REGULARISED JUMP MODELS FOR REGIME IDENTIFICATION AND FEATURE SELECTION

Anonymous authors

006

007

012 013

014

015

016

017

018

019

021

022

023

024

025

026

027

028 029

031 032

033

034

035

037

045

046 047 048

052

Paper under double-blind review

ABSTRACT

A regime modelling framework can be employed to address the complexities of financial markets. Under the framework, market periods are grouped into distinct regimes, each distinguished by similar statistical characteristics. Regimes in financial markets are not directly observable but are often manifested in market and macroeconomic variables. The objective of regime modelling is to accurately identify the active regime from these variables at a point in time, a process known as regime identification. One way to enhance the accuracy of regime identification is to select features that are most responsible for statistical differences between regimes, a process known as *feature selection*. Models based on the Jump Model framework have recently been developed to address the joint problem of regime identification and feature selection. In the following work, we propose a new set of models called *Regularised Jump Models* that are founded upon the Jump Model framework. These models perform feature selection that is more interpretable than that from the Sparse Jump Model, a model proposed in the literature pertaining to the Jump Model framework. Through a simulation experiment, we find evidence that these new models outperform the Standard and Sparse Jump Models, both in terms of regime identification and feature selection.

1 INTRODUCTION

The section outlines and examines the constituent models of the Jump Model framework introduced in (2). A new set of models called *Regularised Jump Models* are then introduced and compared to the existing models.

1.1 STANDARD JUMP MODEL

-

Denote the state or mode sequence associated with the sequence of observations Y by (s_1, s_2, \ldots, s_T) where each $s_t \in \{1, 2, \ldots, K\}$ for $t = 1, 2, \ldots, T$. K is the number of states and is assumed to be known. The K model parameters are given by $\mu_1, \mu_2, \ldots, \mu_K$ where each $\mu_k \in \mathbb{R}^p$ for $k = 1, 2, \ldots, K$. The model parameters are the conditional means of the features assigned to each of the K states, hence the notational choice $\mu_k, k = 1, 2, \ldots, K$.

Definition 1.1 (Standard Jump Model). The *Standard Jump Model* with *K* states is defined by the minimisation of the objective function

$$\sum_{t=1}^{T-1} \left[\| \boldsymbol{y}_t - \boldsymbol{\mu}_{s_t} \|^2 + \lambda \mathbb{1}_{\{s_t \neq s_{t+1}\}} \right] + \| \boldsymbol{y}_T - \boldsymbol{\mu}_{s_T} \|^2$$
(1.1)

over the model parameters $\mu_1, \mu_2, \dots, \mu_K$ and state sequence (s_1, s_2, \dots, s_T) . The term $\|\boldsymbol{y}_t - \boldsymbol{\mu}_{s_t}\|^2$ represents the squared \mathcal{L}_2 -distance between the vectors \boldsymbol{y}_t and $\boldsymbol{\mu}_{s_t}$ and $\lambda \ge 0$ is a hyperparameter.

The objective function in (1.1) can be interpreted as a tradeoff between model fitting and prior assumptions about the tendency of the sequence S to "jump", or change states. This tendency is

su

controlled by the hyperparameter λ . When $\lambda = 0$, the model reduces to splitting the dataset in at most K states and fitting one model per state, thereby generalising the K-means algorithm ((9)). For $\lambda \to \infty$, the Standard Jump Model results in a single-state model since the cost of changing states becomes prohibitive.

059 1.2 Sparse Jump Model

The Sparse Jump Model is an extension of the Standard Jump Model in its incorporation of feature selection. Let $w := (w_1, w_2, \dots, w_p) \in \mathbb{R}^p$ denote a vector of feature weights that are assumed to be the same across all states.

Feature selection is incorporated in the Sparse Jump Model by employing the criterion in (13). Thecriterion is given by

$$\max \quad \boldsymbol{w}' \left(\sum_{k=1}^{K} |C_k| \left(\boldsymbol{\mu}_k - \bar{\boldsymbol{\mu}} \right)^2 \right),$$

$$\text{ch that} \quad \|\boldsymbol{w}\|^2 \le 1, \quad \|\boldsymbol{w}\|_1 \le \kappa,$$

$$w_n \ge 0 \, \forall p,$$

$$(1.2)$$

with respect to the parameters $\mu_1, \mu_2, \dots, \mu_K$, state sequence (s_1, s_2, \dots, s_T) and feature weights w_1, w_2, \dots, w_p . The term $|C_k| (\mu_k - \bar{\mu})^2$ in (1.2) is a vector of size p whose entries are the contributions of each feature to the between-cluster sum of squares (BCSS) in the k^{th} cluster.

1075 If we consider the clusters C_1, C_2, \ldots, C_K fixed, then the feature weights will be assigned to features based on their individual BCSS contributions: features with larger BCSS contributions will be given larger weights which, in turn, optimises the overall spread between the clusters.

 $\begin{array}{ll} \kappa \in [1,\sqrt{p}] \text{ in } (1.2) \text{ is a hyperparameter that controls the degree of sparsity in the feature weights.} \\ \text{The squared \mathcal{L}_2 penalty in (1.2) serves an important role since without it, at most one element of w would be non-zero in general when features are correlated (see (14) for more details). If $w_1 = w_2 = \cdots = w_p$, then 1.2 reduces to the maximisation of the BCSS (equivalent to the minimisation of the within-cluster sum of squares (WCSS)) which is the objective of the K-means clustering algorithm. \\ \end{array}$

Combining 1.2 with a jump penalty, we give the below definition for the Sparse Jump Model:

Definition 1.2 (Sparse Jump Model). The *Sparse Jump Model* with K states is defined by the optimisation programme

087 088

090 091

058

060

070 071

$$\max \quad \boldsymbol{w}' \left(\sum_{k=1}^{K} |C_k| \left(\boldsymbol{\mu}_k - \bar{\boldsymbol{\mu}} \right)^2 \right) - \lambda \sum_{t=1}^{T} \mathbb{1}_{\{s_t \neq s_{t+1}\}}$$
such that $\|\boldsymbol{w}\|^2 \leq 1, \quad \|\boldsymbol{w}\|_1 \leq \kappa,$
 $w_p \geq 0 \ \forall p,$

$$(1.3)$$

092 093 094

096

097 098

099

105

107

with respect to the parameters $\mu_1, \mu_2, \dots, \mu_K$, state sequence (s_1, s_2, \dots, s_T) and feature weights w_1, w_2, \dots, w_p . $\kappa \in [1, \sqrt{p}]$ in (1.2) is a hyperparameter that controls the degree of sparsity in the feature weights.

If $w_1 = w_2 = \dots, w_p$, then Definition 1.2 reduces to the Standard Jump Model in Definition 1.1.

100 1.3 REGULARISED JUMP MODELS

In this section, a new approach to feature selection in the Jump Model framework is introduced. The
 approach is an adaptation of the Regularised K-means algorithm proposed in (11) to the Jump Model
 framework. Therefore, we call the models constituting this approach *Regularised Jump Models*.

106 1.4 REGULARISED K-MEANS

Definition 1.3. The Regularised K-means algorithm is defined by the minimisation of

 $\sum_{k=1}^{K} \left\{ \sum_{t \in C_{k}} \|\boldsymbol{y}_{t} - \boldsymbol{\mu}_{k}\|^{2} \right\} + \gamma \mathcal{P}\left(\underline{\boldsymbol{\mu}}\right), \qquad (1.4)$

with respect to the clusters C_1, C_2, \ldots, C_K and matrix of cluster centres $\underline{\mu} \in \mathbb{R}^{K \times p}$. $\gamma \ge 0$ is a tuning parameter that controls the amount of regularisation applied to the cluster centres. \mathcal{P} : $\mathbb{R}^{K \times p} \to \mathbb{R}$ is a penalty function that depends on $\underline{\mu}$. The first term is the objective function of standard K-means clustering introduced in (8).

From (11), below are several penalty function options which are named after their counterparts from regularised regression:

$$\mathcal{L}_{0}: \quad \mathcal{P}_{0}\left(\underline{\boldsymbol{\mu}}\right) = \sum_{j=1}^{p} \mathbb{1}_{\{\|\boldsymbol{\mu}_{.,j}\| > 0\}}$$
(1.5a)

Lasso:
$$\mathcal{P}_1\left(\underline{\mu}\right) = \sum_{j=1}^p \|\boldsymbol{\mu}_{.,j}\|_1$$
 (1.5b)

Ridge:
$$\mathcal{P}_2(\underline{\mu}) = \sum_{j=1}^p \|\boldsymbol{\mu}_{.,j}\|^2$$
 (1.5c)

Group-Lasso:
$$\mathcal{P}_3\left(\underline{\mu}\right) = \sum_{j=1}^p \|\boldsymbol{\mu}_{.,j}\|,$$
 (1.5d)

130 131

138

139

145 146 147

108

110

111

116

119

121 122 123

124 125

127 128 129

where $\mu_{..j}$ is the *j*th column of $\underline{\mu}$. The penalty on $\underline{\mu}$ balances the size of the cluster centres and their contribution to the objective function in (1.4).

The intuition for penalising the size of the cluster centres lies in the fact that when a variable does not contribute to the partitioning of the data, its estimated cluster centres will be close to the overall mean of the data (see Proposition 1 in (11) for more details).

1.5 REGULARISED JUMP MODEL EQUATION

Combining (1.4) with a jump penalty, we propose the following definition for the Regularised Jump
 Model.

Definition 1.4 (Regularised Jump Model). The *Regularised Jump Model* with K states is defined by the minimisation of the objective function

$$\sum_{t=1}^{T-1} \left\{ \| \boldsymbol{y}_t - \boldsymbol{\mu}_{s_t} \|^2 + \lambda \mathbb{1}_{\{s_t \neq s_{t+1}\}} \right\} + \| \boldsymbol{y}_T - \boldsymbol{\mu}_{s_T} \|^2 + \gamma \mathcal{P}\left(\underline{\boldsymbol{\mu}}\right),$$
(1.6)

with respect to $\underline{\mu}$ and the state sequence (s_1, s_2, \dots, s_T) . $\lambda, \gamma \ge 0$ are hyperparameters and the penalty function $\mathcal{P} : \mathbb{R}^{K \times p} \to \mathbb{R}$ can be chosen from those listed in (1.5).

 $\gamma = 0$ reduces Definition 1.4 to the Standard Jump Model in Definition 1.1. $\lambda = 0$ reduces Definition 1.4 to the Regularised K-Means model in Definition 1.3.

152 153 154

156

157 158

150

151

2 CALIBRATION OF JUMP MODELS

The section details the algorithms that perform calibration of the Jump Models outlined in the previous section. The models are then tested in a simulation experiment.

159 2.1 STANDARD JUMP MODEL CALIBRATION

Denote $\underline{\mu} := (\mu_1, \mu_2, \dots, \mu_K)' \in \mathbb{R}^{K \times p}$ the matrix of model parameters and the state sequence $S := (s_1, s_2, \dots, s_T)$. The calibration algorithm for the Standard Jump Model was proposed in (10)

and is shown in Algorithm 1. The model is calibrated using a coordinate descent algorithm that alternates between finding the model parameters $\mu_1, \mu_2, \dots, \mu_K$ that minimise the objective function (1.1) with a fixed state sequence and finding the state sequence (s_1, s_2, \dots, s_T) that minimise the objective function (1.1) with fixed $\mu_1, \mu_2, \dots, \mu_K$.

Following the algorithm proposed in (10), the process is repeated ten times at the most or until the state sequence does not change after one iteration. However, there is no guarantee that the solution reached is the global solution since the solution depends on the initial state sequence.

Therefore, we adopt the initialisation method in (10): the coordinate descent algorithm is run from ten different state sequences in parallel and the model that achieves the lowest objective function value is chosen. These initial state sequences are generated by the K-means++ seeding technique introduced in (1) which has been shown to improve both the speed and accuracy of standard K-means clustering.

We note that the objective function in Step 5.1 of Algorithm 1 is convex in $\underline{\mu}$ and can be solved in closed-form. $\mu_k \in \mathbb{R}^p$ has the below optimal solution at the j^{th} iteration:

$$\boldsymbol{\mu}_{k}^{(j)} = \frac{\sum_{t=1}^{T} \boldsymbol{y}_{t} \mathbb{1}_{\{\boldsymbol{s}_{k}^{(j-1)} = k\}}}{\sum_{t=1}^{T} \mathbb{1}_{\{\boldsymbol{s}_{k}^{(j-1)} = k\}}}, \ k = 1, 2, \dots, K.$$

 $S^{(j)}$ in Step 5.2 is obtained using the following dynamic programming equations. Define

$$V(T,s) = \|\boldsymbol{y}_T - \boldsymbol{\mu}_s\|^2,$$
(2.1a)

$$V(t,s) = \|\boldsymbol{y}_t - \boldsymbol{\mu}_s\|^2 + \min_{i} \left\{ V(t+1,j) + \lambda \mathbf{1}_{\{s \neq j\}} \right\},$$
(2.1b)

for $t = T - 1, \dots, 2, 1$. The most likely state sequence is then given by

$$s_1 = \operatorname*{arg\,min}_{j} V\left(1, j\right),\tag{2.2a}$$

$$s_{t} = \arg\min_{j} \left\{ V(t, j) + \lambda \mathbb{1}_{\{s_{t-1} \neq j\}} \right\}, \ t = 2, \dots, T.$$
(2.2b)

2.2 Sparse Jump Model Calibration

The calibration of the Sparse Jump Model can be performed using an extension of the coordinate descent algorithm in Algorithm 1. Holding the feature weight vector w fixed, (1.3) is optimised in terms of μ and S. Secondly, holding μ and S fixed, (1.3) is optimised in terms of w.

As noted in (9), this iterative approach is not guaranteed to generate a global optimum because the
 problem is non-convex. Additionally, the first optimisation involves applying the fitting algorithm
 of the Standard Jump Model to a weighted version of the data, which by itself is not guaranteed to
 find a global optimum.

In Step 2(d), solving (1.3) with respect to w, while keeping μ and S fixed, can be done using soft-thresholding. The soft-thresholding operator S is given by $S(\mathbf{x}, c) = \operatorname{sgn}(\mathbf{x}) \odot (|\mathbf{x}| - c)_+$, where \mathbf{x}_+ denotes the positive part of the elements in \mathbf{x} and \odot denotes element-wise multiplication.

The solution to the convex problem (1.3), which follows from the Karush-Kuhn-Tucker conditions (more details can be found in (4)), is $\boldsymbol{w} = \frac{S(\boldsymbol{x},\Delta)}{\|S(\boldsymbol{x},\Delta)\|}$ where $\boldsymbol{x} = \sum_{k=1}^{K} |C_k| (\boldsymbol{\mu}_k - \bar{\boldsymbol{\mu}})^2$ is a vector comprising the between-cluster sum of squares (BCSS) contributions of each feature.

Here, $\Delta = 0$ if that results in $||w||_1 \le \kappa$; otherwise, we choose $\Delta > 0$ to yield $||w||_1 = \kappa$. This assumes that there is a unique maximal element of x and that $1 \le \kappa \le \sqrt{p}$ (13).

212 213

214

183

185

186

187

188 189

190 191

192 193

194

2.3 REGULARISED JUMP MODEL CALIBRATION

215 We propose calibrating the Regularised Jump Models using a coordinate descent algorithm similar to Algorithm 1 for Standard Jump Models. The coordinate descent algorithm alternates between

Algorithm 1 Calibration of Standard Jump Model in (10)

Input: Training dataset $Y = (y_1, y_2, ..., y_T)$, assumed number of states K, and jump penalty λ .

Step 1: Initialise state sequence $S^{(0)} := (s_1^{(0)}, s_2^{(0)}, \dots, s_T^{(0)}).$

Step 2: Iterate for j = 1, 2, ..., 10:

Fit model parameters $\mu^{(j)}$:

$$\underline{\boldsymbol{\mu}}^{(j)} \leftarrow \arg\min_{\underline{\boldsymbol{\mu}}} \sum_{t=1}^{T} \|\boldsymbol{y}_t - \boldsymbol{\mu}_{s_t^{(j-1)}}\|^2.$$
(5.1)

Fit state sequence $S^{(j)}$:

$$S^{(j)} \leftarrow \arg\min_{S} \left\{ \sum_{t=1}^{T-1} \left\{ \| \boldsymbol{y}_{t} - \boldsymbol{\mu}_{s_{t}}^{(j)} \|^{2} + \lambda \mathbf{1}_{\{s_{t} \neq s_{t+1}\}} \right\} + \| \boldsymbol{y}_{T} - \boldsymbol{\mu}_{s_{T}}^{(j)} \|^{2} \right\}.$$
(5.2)

Break if
$$S^{(j-1)} = S^{(j)}$$
.

Output: Estimated model parameters μ^* and state sequence S^* .

 $\alpha(i)$

minimising (1.6) with respect to the state sequence μ while keeping S fixed, and minimising (1.6) with respect to S while keeping μ fixed. These two steps are repeated for ten iterations at the most or until S does not change after two consecutive iterations.

Firstly, seven initial state sequences are generated using the strategy proposed in (11) which the authors show to be the most optimal for the Regularised K-Means model after benchmarking it against popular initialisation strategies. Given the feature selection aspect of the Regularised Jump Model, the initialisation strategy incorporates potential sparsity in the initial cluster centres. The initialisation strategy is given in Algorithm 3.

Each initial state sequence outputted from Algorithm 3 is used as an input in Algorithm 4. The algorithm is run in parallel using these seven initial cluster assignments and the run which ultimately yield the lowest objective function value, is chosen as the final result.

Once an initial state sequence has been generated, the model parameters μ are updated (Step 8.1). In particular, assuming a fixed state sequence $S = (s_1, s_2, \ldots, s_T)$, the following problem is solved:

$$\arg\min_{\underline{\mu}} \sum_{t=1}^{T} ||\boldsymbol{y}_t - \boldsymbol{\mu}_{s_t}||^2 + \gamma \mathcal{P}(\underline{\mu}).$$
(2.3)

The solution to (2.3) for each penalty function introduced in (1.5), is given in B.1. Once μ is updated, the state sequence S is updated.

Since Steps 8.2 and 5.2 are identical, the state sequence is updated using the same dynamic pro-gramming equations in (2.1) and (2.2).

HYPERPARAMETER TUNING

Table 1 shows the hyperparameters of each Jump Model.

The literature on the hyperparameter tuning of Jump Models is limited. (5) suggests tuning based on a grid search of potential hyperparmeter values and picking which combination produces the lowest Fang-Tang Information Criterion (FTIC) value, a criterion introduced in (6).

272 A

 Algorithm 2 Calibration of Sparse Jump Model in (9)

Input: Training dataset $Y = (y_1, y_2, ..., y_T)$, assumed number of states K and hyperparameters λ and κ .

Step 1: Initialise feature weights \boldsymbol{w} as $\boldsymbol{w}^{(0)} = \left(\frac{1}{\sqrt{p}}, \frac{1}{\sqrt{p}}, \dots, \frac{1}{\sqrt{p}}\right)$. Step 2: Iterate for i = 1, 2, ..., until $\|\boldsymbol{w}^{(i)} - \boldsymbol{w}^{(i-1)}\|_1 / \|\boldsymbol{w}^{(i-1)}\|_1 < 10^{-4}$: (a) Compute sequence of weighted features $\boldsymbol{z}_t = \boldsymbol{y}_t \odot \sqrt{\boldsymbol{w}^{(i-1)}}, \ t = 1, 2, \dots, T,$ where $\sqrt{\boldsymbol{w}^{(i)}}$ is the element-wise square root of $\boldsymbol{w}^{(i)}$ (b) Initialise state sequence $S^{(0)} := (s_1^{(0)}, s_2^{(0)}, \dots, s_T^{(0)}).$ (c) Iterate for j = 1, 2, ..., 10: i. Fit model parameters: $\underline{\boldsymbol{\mu}}^{(j)} \leftarrow \arg\min_{\underline{\boldsymbol{\mu}}} \sum_{t=1}^{T} \|\boldsymbol{z}_t - \boldsymbol{\mu}_{s_t^{(j-1)}}\|^2.$ ii. Fit state sequence: $S^{(j)} \leftarrow \arg\min_{S} \left\{ \sum_{t=1}^{T} \left\{ \| \boldsymbol{z}_{t} - \boldsymbol{\mu}_{s_{t}}^{(j)} \|^{2} + \lambda \mathbf{1}_{\{s_{t} \neq s_{t+1}\}} \right\} + \| \boldsymbol{z}_{T} - \boldsymbol{\mu}_{s_{T}}^{(j)} \|^{2} \right\}.$ Break if $S^{(j-1)} = S^{(j)}$. (d) Update weights $w^{(i)}$, while holding the model parameters $\mu^{(j)}$ and $S^{(j)}$ fixed, by solving (1.3) using soft thresholding. **Output**: Estimated model parameters μ^* , state sequence S^* and feature weights w^* . Algorithm 3 Initialisation of state sequences for calibration of Regularised Jump Model in (11) Step 1 Cluster $Y = (y_1, y_2, \dots, y_T)$ using standard K-means, obtaining a matrix of initial cluster centres $\boldsymbol{\mu} = (\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_K)$. **Step 2** Compute the Euclidean norm for each cluster centre $d_j = \|\boldsymbol{\mu}_{.,j}\|$ for j = 1, 2, ..., p and order them in descending order. Step 3 Execute K-means on the subset of variables corresponding to the 1, 2, 5, 10, 25, 50 and 100% largest d_i .

Output: Seven initial state sequences $(s_1^1, s_2^1, \dots, s_T^1), (s_1^2, s_2^2, \dots, s_T^2), \dots, (s_1^7, s_2^7, \dots, s_T^7).$

Algorithm 4 Coordinate descent algorithm for calibration of Regularised Jump Model

Input: Training dataset $Y = (y_1, y_2, \dots, y_T)$, assumed number of states K, penalty function \mathcal{P} , hyperparameters λ and γ , and initial state sequence $S^{(0)} := \left(s_1^{(0)}, s_2^{(0)}, \dots, s_T^{(0)}\right)$ generated using Algorithm 3.

Step 1: Iterate for k = 1, 2, ..., 10:

$$\underline{\boldsymbol{\mu}}^{(k)} \leftarrow \arg\min_{\underline{\boldsymbol{\mu}}} \sum_{t=1}^{T} || \boldsymbol{y}_t - \boldsymbol{\mu}_{s_t^{(k-1)}} ||^2 + \gamma \mathcal{P}\left(\underline{\boldsymbol{\mu}}\right),$$
(8.1)

$$S^{(k)} \leftarrow \arg\min_{S} \left\{ \sum_{t=1}^{T-1} \left\{ \| \boldsymbol{y}_{t} - \boldsymbol{\mu}_{s_{t}}^{(k)} \|^{2} + \lambda \mathbb{1}_{\{s_{t} \neq s_{t+1}\}} \right\} + \| \boldsymbol{y}_{T} - \boldsymbol{\mu}_{s_{T}}^{(k)} \|^{2} \right\}.$$
(8.2)

Break if $S^{(k)} = S^{(k-1)}$.

Output: Estimated matrix of cluster centres (or model parameters) μ^* and state sequence S^* .

Model	λ	κ	γ
Standard Jump	\checkmark		
Sparse Jump	\checkmark	\checkmark	
Regularised Jump	\checkmark		\checkmark

Table 1: Jump Model hyperparameters

In (5), the FTIC equation is an approximation and applies only to the Standard and Sparse Jump Models. Approximating information criteria such as the FTIC would require knowledge of the like-lihood function of the various Jump Models which, for the Regularised Jump Models, is unavailable.

Instead, we opt for a non-parametric criterion that is applicable to all Jump Models. We propose hyperparameter tuning based on a criterion pertaining to *clustering stability*. The key idea of clus-tering stability is that if we repeatedly draw samples from the same distribution as that of the original dataset, and apply the Jump Model calibration algorithms, a good algorithm should produce state sequences that are similar from one sample to another.

Following the notation in (12) to some extent, denote $Z = \{X_1, X_2, \dots, X_T\}$ a random sample of size T from some unknown distribution function $F: \mathbb{R}^p \to \mathbb{R}$. Let θ denotes the vector of active hyperparameters for a given Jump Model which can be referenced in Table 1. A clustering assign-ment $\psi(x)$ is defined as a mapping $\psi: \mathbb{R}^p \to \{1, 2, \dots, K\}$ and a given Jump Model generates a clustering assignment $\Psi(., \theta)$ when applied to a sample Z.

Definition 3.1 (Clustering Distance). The *clustering distance* between any two clustering assign-ments $\psi_1(x)$ and $\psi_2(x)$ is defined as

 $d(\psi_1, \psi_2) = \mathbb{P}\left(\{\psi_1(X) = \psi_1(Y)\} \cap \{\psi_2(X) = \psi_2(Y)\}\right),\$

where $X, Y \in \mathbb{R}^p$ are realisations sampled from F.

The distance between ψ_1 and ψ_2 measures the probability of their disagreement. The clustering instability of a given Jump Model is given in the following definition:

Definition 3.2 (Clustering Instability). The *clustering instability* of a clustering algorithm Ψ is given bv

- $S(\Psi; \boldsymbol{\theta}, T) = \mathbb{E}[d(\Psi(Z_1; \boldsymbol{\theta}), \Psi(Z_2; \boldsymbol{\theta}))].$ (3.1)
- where $\Psi(Z_1; \theta)$ and $\Psi(Z_2; \theta)$ are two clustering assignments obtained by applying $\Psi(.; \theta)$ to Z_1 and Z_2 which are two independent samples from F of size T.

Hyperparameter tuning is thus performed by finding a set of hyperparameter values that minimise
(3.1). Various methods of estimating (3.1) for hyperparameter tuning in clustering problems have
been proposed ((12), (7) and (3)).

We opt for a simple method of estimating (3.1) that is similar to the one proposed in (7); ten bootstrapped samples of the dataset Y are generated, each of which are size T. For a fixed combination of hyperparameter values and a given Jump Model, ten clustering assignments (or state sequences) are estimated once the model is fitted onto each bootstrapped sample. The number of times any two state sequences are not equal, are then counted and averaged across all $\binom{10}{2} = 45$ comparisons. This estimate can be expressed mathematically as

388 389

390 391

$$\hat{S}(\Psi; \boldsymbol{\theta}, T) = {\binom{10}{2}}^{-1} \sum_{i=1}^{10} \sum_{j=i+1}^{10} \Big(\sum_{t=1}^{T} \mathbb{1}_{\{\hat{s}_t^i \neq \hat{s}_t^j\}} \Big),$$
(3.2)

where \hat{s}_t^i is the estimated state at time t for the ith boostrapped sample.

392 393 394

395

408

409

4 SIMULATION STUDY

We conduct a simulation study to compare the accuracy of the Standard, Sparse and Regularised Jump Models, with respect to both state estimation and feature selection. In the simulation study, the true state sequence and array of relevant features are both known, which makes it possible to evaluate the ability of each model to correctly identify the state sequence and set of relevant features.

We adopt the same data generating process as in (9). In their simulation study, the authors test the Standard and Sparse Jump Models against several popular regime-switching models such as the Hidden Markov Model (HMM) and its various extensions, and determine that these two Jump Models deliver superior performance.

Instead of duplicating the simulation study by testing the same regime-switching models, we test the
 proposed Regularised Jump Models and compare their performances with the Standard and Sparse
 Jump Models.

4.1 RESULTS

Tables 2 and 3 compare the respective state estimation and feature selection performances of the Standard Jump Model (Standard), Sparse Jump Model (Sparse) and Regularised Jump Model with \mathcal{P}_k Penalty Function (Regularised \mathcal{P}_k) for $k \in \{0, 1, 2, 3\}$, for different values of μ and for different numbers of features P. The reported values of Tables 2 and 3 are the mean (and standard deviation) of the BAC of the estimated state sequence and of (3.2) respectively; the means and standard deviations are calculated over 100 simulations of T = 500 observations for each combination of μ and P.

Similar to (9), we use the Wilcoxon signed-rank test to determine whether the differences between the BACs in the Regularised Jump Models and those in the Sparse Jump Model are statistically significant. Bold entries in Tables 2 and 3 denote BACs of a Regularised Jump Model that are higher than that of the Sparse Jump Model with statistical significance $\alpha = 0.05$.

421 Comparing Standard and the remaining models in Table 2, it is evident that feature selection im-422 proves accuracy of state estimation when the number of features P is increased. In Table 2, the BAC 423 of the Standard Jump Model, a model that does not perform feature selection, generally decreases 424 when P is increased. On the other hand, the BAC of the Sparse and Regularised Jump Models stay 425 approximately level compared to the case of there being no irrelevant features (when P = 15).

The feature selection is more efficient for higher values of μ since for higher values of μ , the relevance of the first fifteen features becomes more apparent and hence easier for the model to identify. This logic is reinforced by the results in Table 3; the feature selection performance of the models improves for higher values of μ .

431 When comparing the relative performances of the Sparse, Regularised \mathcal{P}_0 , \mathcal{P}_1 and \mathcal{P}_3 Jump Models in Tables 2 and 3, we can see a correspondence between classification accuracy in terms of state estimation, and classification accuracy in terms of feature selection: models that are more accurate in terms of feature selection, are also more accurate in terms of state estimation. This further highlights the importance of feature selection in accurately estimating the true state sequence.

If we unbundle the BACs in Tables 2 and 3 in terms of μ , model and number of features P, and calculate the Spearman rank correlation coefficient between these two unbundled series of BAC values, we compute an estimate of 91% with a *p*-value of $\mathcal{O}(10^{-24})$. This leads us to reject the null hypothesis of there being no monotonic relationship between state estimation accuracy and feature selection accuracy.

Tables 2 and 3 demonstrate outperformance of the Regularised $\mathcal{P}_1, \mathcal{P}_2$ and \mathcal{P}_3 Models compared to both the Standard and Sparse Jump Models. The Regularised \mathcal{P}_0 model performs roughly in line with, or slight worse than, the Sparse Jump Model in terms of state estimation and feature selection.

In Table 2, the outperformance of the Regularised $\mathcal{P}_1, \mathcal{P}_2$ and \mathcal{P}_3 is most evident for the cases $\mu \ge 0.5$ and P < 300. For $\mu = 0.25$ and P = 300, all models produce roughly the same BACs that range between 0.334 and 0.344.

Similar results can be observed in Table 3. All models perform at approximately the same level for $\mu = 0.25$ and for $\mu \ge 0.5$, the Regularised \mathcal{P}_1 and \mathcal{P}_3 outperform both the Sparse and Regularised \mathcal{P}_0 models. The outperformance of the Regularised $\mathcal{P}_1, \mathcal{P}_2$ and \mathcal{P}_3 over the Sparse Jump Model is statistically significant for some cases, as shown by the bold entries in Tables 2 and 3.

451 452

453

468

469 470

471

472

473

474

475 476

477

478

479

480 481

482

483

484

485

5 CONCLUSION

In the new class of models called *Regularised Jump Models*, the regularised K-means model proposed in (11) has been adapted to the Jump Model Framework. These models perform joint state and parameter estimation, and feature selection. The feature selection process performed by the Regularised Jump Models is more direct and interpretable than that from the Sparse Jump Model, which performs feature selection by proxy. These models were tested in a simulation experiment and demonstrate evidence of outperformance over existing Jump Models.

One avenue of future research is adapting the asymptotic properties of the Regularised K-Means
model proven in (11) to the Regularised Jump Models. Some of these properties include consistency
and strong consistency in terms of the Hausdorff distance (see (11) Theorem 6 and the proof in
Appendix A. A.3).

 We've proposed a new hyperparameter tuning method for Jump Models, based on the idea of clustering stability. There is scope for future research on developing and testing alternative hyperparameter tuning methods for Jump Models.

REFERENCES

- [1] D. ARTHUR AND S. VASSILVITSKII, *k-means++: the advantages of careful seeding*, in Proceedings of the Eighteenth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA '07, USA, 2007, Society for Industrial and Applied Mathematics, p. 1027–1035.
- [2] A. BEMPORAD, V. BRESCHI, D. PIGA, AND S. P. BOYD, *Fitting jump models*, Automatica, 96 (2018), pp. 11–21.
- [3] A. BEN-HUR, A. ELISSEEFF, AND I. GUYON, A stability based method for discovering structure in clustered data, Pacific Symposium on Biocomputing. Pacific Symposium on Biocomputing, 2002 (2002), pp. 6–17.
- [4] S. BOYD AND L. VANDENBERGHE, Convex Optimization, Cambridge University Press, 2004.
- [5] F. CORTESE, P. KOLM, AND E. LINDSTROM, Generalized information criteria for highdimensional sparse statistical jump models, SSRN Electronic Journal, (2024).
- [6] Y. FAN AND C. Y. TANG, *Tuning parameter selection in high dimensional penalized likelihood*, Journal of the Royal Statistical Society. Series B (Statistical Methodology), 75 (2013), pp. 531–552.

- [7] J. HASLBECK AND D. WULFF, *Estimating the number of clusters via a corrected clustering instability*, Computational Statistics, 35 (2020).
 - [8] S. P. LLOYD, Least squares quantization in pcm, IEEE Trans. Inf. Theory, 28 (1982), pp. 129– 136.
 - [9] P. NYSTRUP, P. KOLM, AND E. LINDSTRÖM, Feature selection in jump models, Expert Systems with Applications, 184 (2021), p. 115558.
- [10] P. NYSTRUP, E. LINDSTRÖM, AND H. MADSEN, Learning hidden markov models with persistent states by penalizing jumps, Expert Systems with Applications, 150 (2020), p. 113307.
- [11] J. RAYMAEKERS AND R. H. ZAMAR, Regularized k-means through hard-thresholding, J. Mach. Learn. Res., 23 (2022).
- [12] W. SUN, J. WANG, AND Y. FANG, Regularized k-means clustering of high-dimensional data and its asymptotic consistency, Electronic Journal of Statistics, 6 (2012), pp. 148 167.
- [13] D. M. WITTEN AND R. TIBSHIRANI, A framework for feature selection in clustering, Journal of the American Statistical Association, 105 (2010), pp. 713–726. PMID: 20811510.
- [14] H. ZOU AND T. HASTIE, *Regularization and Variable Selection Via the Elastic Net*, Journal of the Royal Statistical Society Series B: Statistical Methodology, 67 (2005), pp. 301–320.

A SUPPLEMENTARY MATERIAL

B PARAMETER UPDATE EQUATIONS FOR REGULARISED JUMP MODEL CALIBRATION

Proposition B.1. Suppose that we have an assignment of the elements of Y into K clusters C_1, C_2, \ldots, C_K . Let $|C_k|$ denote the number of elements in cluster k. Furthermore, let M be a $T \times K$ cluster assignment matrix with elements $m_{t,k} = \mathbb{1}_{\{\boldsymbol{y}_t \in C_k\}}$.

Let $\underline{\mu}^*$ be the corresponding $K \times p$ matrix of cluster centres and $\underline{\mu}^*_{k,j}$ be the (k, j) element of $\underline{\mu}^*$ (k^{th} cluster centre along the j^{th} dimension). Keeping the cluster assignment matrix M fixed, solving (2.3) gives the following expressions for each penalty function $\mathcal{P}(\mu)$:

$$\mathcal{P}_{0}\left(\underline{\boldsymbol{\mu}}\right): \ \underline{\boldsymbol{\mu}}_{k,j} = \begin{cases} \underline{\boldsymbol{\mu}}_{k,j}^{*} & \text{if } \|\boldsymbol{y}_{.,j}\|^{2} > \|\boldsymbol{y}_{.,j} - \boldsymbol{M}\underline{\boldsymbol{\mu}}_{.,j}^{*}\| + T\gamma\\ 0 & \text{otherwise} \end{cases}$$
(B.1a)

$$\mathcal{P}_{1}\left(\underline{\mu}\right):\underline{\mu}_{k,j} = \max\left(0, 1 - \frac{T\gamma}{2|C_{k}||\underline{\mu}_{k,j}^{*}|}\right)\underline{\mu}_{k,j}^{*}$$
(B.1b)

$$\mathcal{P}_{2}\left(\underline{\mu}\right): \underline{\mu}_{k,j} = \frac{1}{1 + \frac{T\gamma}{|C_{k}|}} \underline{\mu}_{k,j}^{*}$$
(B.1c)

$$\mathcal{P}_{3}\left(\underline{\boldsymbol{\mu}}\right):\underline{\boldsymbol{\mu}}_{k,j}=\frac{1}{1+\frac{T\gamma}{2|C_{k}|||\underline{\boldsymbol{\mu}}_{...j}||}}\underline{\boldsymbol{\mu}}_{k,j}^{*}\quad\text{if }\underline{\boldsymbol{\mu}}_{...j}\neq\mathbf{0},\tag{B.1d}$$

where y_{ij} is the j^{th} column of Y.

Proof. Proof can be found in A.2. of the Appendix in (11).

Remark B.2. The update equations in (B.1) lend insight into the effects of each penalty function.

 \mathcal{P}_0 leads to *hard-thresholding*, the process of setting to zero the elements of an input vector whose absolute values are lower than an input threshold value; elements whose absolute values exceed

the threhold are left unchanged. The threshold value in (B.1a) is the sum of $T\gamma$ and the WCSS contribution from each feature. The interpretation of (B.1a) is that variables are included in the clustering if it sufficiently contributes to a decrease of the WCSS. \mathcal{P}_1 leads to soft-thresholding, whereby some elements are set to zero and the rest are shrunk towards zero. This resembles the solution of \mathcal{L}_1 -regularised (or Lasso) regression with orthonormal covariates. \mathcal{P}_2 is a ridge-type penalty that scales down each element of the array uniformly. \mathcal{P}_2 is the only penalty function that does not directly induce sparsity.

 \mathcal{P}_3 does not have an updating equation since the right-hand side of (B.1d) includes the \mathcal{L}_2 norm of $\|\boldsymbol{\mu}_{...i}\|$. The solution is thus found through an iterative algorithm.

	Р	15	30	60	150	300
μ	Model					
0.25	Standard	0.351	0.338	0.332	0.343	0.334
		(0.06)	(0.05)	(0.05)	(0.04)	(0.03
	Sparse	0.349	0.343	0.343	0.337	0.336
		(0.06)	(0.05)	(0.05)	(0.04)	(0.03
	Regularised \mathcal{P}_0	0.354	0.345	0.336	0.341	0.334
		(0.06)	(0.05)	(0.05)	(0.05)	(0.05
	Regularised \mathcal{P}_1	0.359	0.345	0.341	0.335	0.344
		(0.08)	(0.06)	(0.06)	(0.04)	(0.04
	Regularised P_2	0.343	(0.342)	(0.340)	0.335	0.344
	Deculorized D.	(0.07)	(0.00)	(0.03)	(0.05)	0.22/
	Regularised P3	(0.07)	(0.06)	(0.05)	(0.0542)	(0.04
		()	()	()	()	(
0.50	Standard	0.362	0.368	0.341	0.341	0.336
	~	(0.10)	(0.08)	(0.07)	(0.05)	(0.04
	Sparse	0.376	0.356	0.345	0.345	0.344
		(0.08)	(0.08)	(0.06)	(0.05)	(0.04
	Regularised \mathcal{P}_0	0.358	0.374	0.365	0.349	0.339
		(0.09)	(0.09)	(0.09)	(0.07)	(0.00
	Regularised P_1	(0.362)	(0.393)	(0.387)	(0.09)	0.349
	Deculorized D.	(0.11)	(0.15)	(0.12)	(0.09)	0.260
	Regularised P2	(0.14)	(0.14)	(0.13)	(0.07)	(0.07
	Regularised \mathcal{P}_{2}	0 383	0.382	0 374	0.358	0.354
	Regularised / 3	(0.14)	(0.14)	(0.13)	(0.08)	(0.06
0.75	Standard	0.385	0.401	0.366	0.353	0.354
0.75	Stanuaru	(0.15)	(0.13)	(0.300)	(0.06)	(0.06
	Sparse	0 380	0 363	0 364	0 347	0.33
	opuise	(0.12)	(0.11)	(0.09)	(0.07)	(0.06
	Regularised \mathcal{P}_0	0.418	0.386	0.389	0.387	0.408
		(0.14)	(0.14)	(0.14)	(0.13)	(0.14
	Regularised \mathcal{P}_1	0.418	0.406	0.432	0.394	0.458
		(0.18)	(0.18)	(0.18)	(0.17)	(0.18
	Regularised \mathcal{P}_2	0.400	0.420	0.447	0.399	0.447
		(0.18)	(0.16)	(0.18)	(0.18)	(0.18
	Regularised \mathcal{P}_3	0.407	0.435	0.391	0.403	0.468
		(0.18)	(0.19)	(0.17)	(0.15)	(0.17
1.00	Standard	0.448	0.423	0.391	0.357	0.347
		(0.20)	(0.15)	(0.13)	(0.09)	(0.08
	Sparse	0.408	0.393	0.363	0.374	0.347
		(0.15)	(0.14)	(0.12)	(0.11)	(0.08
	Regularised \mathcal{P}_0	0.423	0.409	0.464	0.411	0.494
		(0.19)	(0.18)	(0.19)	(0.16)	(0.20
	Regularised \mathcal{P}_1	0.483	0.446	0.524	0.499	0.56
		(0.22)	(0.24)	(0.25)	(0.24)	(0.22
	Regularised \mathcal{P}_2	0.495	0.465	0.523	0.462	0.498
		(0.23)	(0.24)	(0.23)	(0.23)	(0.23
	Regularised \mathcal{P}_3	0.444	0.447	0.507	0.492	0.59
		(11.77)	111 2/11	(11:7/1)	111 / / /	111.12

646

Table 2: Average BAC of state sequence (3 d.p.) for each Jump Model, across different values of μ and number of features *P*. Values in brackets are standard deviations (2 d.p.).

653						
654						
655						
656		Р	30	60	150	300
657		N 11				
658	μ	Model				
659	0.25	Sparse	0.499	0.502	0.502	0.506
660		I	(0.08)	(0.05)	(0.03)	(0.03)
661		Regularised \mathcal{P}_0	0.507	0.498	0.502	0.501
662			(0.04)	(0.02)	(0.01)	(0.02)
663		Regularised \mathcal{P}_1	0.511	0.505	0.496	0.504
664			(0.05)	(0.03)	(0.02)	(0.03)
665		Regularised \mathcal{P}_3	0.503	0.506	0.499	0.498
666			(0.05)	(0.03)	(0.03)	(0.04)
667	0.5	Sparse	0.509	0.507	0.509	0.505
668	0.5	opuise	(0.07)	(0.05)	(0.04)	(0.04)
669		Regularised \mathcal{P}_0	0.530	0.529	0.513	0.511
670		8	(0.05)	(0.04)	(0.03)	(0.03)
671		Regularised \mathcal{P}_1	0.560	0.573	0.535	0.509
672			(0.10)	(0.10)	(0.08)	(0.04)
673		Regularised \mathcal{P}_3	0.582	0.588	0.528	0.513
674			(0.10)	(0.11)	(0.06)	(0.05)
675	0.75	C	0.529	0.520	0.521	0.522
676	0.75	Sparse	0.538	(0.039)	(0.031)	0.522
677		Regularised \mathcal{D}_{α}	0.540	0.538	0 544	0.568
678		Regularised 7 ()	(0.07)	(0.05)	(0.07)	(0.05)
679		Regularised \mathcal{P}_1	0.594	0.608	0.621	0.650
680		0	(0.13)	(0.14)	(0.15)	(0.13)
681		Regularised \mathcal{P}_3	0.621	0.603	0.631	0.665
682			(0.15)	(0.13)	(0.16)	(0.14)
003		~	. .	0.600		0.7.0
004	1.00	Sparse	0.574	(0.12)	0.587	0.563
000		Deculorized D	(0.11)	(0.12)	(0.12)	(0.10)
000		Regularised P_0	(0.10)	(0.049	(0.048)	(0.024)
00/		Regularised \mathcal{D}_{i}	0.10)	0.660	0.695	0 743
000		r_1	(0.16)	(0.18)	(0.18)	(0.17)
009		Regularised \mathcal{P}_3	0.603	0.629	0.637	0.782
090			(0.14)	(0.15)	(0.16)	(0.17)
031						

Table 3: Average BAC of relevant features (3 d.p.) for each Jump Model, across different values of μ and number of features *P*. Values in brackets are standard deviations (2 d.p.).