50 (0.11)	-0.49 (0.30)	0.36 (0.07)
rB	3.36 (0.85)	-0.41 (0.08)	0.47 (0.09)	-0.60 (0.28)	0.22 (0.06)	3.37 (0.85)	-0.38 (0.07)	0.47 (0.09)	-0.62 (0.28)	0.20 (0.07)

Table 3: Regressions coefficient estimates. ”pL” refers to the profile likelihood, ”rB” to the regularized Bayesian framework. Standard errors in parentheses below the estimates.

## 454 7.2 Effects of climate on snowshoe hare survival

455 In our final example, we study a famous dataset of snowshoe hare abundance in the  
456 main drainage of Hudson Bay in Canada. It consists of annual observations starting in  
457 the 19th century. A preeminent feature of the data is cyclical fluctuations in the hare  
458 population, see Fig. 5. These have been ascribed to the predator-prey relationship between  
459 lynx and snowshoe hares. Samia and Chan (2011) highlight selected references and further  
460 investigate one strand of the discussion focusing on the effect of snow conditions on  
461 hunting efficiency in different phases of the cycle. To this end, they estimate a generalized  
462 threshold regression model with the hare count  $y_t$  as a Poisson distributed response whose  
463 mean is related to the explanatory variables via a log-link,

$$\log(\mu_t) = \beta_0 + \beta_1 D_t + \begin{cases} \sum_{i=1}^3 \beta_{1,i} \log(y_{t-i} + 1) + \beta_{1,4} w_{t-1} & y_{t-d} \leq \psi, \\ \sum_{i=1}^3 \beta_{2,i} \log(y_{t-i} + 1) + \beta_{2,4} w_{t-1} & y_{t-d} > \psi \end{cases}$$



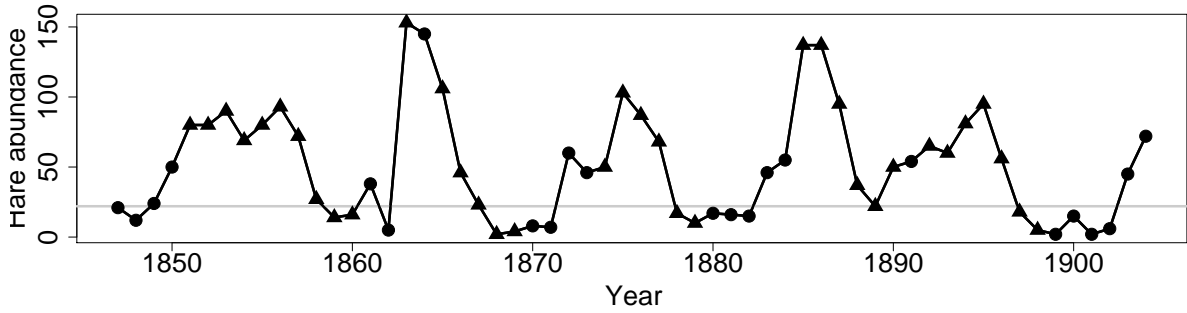


Figure 5: Annual hare abundance. Observations estimated to belong to the lower regime are plotted as dots, observations estimated to belong to the upper regime as triangles. The horizontal grey line indicates the location of the estimated threshold,  $\hat{\psi}_{rB} = 22$ .

464 for the years  $t = 1844, \dots, 1904$ . Apart from the regression coefficients and the threshold,  
 465 the delay of the transition variable  $d$  is included as an additional parameter,  $d \in \{1, 2, 3\}$ .  
 466 As the count for the year  $t = 1863$  is considered an outlier, the model contains a dummy  
 467 variable  $D_t = I(t = 1863)$ . The covariate  $w_t$  denotes the detrended annual winter climate  
 468 index of the North Atlantic Oscillation, published at [www.cru.uea.ac.uk/cru/data/nao](http://www.cru.uea.ac.uk/cru/data/nao).  
 469 We follow Samia and Chan (2011) in estimating this model. Our analysis is based on the  
 470 series of hare abundance initially presented graphically by MacLulich (1937) which we  
 471 calibrate with data available online; it is included in the supplementary material to this  
 472 paper.

473 The series of 61 observations is rather short and maximizing out regression coefficients  
 474 leaves us with a profile likelihood function for  $(d, \psi)$  which is characterized by various  
 475 local maxima; it is displayed in the upper row of Fig. 6 for  $d = 1, 2, 3$  and  $\psi \in \Psi^*(1)$ . In  
 476 addition, we cannot rule out overdispersion. Hence, we are confronted with a setting in  
 477 which the regularized Bayesian estimate can be more reliable than the profile likelihood  
 478 estimate. This becomes evident in the second row of Fig. 6, which shows the posterior  
 479 densities for  $\psi$  corresponding to  $d = 1, 2, 3$ . While we obtain a profile likelihood estimate  
 480  $(\hat{d}_{pL}, \hat{\psi}_{pL}) = (3, 55)$ , the regularized Bayesian estimator yields  $(\hat{d}_{rB}, \hat{\psi}_{rB}) = (2, 22)$  with  $\hat{d}_{rB}$

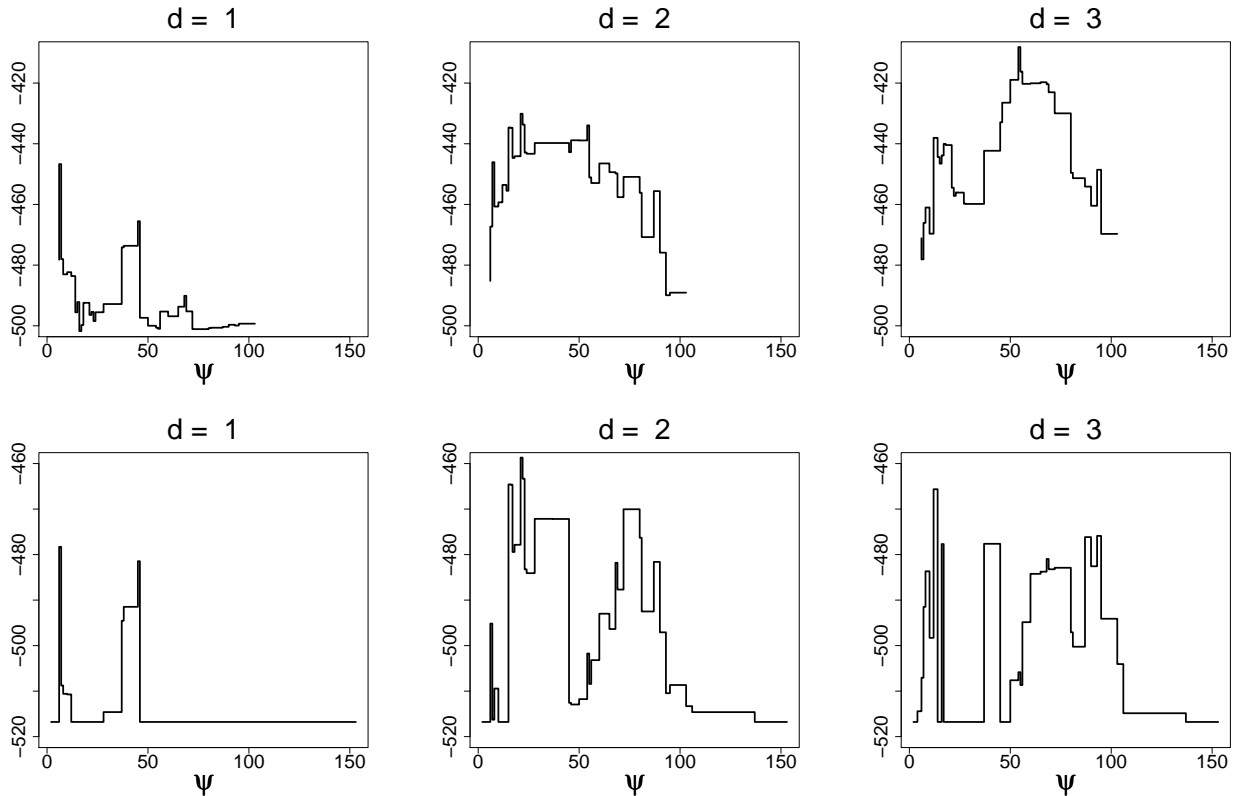


Figure 6: Log-likelihood functions (upper row) and log-posterior densities (lower row) for different delays of the transition variable.

481 calculated as the posterior median based on a flat prior on  $\{1, 2, 3\}$ .

482 When referring to Samia and Chan (2011) we have to keep in mind that their results  
 483 diverge slightly from ours and are not directly comparable as we were not able to obtain the  
 484 data they used. Yet, their profile likelihood estimate is still very close,  $(\hat{d}_{pL}, \hat{\psi}_{pL}) = (3, 69)$ .

485 However, they discard this estimate in favor of  $(\hat{d}, \hat{\psi}) = (2, 25)$ , giving heuristic arguments  
 486 based on residual analysis. The latter also allows for a very plausible interpretation.

487 Apparently, our regularized Bayesian estimate  $(\hat{d}_{rB}, \hat{\psi}_{rB}) = (2, 22)$  is close to the preferred  
 488 estimate in Samia and Chan (2011). In fact, the difference in estimated thresholds only  
 489 has implications for a single observation ( $t = 1869$ ). Except for this, thresholds induce  
 490 identical allocations of observations to regimes (in the respective datasets), as is clearly  
 491 visible when comparing our Fig. 5 with Fig. 1 in Samia and Chan (2011). Hence,  
 492 the regularized Bayesian estimator enables us to attain a meaningful estimate directly

493 avoiding any arbitrary modification of the suggested estimation method as done by Samia  
494 and Chan (2011). Coefficient estimates are similar in both modeling frameworks.

## 495 **8 Conclusions**

496 In this work we describe settings in which estimation of generalized threshold regression  
497 models can be problematic. We suggest a new regularized Bayesian estimator which out-  
498 performs standard estimators. In particular, the suggested threshold estimator is defined  
499 on the whole parameter space and thus circumvents the subjective and often misleading  
500 restriction of the threshold domain which standard estimators require. Moreover, regu-  
501 larizing the posterior density at the boundary of its domain helps to improve estimation,  
502 especially if the true threshold is close to this boundary. Employing the empirical Bayes  
503 approach, we can use built-in functions for generalized linear mixed models in statistics  
504 software and obtain estimates with little additional numerical effort and without the use  
505 of Markov chain Monte Carlo or other sampling techniques. Inference about the estimated  
506 parameter can be carried out in the standard Bayesian manner. Simulation studies and  
507 a real-data example confirm the effectiveness and relevance of our method.

## 508 **Acknowledgments**

509 The authors are very grateful to the responsible editor, the associate editor and two  
510 anonymous referees for numerous constructive comments, which have greatly improved  
511 the paper. The support of the German Research Foundation (Deutsche Forschungsgemein-  
512 schaft) as part of the Institutional Strategy of the University of Göttingen is acknowledged  
513 by all the authors. The third author also acknowledges funding through FOR 916 and  
514 the VW foundation.

## 515 Appendix

### 516 Derivation of equation (5)

517 We obtain the approximate posterior (5) as follows. Laplace approximation produces

$$\begin{aligned}
& \int_{\mathbb{R}^p} \int_{\mathbb{R}^p} p(\mathbf{y}|\boldsymbol{\beta}_1, \boldsymbol{\delta}, \psi, \phi, \sigma_\delta^2, \mathbf{X}, \mathbf{q}) p(\boldsymbol{\delta}|\sigma_\delta^2) d\boldsymbol{\delta} d\boldsymbol{\beta}_1 \\
&= (2\pi)^{-p/2} |\sigma_\delta^2 \mathbf{I}_p|^{-1/2} \int_{\mathbb{R}^p} \int_{\mathbb{R}^p} \exp\{-\kappa(\boldsymbol{\delta}, \boldsymbol{\beta}_1)\} d\boldsymbol{\delta} d\boldsymbol{\beta}_1 \\
&= (2\pi)^{p/2} |\sigma_\delta^2 \mathbf{I}_p|^{-1/2} \exp\left\{-\kappa(\hat{\boldsymbol{\delta}}, \hat{\boldsymbol{\beta}}_1)\right\} \left| \frac{\partial^2 \kappa}{\partial(\boldsymbol{\delta}, \boldsymbol{\beta}_1) \partial(\boldsymbol{\delta}, \boldsymbol{\beta}_1)^T}(\hat{\boldsymbol{\delta}}, \hat{\boldsymbol{\beta}}_1) \right|^{-1/2} + \mathcal{O}(n^{-1})
\end{aligned}$$

518 for  $\kappa(\boldsymbol{\delta}, \boldsymbol{\beta}_1) = -\sum_{i=1}^n \frac{y_i \theta_i - b(\theta_i)}{\phi} - c(y_i, \phi) + \frac{1}{2\sigma_\delta^2} \boldsymbol{\delta}^T \boldsymbol{\delta}$  and  $(\hat{\boldsymbol{\delta}}, \hat{\boldsymbol{\beta}}_1) = \underset{(\boldsymbol{\delta}, \boldsymbol{\beta}_1) \in \mathbb{R}^{2p}}{\operatorname{argmax}} -\kappa(\boldsymbol{\delta}, \boldsymbol{\beta}_1)$ .

519 Given the derivatives

$$\frac{\partial \kappa}{\partial \boldsymbol{\delta}}(\boldsymbol{\delta}) = -\sum_{i=1}^n \frac{(y_i - \mu_i)(\mathbf{X}_2)_i}{\phi b''(\theta_i) g'(\mu_i)} + \frac{1}{\sigma_\delta^2} \boldsymbol{\delta} = -\mathbf{X}_2^T \mathbf{W} \mathbf{G}(\mathbf{y} - \boldsymbol{\mu}) + \frac{1}{\sigma_\delta^2} \boldsymbol{\delta},$$

520

$$\frac{\partial \kappa}{\partial \boldsymbol{\beta}_1}(\boldsymbol{\beta}_1) = -\sum_{i=1}^n \frac{(y_i - \mu_i)(\mathbf{X})_i}{\phi b''(\theta_i) g'(\mu_i)} = -\mathbf{X}^T \mathbf{W} \mathbf{G}(\mathbf{y} - \boldsymbol{\mu}),$$

521 and

$$\frac{\partial^2 \kappa}{\partial(\boldsymbol{\delta}, \boldsymbol{\beta}_1) \partial(\boldsymbol{\delta}, \boldsymbol{\beta}_1)^T} = \begin{pmatrix} \mathbf{X}_2^T \mathbf{W} \mathbf{X}_2 + (1/\sigma_\delta^2) \mathbf{I}_p & \mathbf{X}_2^T \mathbf{W} \mathbf{X} \\ \mathbf{X}^T \mathbf{W} \mathbf{X}_2 & \mathbf{X}^T \mathbf{W} \mathbf{X} \end{pmatrix} \quad (11)$$

522 for  $\mathbf{W}^{-1} = \operatorname{diag}\{\phi b''(\theta_i) g'(\mu_i)^2\}$  and  $\mathbf{G} = \operatorname{diag}\{g'(\mu_i)\}$ , we obtain

$$\left| \frac{\partial^2 \kappa}{\partial(\boldsymbol{\delta}, \boldsymbol{\beta}_1) \partial(\boldsymbol{\delta}, \boldsymbol{\beta}_1)^T} \right| = \left| \mathbf{X}_2^T \mathbf{W} \mathbf{X}_2 + (1/\sigma_\delta^2) \mathbf{I}_p \right| \left| \mathbf{X}^T \mathbf{V}^{-1} \mathbf{X} \right|$$

523 using basic matrix algebra.

524 To find  $\hat{\boldsymbol{\delta}}$  and  $\hat{\boldsymbol{\beta}}_1$ , we iteratively solve

$$\mathbf{X}_2^T \mathbf{W} \mathbf{G}(\mathbf{y} - \boldsymbol{\mu}) = \frac{1}{\sigma_\delta^2} \boldsymbol{\delta} \text{ and } \mathbf{X}^T \mathbf{W} \mathbf{G}(\mathbf{y} - \boldsymbol{\mu}) = 0$$

525 via Fisher-scoring: Starting at  $\hat{\boldsymbol{\delta}} = \boldsymbol{\delta}_0$  and  $\hat{\boldsymbol{\beta}}_1 = (\boldsymbol{\beta}_1)_0$ , we solve

$$\mathcal{I}(\boldsymbol{\delta}_m, \boldsymbol{\beta}_m) \begin{pmatrix} \boldsymbol{\delta}_{m+1} \\ (\boldsymbol{\beta}_1)_{m+1} \end{pmatrix} = \mathcal{I}(\boldsymbol{\delta}_m, \boldsymbol{\beta}_m) \begin{pmatrix} \boldsymbol{\delta}_m \\ (\boldsymbol{\beta}_1)_m \end{pmatrix} + s(\boldsymbol{\delta}_m, (\boldsymbol{\beta}_1)_m),$$

526  $\mathcal{I} = \partial^2 \kappa / \partial(\boldsymbol{\delta}, \boldsymbol{\beta}_1) \partial(\boldsymbol{\delta}, \boldsymbol{\beta}_1)^T$  and  $s = -\partial \kappa / \partial(\boldsymbol{\delta}, \boldsymbol{\beta}_1)$ , or, more explicitly,

$$\left\{ \mathbf{X}_2^T \mathbf{W}_m \mathbf{X}_2 + \frac{1}{\sigma_\delta^2} \mathbf{I}_p \right\} \boldsymbol{\delta}_{m+1} + \mathbf{X}_2^T \mathbf{W}_m \mathbf{X}(\boldsymbol{\beta}_1)_{m+1} = \mathbf{X}^T \mathbf{W}_m \mathbf{z}_m$$

527 and

$$\mathbf{X}^T \mathbf{W}_m \mathbf{X}_2 \boldsymbol{\delta}_{m+1} + \mathbf{X}^T \mathbf{W}_m \mathbf{X}(\boldsymbol{\beta}_1)_{m+1} = \mathbf{X}^T \mathbf{W}_m \mathbf{z}_m,$$

528 where  $\mathbf{z}_m = \mathbf{X}_2 \boldsymbol{\delta}_m + \mathbf{X}(\boldsymbol{\beta}_1)_m + \mathbf{G}_m(\mathbf{y} - \boldsymbol{\mu}_m)$ . This yields

$$\hat{\boldsymbol{\beta}}_1 = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \tilde{\mathbf{z}} \quad \text{and} \quad \hat{\boldsymbol{\delta}} = \sigma_\delta^2 \mathbf{X}_2^T \mathbf{V}^{-1} (\tilde{\mathbf{z}} - \mathbf{X} \hat{\boldsymbol{\beta}}_1),$$

529 where  $\mathbf{V} = \mathbf{W}^{-1} + \sigma_\delta^2 \mathbf{X}_2 \mathbf{X}_2^T$  and  $\tilde{\mathbf{z}} = \mathbf{X}_2^T \hat{\boldsymbol{\delta}} + \mathbf{X} \hat{\boldsymbol{\beta}}_1 + \mathbf{G}(\mathbf{y} - \boldsymbol{\mu})$ , with  $\mathbf{W}$ ,  $\mathbf{G}$  and  $\boldsymbol{\mu}$

530 evaluated at  $\boldsymbol{\delta} = \hat{\boldsymbol{\delta}}$  and  $\boldsymbol{\beta}_1 = \hat{\boldsymbol{\beta}}_1$  (Harville, 1977).

531 With this, we can now further simplify the posterior. Following Breslow and Clayton

532 (1993) in replacing

$$-2 \sum_{i=1}^n \{y_i \theta_i - b(\theta_i)\} \quad \text{by the chi-squared statistic} \quad \sum_{i=1}^n \frac{(y_i - \mu_i)^2}{b''(\theta_i)}$$

533 we can exploit the identity

$$\mathbf{V}^{-1} (\tilde{\mathbf{z}} - \hat{\boldsymbol{\beta}}_1) = \mathbf{W} (\tilde{\mathbf{z}} - \mathbf{X} \hat{\boldsymbol{\beta}}_1 - \mathbf{X}_2 \hat{\boldsymbol{\delta}}),$$

534 which results in

$$(\tilde{\mathbf{z}} - \mathbf{X} \hat{\boldsymbol{\beta}}_1 - \mathbf{X}_2 \hat{\boldsymbol{\delta}})^T \mathbf{W} (\tilde{\mathbf{z}} - \mathbf{X} \hat{\boldsymbol{\beta}}_1 - \mathbf{X}_2 \hat{\boldsymbol{\delta}}) = (\tilde{\mathbf{z}} - \hat{\boldsymbol{\beta}}_1)^T \mathbf{V}^{-1} (\tilde{\mathbf{z}} - \hat{\boldsymbol{\beta}}_1) - \frac{1}{\sigma_\delta^2} \hat{\boldsymbol{\delta}}^T \hat{\boldsymbol{\delta}},$$

535 and, hence,

$$\begin{aligned} & \exp \left\{ \sum_{i=1}^n \frac{y_i \theta_i - b(\theta_i)}{\phi} + c(y_i, \phi) - \frac{1}{2\sigma_\delta^2} \hat{\boldsymbol{\delta}}^T \hat{\boldsymbol{\delta}} \right\} \\ & \approx \exp \left\{ -\frac{1}{2} (\tilde{\mathbf{z}} - \mathbf{X} \hat{\boldsymbol{\beta}}_1 - \mathbf{X}_2 \hat{\boldsymbol{\delta}})^T \mathbf{W} (\tilde{\mathbf{z}} - \mathbf{X} \hat{\boldsymbol{\beta}}_1 - \mathbf{X}_2 \hat{\boldsymbol{\delta}}) + \sum_{i=1}^n c(y_i, \phi) - \frac{1}{2\sigma_\delta^2} \hat{\boldsymbol{\delta}}^T \hat{\boldsymbol{\delta}} \right\} \\ & = \exp \left\{ -\frac{1}{2} (\tilde{\mathbf{z}} - \hat{\boldsymbol{\beta}}_1)^T \mathbf{V}^{-1} (\tilde{\mathbf{z}} - \hat{\boldsymbol{\beta}}_1) + \sum_{i=1}^n c(y_i, \phi) \right\}. \end{aligned}$$

536 Altogether, this leaves us with

$$\begin{aligned}
& \int_{\mathbb{R}^p} \int_{\mathbb{R}^p} p(\mathbf{y}|\boldsymbol{\beta}_1, \boldsymbol{\delta}, \psi, \phi, \sigma_\delta^2, \mathbf{X}, \mathbf{q}) p(\boldsymbol{\delta}|\sigma_\delta^2) d\boldsymbol{\delta} d\boldsymbol{\beta}_1 \\
&= (2\pi)^{p/2} |\sigma_\delta^2 \mathbf{I}_p|^{-1/2} \exp \left\{ \sum_{i=1}^n \frac{y_i \theta_i - b(\theta_i)}{\phi} + c(y_i, \phi) - \frac{1}{2\sigma_\delta^2} \hat{\boldsymbol{\delta}}^T \hat{\boldsymbol{\delta}} \right\} |\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X}|^{-1/2} \\
&\quad \cdot \left| \mathbf{X}_2^T \mathbf{W} \mathbf{X}_2 + (1/\sigma_\delta^2) \mathbf{I}_p \right|^{-1/2} + \mathcal{O}(n^{-1}) \\
&\approx (2\pi)^{p/2} \exp \left\{ -\frac{1}{2} (\tilde{\mathbf{z}} - \hat{\boldsymbol{\beta}}_1)^T \mathbf{V}^{-1} (\tilde{\mathbf{z}} - \hat{\boldsymbol{\beta}}_1) + \sum_{i=1}^n c(y_i, \phi) \right\} |\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X}|^{-1/2} \\
&\quad \cdot \left| \sigma_\delta^2 \mathbf{X}_2^T \mathbf{W} \mathbf{X}_2 + \mathbf{I}_p \right|^{-1/2} + \mathcal{O}(n^{-1}).
\end{aligned}$$

### 537 Details for Theorem 1

538 Basic matrix algebra yields a representation of the regularized Bayesian estima-  
539 tors  $\hat{\boldsymbol{\beta}}_1 = (\mathbf{X}^T \mathbf{V} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{z}$  and  $\hat{\boldsymbol{\beta}}_2 = \hat{\boldsymbol{\beta}}_1 + \hat{\boldsymbol{\delta}} = \hat{\boldsymbol{\beta}}_1 + \sigma_\delta^2 \mathbf{X}_2^T \mathbf{V}^{-1} (\mathbf{z} - \mathbf{X} \hat{\boldsymbol{\beta}}_1)$ , where  
540  $\mathbf{V} = \mathbf{W}^{-1} + \sigma_\delta^2 \mathbf{X}_2 \mathbf{X}_2^T$ , as  $\hat{\boldsymbol{\beta}}_{r_B} = (\mathbf{X}_\psi^T \mathbf{W} \mathbf{X}_\psi + \mathbf{H})^{-1} \mathbf{X}_\psi^T \mathbf{W} \mathbf{z}$ . To obtain equivalence (i),  
541 we employ a theorem by Theobald (1974, theorem 1) stating that for two estimators  $\hat{\boldsymbol{\beta}}_*$   
542 and  $\hat{\boldsymbol{\beta}}_{**}$

$$\begin{aligned}
& M_A(\hat{\boldsymbol{\beta}}_*) - M_A(\hat{\boldsymbol{\beta}}_{**}) \geq 0 \text{ for all non-negative definite matrices } \mathbf{A} \\
& \Leftrightarrow E(\hat{\boldsymbol{\beta}}_* - \boldsymbol{\beta})(\hat{\boldsymbol{\beta}}_* - \boldsymbol{\beta})^T - E(\hat{\boldsymbol{\beta}}_{**} - \boldsymbol{\beta})(\hat{\boldsymbol{\beta}}_{**} - \boldsymbol{\beta})^T \text{ is non-negative definite.}
\end{aligned}$$

543 The equivalence then follows from

$$\begin{aligned}
& E(\mathbf{X}_\psi \hat{\boldsymbol{\beta}}_{r_B} - \mathbf{X}_\psi \boldsymbol{\beta})(\mathbf{X}_\psi \hat{\boldsymbol{\beta}}_{r_B} - \mathbf{X}_\psi \boldsymbol{\beta})^T = \mathbf{X}_\psi (\mathbf{X}_\psi^T \mathbf{W} \mathbf{X}_\psi + \mathbf{H})^{-1} \mathbf{X}_\psi^T \mathbf{W} \mathbf{X}_\psi (\mathbf{X}_\psi^T \mathbf{W} \mathbf{X}_\psi + \mathbf{H})^{-1} \mathbf{X}_\psi^T \\
& \quad + \mathbf{X}_\psi (\mathbf{X}_\psi^T \mathbf{W} \mathbf{X}_\psi + \mathbf{H})^{-1} (\mathbf{B} + \mathbf{H}) \boldsymbol{\beta} \boldsymbol{\beta}^T (\mathbf{B}^T + \mathbf{H}) (\mathbf{X}_\psi^T \mathbf{W} \mathbf{X}_\psi + \mathbf{H})^{-1} \mathbf{X}_\psi^T.
\end{aligned}$$

544 Using  $E(\hat{\boldsymbol{\beta}}_{r_B} - \boldsymbol{\beta})^T (\hat{\boldsymbol{\beta}}_{r_B} - \boldsymbol{\beta}) = \text{tr} \left\{ E(\hat{\boldsymbol{\beta}}_{r_B} - \boldsymbol{\beta})(\hat{\boldsymbol{\beta}}_{r_B} - \boldsymbol{\beta})^T \right\}$  then yields equivalence (ii).

545 For remark 1,  $\psi = \psi_0$  implies  $\mathbf{B} = 0$ . Consequently,

$$\mathbf{D} \{ (\mathbf{I} + \mathbf{C}) \mathbf{H} - (\mathbf{B} + \mathbf{H}) \boldsymbol{\beta} \boldsymbol{\beta}^T (\mathbf{B}^T + \mathbf{H}) + \mathbf{C} \mathbf{B} \boldsymbol{\beta} \boldsymbol{\beta}^T \mathbf{B}^T \mathbf{C}^T \} \mathbf{D}^T \geq 0$$

546 reduces to

$$\mathbf{D} \{ (\mathbf{I} + \mathbf{C}) \mathbf{H} - \mathbf{H} \boldsymbol{\beta} \boldsymbol{\beta}^T \mathbf{H} \} \mathbf{D}^T \geq 0.$$

547 Assuming that  $\text{rank}(\mathbf{X}) = p$ , this is equivalent to

$$\begin{aligned} & (\mathbf{X}_\psi^T \mathbf{W} \mathbf{X}_\psi + \mathbf{H})^{-1} \{(\mathbf{I} + \mathbf{C}) \mathbf{H} - \mathbf{H} \boldsymbol{\beta} \boldsymbol{\beta}^T \mathbf{H}\} (\mathbf{X}_\psi^T \mathbf{W} \mathbf{X}_\psi + \mathbf{H})^{-1} \geq 0 \\ \Leftrightarrow & (\mathbf{I} + \mathbf{C}) \mathbf{H} - \mathbf{H} \boldsymbol{\beta} \boldsymbol{\beta}^T \mathbf{H} = 2\mathbf{H} + \mathbf{H} (\mathbf{X}_\psi^T \mathbf{W} \mathbf{X}_\psi)^{-1} \mathbf{H} - \mathbf{H} \boldsymbol{\beta} \boldsymbol{\beta}^T \mathbf{H} \geq 0 \end{aligned}$$

548 since  $\mathbf{X}_\psi^T \mathbf{W} \mathbf{X}_\psi + \mathbf{H}$  is positive definite and symmetric. Taking advantage of a result by

549 Gruber (1990, theorem 2.5.3), this amounts to

$$\begin{aligned} & \boldsymbol{\beta}^T \mathbf{H} \left( 2\sigma_\delta^2 \mathbf{H} + \sigma_\delta^2 \mathbf{H} (\mathbf{X}_\psi^T \mathbf{W} \mathbf{X}_\psi)^{-1} \mathbf{H} \right)^+ \mathbf{H} \boldsymbol{\beta} \leq 1/\sigma_\delta^2 \\ \Leftrightarrow & \boldsymbol{\delta}^T \left\{ 2\sigma_\delta^2 \mathbf{I} + (\mathbf{X}_1^T \mathbf{W} \mathbf{X}_1)^{-1} + (\mathbf{X}_2^T \mathbf{W} \mathbf{X}_2)^{-1} \right\}^{-1} \boldsymbol{\delta} \leq 1. \end{aligned}$$

For  $\mathbf{W} = 1/\sigma_\delta^2 \mathbf{I}$  this is equivalent to

$$\boldsymbol{\delta}^T \left\{ 2\sigma_\delta^2/\sigma^2 \mathbf{I} + (\mathbf{X}_1^T \mathbf{X}_1)^{-1} + (\mathbf{X}_2^T \mathbf{X}_2)^{-1} \right\}^{-1} \boldsymbol{\delta} \leq \sigma^2.$$

550 Basic matrix calculations suffice to obtain the rest of this as well as the following remarks.

## 551 References

- 552 Andrews, D. (1993). Tests for parameter instability and structural change with unknown  
553 change point. *Econometrica*, 61(4):821–856.
- 554 Bacon, D. and Watts, D. (1971). Estimating the transition between two intersecting  
555 straight lines. *Biometrika*, 58(3):525–534.
- 556 Bai, J. (1997). Estimation of a change point in multiple regression models. *Review of*  
557 *Economics and Statistics*, 79(4):551–563.
- 558 Breiman, L., Friedman, J., Olshen, R., and Stone, C. (1984). *Classification and Regression*  
559 *Trees*. Wadsworth, Belmont, California.
- 560 Breslow, N. and Clayton, D. (1993). Approximate inference in generalized linear mixed  
561 models. *Journal of the American Statistical Association*, 88(421):9–25.
- 562 Chan, K. and Tsay, R. (1998). Limiting properties of the least squares estimator of a  
563 continuous threshold autoregressive model. *Biometrika*, 85(2):413–426.
- 564 Chan, N. H. and Kutoyants, Y. A. (2012). On parameter estimation of threshold autore-  
565 gressive models. *Statistical inference for stochastic processes*, 15(1):81–104.
- 566 Doodson, A. (1917). Relation of the mode, median and mean in frequency curves.  
567 *Biometrika*, 11(4):425–429.
- 568 Durlauf, S. and Johnson, P. (1995). Multiple regimes and cross-country growth behaviour.  
569 *Journal of Applied Econometrics*, 10(4):365–384.
- 570 Feder, P. (1975). On asymptotic distribution theory in segmented regression problems–  
571 identified case. *The Annals of Statistics*, 3(1):49–83.
- 572 Geweke, J. and Terui, N. (1993). Bayesian threshold autoregressive models for nonlinear  
573 time series. *Journal of Time Series Analysis*, 14(5):441–454.



- 574 Greb, F., von Cramon-Taubadel, S., Krivobokova, T., and Munk, A. (2013). The estima-  
575 tion of threshold models in price transmission analysis. *American Journal of Agricul-  
576 tural Economics*. First published online: March 28, 2013.
- 577 Gruber, M. H. (1990). *Regression estimators: A comparative study*. Academic Press  
578 (Boston).
- 579 Hansen, B. (2000). Sample splitting and threshold estimation. *Econometrica*, 68(3):575–  
580 603.
- 581 Hansen, B. (2011). Threshold autoregression in economics. *Statistics and Its Interface*,  
582 4(2):123–128.
- 583 Hansen, B. and Seo, B. (2002). Testing for two-regime threshold cointegration in vector  
584 error-correction models. *Journal of Econometrics*, 110(2):293–318.
- 585 Harville, D. (1977). Maximum likelihood approaches to variance component estimation  
586 and to related problems. *Journal of the American Statistical Association*, 72(358):320–  
587 338.
- 588 Kendall, M. G. (1943). *The advanced theory of statistics, Vol. 1*. J.B. Lippincott company.
- 589 Lee, S., Seo, M., and Shin, Y. (2011). Testing for threshold effects in regression models.  
590 *Journal of the American Statistical Association*, 106(493):220–231.
- 591 MacLulich, D. (1937). *Fluctuations in the numbers of the varying hare (Lepus ameri-  
592 canus)*. University of Toronto Press.
- 593 Nelder, J. (1972). Discussion of a paper by D.V. Lindley and A.F.M. Smith. *Journal of  
594 the Royal Statistical Society. Series B (Methodological)*, 34:18–20.
- 595 Samia, N. and Chan, K. (2011). Maximum likelihood estimation of a generalized threshold  
596 stochastic regression model. *Biometrika*, 98(2):433–448.

- 597 Samia, N., Chan, K., and Stenseth, N. (2007). A generalized threshold mixed model for  
598 analyzing nonnormal nonlinear time series, with application to plague in Kazakhstan.  
599 *Biometrika*, 94(1):101–118.
- 600 Severini, T. (2000). *Likelihood methods in statistics*. Oxford University Press, USA.
- 601 Shun, Z. and McCullagh, P. (1995). Laplace approximation of high dimensional integrals.  
602 *Journal of the Royal Statistical Society. Series B (Methodological)*, 57(4):749–760.
- 603 Theobald, C. (1974). Generalizations of mean square error applied to ridge regression.  
604 *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 103–106.
- 605 Tikhonov, A., Arsenin, V., and John, F. (1977). *Solutions of Ill-posed Problems*. Vh  
606 Winston Washington, DC.
- 607 Tong, H. (2011). Threshold models in time series analysis – 30 years on. *Statistics and*  
608 *Its Interface*, 4(2):107–118.
- 609 Tong, H. and Lim, K. (1980). Threshold autoregression, limit cycles and cyclical data.  
610 *Journal of the Royal Statistical Society. Series B (Methodological)*, 42(3):245–292.
- 611 Van Dijk, D., Teräsvirta, T., and Franses, P. (2002). Smooth transition autoregressive  
612 models—a survey of recent developments. *Econometric Reviews*, 21(1):1–47.
- 613 Yu, P. (2012). Likelihood estimation and inference in threshold regression. *Journal of*  
614 *Econometrics*, 167(1):274–294.