

# Protocol for Comparative Methods Study on Penalisation Methods for Multi-State Models

Ms. Chantelle Cornett (CC), Dr. Glen Martin (GM), Dr. Niels Peek (NP)

## Overview

Multi-state models (MSMs) are useful for prediction of complex clinical settings where patients may transition through several intermediate states before experiencing a terminal event [1-5]. A challenge of MSMs is the combinatorial complexity in models as the number of states being modelled grows. This increases computation time, but also risks overfitting of some of the transition models within the MSM, especially those with fewer events/ observations.

Overfitting refers to when a model is too complex and captures noise/random fluctuations in the data rather than the true relationships between variables. In the case of MSMs, this usually manifests as too many parameters and overly complex model structures. One way overfitting can be diagnosed for MSMs is to assess the model's performance on data that was not used in the training process. Techniques such as cross-validation and comparison of prediction accuracy metrics such as the C-index between the training and test sets can aid in diagnosing overfitting in an MSM [6]. Overfitting in this context has been recognised as an issue in multi-state models [7][8].

One potential solution to this is to share parameters between transitions within the model [5]. For example, one could assume that certain groups of transitions share a common baseline and/or predictor effects. The properties of doing this in a prediction context are unknown. Likewise, in a prediction context, penalisation methods are often advocated to help minimise the risk of overfitting. Therefore, it is of interest to examine the impact of penalisation methods on this situation with or without sharing of parameter terms across the transitions. While the use of shared parameters across transitions to reduce the parameter space is well-documented, the concept of penalisation in this context is less well-studied.

Although sample size requirements for MSMs have not yet been developed, it is known that sample size is an important factor in determining a well-functioning model [9][10][11]. While penalisation may prove to be a useful tool to avoid overfitting, many papers have noted that penalisation is not a substitute for an adequate sample size [12][13][14]. Therefore, we will not go below the minimum sample size for a Cox model developed for each transition (in the absence of the MSM sample size criteria). We will explore the use of penalisation and parameter sharing as further protection against overfitting.

## Definitions

We broadly define the terms relevant to this review:

Term	Definition
MLE	Maximum likelihood estimate
LASSO	Least absolute shrinkage and selection operator
CPRD	Clinical practice research datalink
Fold	A fold refers to the subset of the dataset that is used for training and validation
Cross validation loop	A cross validation loop is the process by which a dataset is split into k folds. The model would then be trained and validated k times using a different fold as the validation set each time
Horseshoe prior	A type of Bayesian prior distribution, designed to encourage shrinkage of parameter estimates toward zero while allowing for the possibility of some parameters to remain large
Metropolis-Hasting method	A Markov Chain Monte Carlo method (MCMC) used for generating a sequence of random samples from a probability distribution
Burn-in period	The initial phase of the simulation where the Markov chain is allowed to converge to the target distribution
Bootstrap sampling	A sample made by randomly drawing observations with replacement from the original dataset. This creates new datasets of the same size as the original but with some observations omitted and others repeated

Table 1: A list of definitions used in the protocol.

## Methodology

All analysis for this study will be conducted using R (version 4.3.1) in the RStudio environment under the MacOS Ventura operating system on a MacBook Pro (with Apple M1 Pro processor). Code for simulations and models will be provided on GitHub to ensure reproducibility.

## Aims

This study aims to compare the predictive performance of current penalisation methodology for multi-state models. We will be considering the following methods of implementing the penalisation within prediction models:

- MLE no shrinkage
- MLE with uniform shrinkage factor applied to each transition
- LASSO penalised likelihood
- Fused-LASSO penalised likelihood [5]
- Reduced-rank methods [15]
- Bayesian approach with penalising prior to pull coefficients down towards an overall model-level value [16]

These methods will be applied to simulated data, and then applied to CPRD data [17].

## Data Simulation

We will consider two data generating models: first, a simple illness-death model, and second, a more complex 5-state model. This is done so that we can investigate the impact of the number of states (and transitions) on the predictive performance results. Each simulation will have  $n=200,000$  individuals. In order to compromise between reducing the Monte-Carlo standard errors and having a reasonable computation time, we will repeat the simulation 500 times ( $n_{sim} = 500$ ).

Data will need to be passed in a long data format. We will first simulate the model in wide format for ease of computation, and transpose this to create the long format we require.

### Simple Illness-Death Model

We will start by simulating data for a simple illness-death model (Figure 1). We will use the following Q-matrix:

$$Q = \begin{bmatrix} -0.02 & 0.015 & 0.005 \\ 0 & -0.01 & 0.01 \\ 0 & 0 & 0 \end{bmatrix},$$

using the *simmulti.msm* function from the *MSM* package to get a list of states and their transition times (starting in the healthy state at time 0) for each subject. The numbers in this Q-matrix were chosen in order to get sufficient progress through the states in the model while compromising on computation time.

The survival distributions for each transition were drawn from a hazard function of:

$$\exp[0.5 * age_i - 0.5 * gender_i + 0.5 * BMI_i] * \lambda_{j,k}$$

The baseline hazard for each transition ( $\lambda_{j,k}$ ) was assumed to follow a Weibull distribution, where parameters were chosen based on the mean event time:

$$\lambda_{1,2} = t^{-0.791} (\lambda = 1, k = 0.219)$$

$$\lambda_{1,3} = t^{-0.791} (\lambda = 1, k = 0.219)$$

$$\lambda_{2,3} = t^{-0.766} (\lambda = 1, k = 0.233)$$

We will base the structure of the data on the structure of the European Society for Blood and Marrow Transplantation dataset (ebmt4)[18]. The final dataset will have the following columns before being passed into *msprep* to transform the dataset into a format suitable for use by the *mstate* package:

Column	Description
Subject	The subject IDs.
timeToIll	The time passed since starting in the healthy state until being in the ill state.
timeToDeath	The time passed since starting in the healthy state until being in the death state (if death state is not reached this will take the value of 200, indicating censoring).
illStat	An indicator variable taking a value of 1 if the subject reaches the ill state, and 0 if they do not.
deathStat	An indicator variable taking a value of 1 if the subject reaches the death state, and 0 if they do not.
Age	A factor variable, taking values $\leq 20$ , 20-40, $\geq 40$ . These are assigned to subjects with probabilities 0.332, 0.531, 0.137 in accordance with the estimated worldwide age distribution [19].
Gender	A factor variable, taking values 0 and 1 and assigned with equal probability.
BMI	A continuous variable, indicating the body mass index of an individual. This will be drawn from a standard Gaussian distribution.

Table 2: A list of columns for the simple illness-death model and their descriptions.

### 5-State Model

The next model we will simulate is a 5-state model (Figure 2). We will use the following Q-matrix:

$$Q = \begin{bmatrix} -0.025 & 0.02 & 0 & 0 & 0.005 \\ 0 & -0.035 & 0.025 & 0 & 0.01 \\ 0 & 0 & -0.035 & 0.025 & 0.01 \\ 0 & 0 & 0 & -0.035 & 0.035 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix},$$

using the *simmulti.msm* function from the *MSM* package to get a list of states and their transition times (starting in the healthy state at time 0) for each subject. The numbers in this Q-matrix were chosen in order to get sufficient progress through the states in the model while compromising on computation time.

The survival distributions for each transition were drawn from a hazard function of:

$$\exp[0.5 * age_i - 0.5 * gender_i + 0.5 * BMI_i] * \lambda_{j,k}$$

The baseline hazard for each transition ( $\lambda_{j,k}$ ) was assumed to follow a Weibull distribution, where parameters were chosen based on the mean event time:

$$\lambda_{1,2} = t^{-0.705} (\lambda = 1, k = 0.295)$$

$$\lambda_{2,3} = t^{-0.690} (\lambda = 1, k = 0.310)$$

$$\lambda_{3,4} = t^{-0.666} (\lambda = 1, k = 0.333)$$

$$\lambda_{4,5} = t^{-0.625} (\lambda = 1, k = 0.375)$$

$$\lambda_{1,5} = t^{-0.752} (\lambda = 1, k = 0.248)$$

$$\lambda_{2,5} = t^{-0.742} (\lambda = 1, k = 0.258)$$

$$\lambda_{3,5} = t^{-0.729} (\lambda = 1, k = 0.271)$$

We will base the structure of the data on the structure of the ebmt4 dataset [18]. The final dataset will have the following columns, before being passed into *msprep* to transform the dataset into a format suitable for use by the *mstate* package:

Column	Description
Subject	The subject IDs.
timeToMild	The time passed since starting in the healthy state until being in the mild state.
timeToModerate	The time passed since starting in the healthy state until being in the moderate state (if moderate state is not reached this will take the value of 200, indicating censoring).
timeToSevere	The time passed since starting in the healthy state until being in the severe state (if severe state is not reached this will take the value of 1000, indicating censoring).
timeToDeath	The time passed since starting in the healthy state until being in the death state (if death state is not reached this will take the value of 1000, indicating censoring).
mildStat	An indicator variable taking a value of 1 if the subject reaches the mild state, and 0 if they do not.
modStat	An indicator variable taking a value of 1 if the subject reaches the moderate state, and 0 if they do not.
sevStat	An indicator variable taking a value of 1 if the subject reaches the severe state, and 0 if they do not.
deathStat	An indicator variable taking a value of 1 if the subject reaches the death state, and 0 if they do not.
Age	A factor variable, taking values $\leq 20$ , 20-40, $\geq 40$ . These are assigned to subjects with probabilities 0.332, 0.531, 0.137 in accordance with the estimated worldwide age distribution [19].
Gender	A factor variable, taking values 0 and 1 and assigned with equal probability.
BMI	A continuous variable, indicating the body mass index of an individual. This will be drawn from a standard Gaussian distribution.

Table 3: A list of columns for the 5-state model and their descriptions.

## Penalisation Methods

Since we are using simulated survival data, a Cox proportional hazards model will be implemented using each of the penalisation methods in turn.

### *MLE No Shrinkage*

In order to apply the MLE without shrinkage, the *msprep* function of the *mstate* package will be used to prepare the data for the *coxph* function. This facilitates fitting a cox proportional hazards model for each transition in turn.

### *MLE with Uniform Shrinkage Factor Applied to Each Transition*

To apply a uniform shrinkage factor to each transition, we will estimate the shrinkage factor post model fit using the heuristic shrinkage factor:

$$S_{VH} = 1 - \frac{p}{LR},$$

where  $p$  is the total number of predictor parameters for the full set of candidate predictors and  $LR$  is the likelihood ratio (chi-squared) statistic for the fitted model defined as:

$$LR = -2(\ln L_{null} - \ln L_{model}),$$

where  $\ln L_{null}$  is the log-likelihood of a model with no predictors (e.g. intercept-only cox proportional hazards model), and  $\ln L_{model}$  is the log-likelihood of the final model [20]. In the model, we will multiply each covariate by the shrinkage factor, allowing the shrinkage factor to be applied uniformly to both covariates across each transition in the model. This will be fitted using the *mscoxph* function from the *mstate* package.

### *LASSO Penalised Likelihood*

In order to apply the LASSO penalised likelihood, we will first concentrate on finding the optimum tuning parameter for the LASSO. This will be done using cross-validation:

1. Split the dataset into 10 folds.
2. For each potential value of the tuning parameter, perform a cross-validation loop.
3. Choose the tuning parameter that provides the best model performance across all folds.

In R, this will be done using the *penMSM* package, using the *penMSM* function with *type = lasso* to perform the cross-validated LASSO.

### *Fused-LASSO Penalised Likelihood (Shared Parameters)*

Just like in the LASSO penalised likelihood, we will first find the optimum tuning parameter for the fused-LASSO penalisation. When using the *penMSM* function, we will set the *type* argument to *fused* to indicate a fused-LASSO approach.

### *Reduced-Rank Method*

To implement the reduced-rank method of penalisation, we will use the *redrank* function of the *mstate* package. We will choose the value of  $R$  (the rank of the matrix used for the reduced-rank problem) based on cross-validation and choose the rank that minimizes prediction error on unseen data.

### Bayesian Approach with Penalising Prior (Shared Parameters)

We have decided to use a horseshoe prior as described in Beesley (2020). A horseshoe prior is a continuous shrinkage prior, that strongly shrinks moderate and weak signals to zero, while still allowing very large signals. The prior takes the form:

$$\begin{aligned}f(\theta_k|\lambda_k, \tau) &= N(\alpha, \lambda_k^2 \tau^2) \\f(\lambda_k) &= \text{Cauchy}^+(0,1) \\f(\tau) &= \text{Cauchy}^+(0,1),\end{aligned}$$

where smaller  $\tau$  implements a greater global shrinkage,  $\theta_k$  are the log-hazard ratios for each of the  $k$  transition models,  $\alpha$  is a non-zero constant (assuming that the ‘shared’ coefficient is non-zero across all transitions), and  $\text{Cauchy}^+$  indicated the half-Cauchy distribution.  $\tau$ ,  $\lambda_k$ , and  $\theta_k$  can be sampled in turn using the Metropolis-Hasting method. As we are applying the horseshoe prior to an MSM, we will define  $\tau$  separately for each transition.

We will be using the *brm* function from the *brms* package to implement this type of penalisation. We will specify priors for the intercept, standard deviation, and each covariate in the model. The models using the simulated data will not be complex, so we will only use 2 chains. We will start with 2000 (with a burn-in period of 1000 iterations) iterations, as this provides a good balance between computational efficiency. Using trace-plots and the Gelman-Rubin statistic, we will assess convergence; if these diagnostic tools show that insufficient convergence has occurred, we will increase the number of iterations.

### Packages to Use

Package	Use
mstate	Contains functions needed for data preparations, description, and hazard estimation for multi-state models [21]. This includes reduced rank proportional hazards model for multi-state models.
tidyverse	A collection of packages for preparing, wrangling, and visualising data [22].
genSurv	Generation of survival data with one (binary) time-dependent covariate [23].
shrink	Used for uniform (global) shrinkage [24].
brms	Used for fitting a Bayesian Cox model with a horseshoe prior [25].
penMSM	Used for providing efficient LASSO penalisation [7].

Table 4: A list of packages and their uses.

### Performance Measures

- Computation time – As combinatorial complexity increases, so does complexity.
- Time-dependent ROC curves – To quantify discrimination.
- C-Index – To evaluate predictive performance.
- Calibration plots – Show any potential mismatches between observed and predicted probabilities in the data.
- Transition-specific calibration plots – To evaluate overfitting.
- Monte-Carlo standard error.

## Validation

### Internal Validation

For each iteration of the simulation, we will be generating a large ( $n \geq 500,000$ ) sample size and will split this into two parts: a development subset and a validation subset. We will ensure that the validation subset is large ( $n \geq 200,000$ ). For each penalisation method, we will fit a model using the validation subset using the same model specification as we used in the original analysis.

For each model fitted using the validation subset, we will calculate the following performance measures:

- Time-dependent ROC curves [26]
- C-Index [27]
- Calibration plots [28]
- Transition-specific calibration plots

These will be compared with the performance measures of the original model to check that the results are consistent between the development and validation subsets.

### Author Contributions

CC and GM contributed to the conception and design of the study. All authors contributed to discussions around updates to the protocol.

CC wrote the protocol.

CC performed the data simulation and method implementation, iteratively updating the protocol where necessary.



## Bibliography

- [1] Heinz Freisling et al, 2020. Lifestyle factors and risk of multimorbidity of cancer and cardiometabolic diseases: a multinational cohort study. *National Library of Medicine*, 10 Jan.
- Harrison Reeder et al, 2022. Penalized Estimation of Frailty-Based Illness-Death Models for Semi-Competing Risks. *Biometrics*.
- [2] Hein Putter et al, 2006. Estimation and Prediction in a Multi-State Model for Breast Cancer. *Biometrical Journal*, 16 June, 48(3), pp. 366-380.
- [3] Niels Keiding et al, 2001. Multi-state models and outcome prediction in bone marrow transplantation. *Statistics in Medicine*, 30 May, 20(12), pp. 1871-1885.
- [4] Jeffrey Cannon et al, 2017. Rheumatic Heart Disease Severity, Progression and Outcomes: A Multi-State Model. *Journal of the American Heart Association*, 2 March.
- [5] Holger Sennhenn-Reulen, 2016. Structured fusion lasso penalized multi-state models. *Statistics in Medicine*, 23 June, 35(25), p. 4637-4659.
- [6] Jake Lever et al, 2016. Model selection and overfitting. *Nature Methods*, 30 August, 13, p703-704.
- [7] Holger Reulen, 2023. Package 'penMSM'. [Online]  
Available at: <https://cran.r-project.org/web/packages/penMSM/penMSM.pdf>  
[Accessed 15 11 2023]
- [8] Michael Crowther et al, 2017. Parametric multistate survival models: Flexible modelling allowing transition-specific distributions with application to estimating clinically useful measures of effect differences. *Statistics in Medicine*, 05 September, 36(29), p4719-4742.
- [9] Richard Riley et al, 2018. Minimum sample size for developing a multivariable prediction model: Part I - Continuous Outcomes. *Statistics in Medicine*, 22 October, 38(7), p1262-1275.
- [10] Richard Riley et al, 2020. Calculating the sample size required for developing a clinical prediction model. *BMJ*, 18 March, 8238.
- [11] Alexander Pate et al, 2023. Minimum sample size for developing a multivariable prediction model using multinomial logistic regression. *Statistical Methods in Medical Research*, 19 January, 32(3), p555-571.
- [12] Ben Van Calster et al, 2020. Regression shrinkage methods for clinical prediction models do not guarantee improved performance: Simulation study. *Statistical Methods in Medical Research*, 13 May.
- [13] Richard Riley et al, 2021. Penalization and shrinkage methods produced unreliable clinical prediction models especially when sample size was small. *Journal of Clinical Epidemiology*, April, 132, p88-96.
- [14] Glen Martin et al, 2021. Developing clinical prediction models when adhering to minimum sample size recommendations: The importance of quantifying bootstrap variability in tuning parameters and predictive performance. *Statistical Methods in Medical Research*, December, 30(12), p2545-2561.
- [15] Marta Fiocco et al, 2009. Reduced-rank proportional hazards regression and simulation-based prediction for multi-state models. *Statistics in Medicine*, 21 April, 27(21), pp. 4340-4358.
- [16] Lauren Beesley et al, 2020. Bayesian variable selection and shrinkage strategies in a complicated modelling setting with missing data: A case study using multistate models. *Statistical Modelling*, 21(1-2), p11-29.
- [17] Shivani Padmanabhan, 2011. CPRD GOLD Data Dictionary [Restricted]
- [18] Gustavo Soutinho, 2021. ebmt4: Data from the European Society for Blood and Marrow Transplantation. Available at: <https://rdr.io/github/gsoutinho/markovMSM/man/ebmt4.html>  
[Accessed 19 11 2023]
- [19] Carmen Ang, 2021. Visualizing the World's Population by Age Group. Available at: <https://www.visualcapitalist.com/the-worlds-population-2020-by-age/>  
[Accessed 19 11 2023]

- [20] Richard Riley et al, 2018. Minimum sample size for developing a multivariable prediction model: PART II - binary. *Statistics in Medicine*, 24 Oct, 38(7), p1276–1296.
- [21] Hein Putter et al, 2021. *mstate: Data Preparation, Estimation and Prediction in Multi-State Models*. [Online]  
Available at: <https://cran.r-project.org/web/packages/mstate/index.html>  
[Accessed 9 11 2023].
- [22] Hadley Wickham, 2023. *tidyverse: Easily Install and Load the 'Tidyverse'*. [Online]  
Available at: <https://cran.r-project.org/web/packages/tidyverse/index.html>  
[Accessed 9 11 2023].
- [23] Arthur Araujo, 2021. *genSurv: Generating Multi-State Survival Data*. [Online]  
Available at: <https://cran.r-project.org/web/packages/genSurv/index.html>  
[Accessed 9 11 2023].
- [24] Daniela Dunkler, 2023. *Package 'shrink'*. [Online]  
Available at: <https://cran.r-project.org/web/packages/shrink/shrink.pdf>  
[Accessed 9 11 2023].
- [25] Paul Christian Bürkner, 2023. *Package 'brms'*. [Online]  
Available at: <https://cran.r-project.org/web/packages/brms/brms.pdf>  
[Accessed 9 11 2023].
- [26] Frank Harrell et al, 1996. Multivariable Prognostic Models: Issue in Developing Models, Evaluating Assumptions and Adequacy, and Measure and Reducing Errors. *Statistics in Medicine* 15(4), p361-387.
- [27] Patrick Heagerty et al, 2005. Survival Model Predictive Accuracy and ROC Curves. *Biometrics* 61(1), p92-105.
- [28] Ben Van Calster et al, 2019. Calibration: the Achilles heel of predictive analytics. *BMC Medicine* 17(230).

## Figures

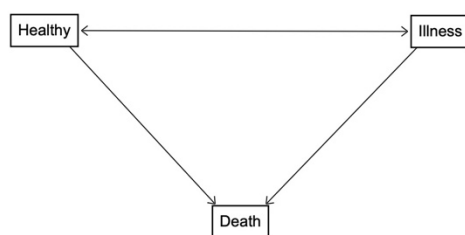


Figure 1 - Simple illness-death model.

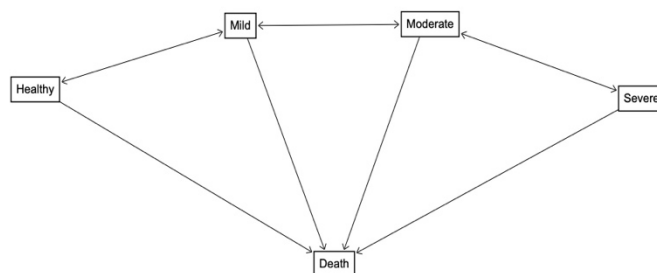


Figure 2 – An example of a 5-state model.