# SUPPLEMENTARY MATERIALS: The Convex Mixture Distribution: Granger Causality for Categorical Time Series[*]

Alex Tank[†‡], Xiudi Li[†§], Emily B. Fox[¶], and Ali Shojaie[‖]

## SM1. Experiments.

### SM1.1. mLTD Bach Analysis.
For the mLTD Bach analysis, we performed a 5-fold cross validation to select the tuning parameter $\lambda$, then thresholded the final connection weights, given by the standardised $L_2$ norm of $\mathbf{Z}^{ij}$, at .01, as in the MTD case. First, we note that with only 5 total zero weights the final mLTD model is much less sparse than the MTD model. We display the final graph in Figure SM1, where, for interpretability, we bold edges with total weight greater than .45. In this graph there are strong connections in the counter-clockwise direction between G#, C#, F#, and B. However, the other connections on the circle of fifths are relatively weaker, and there are many more connections between notes far away on the circle of fifths. The mLTD graph also shows that the chord note both affects and is affected by many harmony notes. Furthermore, we see that the bass category is effected by most harmony notes as well. Overall, however, this graph is much less interpretable than the MTD graph and fails to find the full circle of fifths structure.

### SM1.2. iEEG Segmentation.
To segment the iEEG time series into a sequence of categorical states, we use a Markov switching autoregressive model. The model assumes that each channel in the $d$-dimensional EEG signal, $\mathbf{y}_t \in \mathbb{R}^d$, follows a Markov switching univariate autoregressive process (AR) each with the same $m$ dynamic regimes. Specifically, let $\mathbf{a}^1, \ldots, \mathbf{a}^m$, where $\mathbf{a}^i = (a_1^i, \ldots, a_h^i)$, denote the lag $h$ $\mathbf{AR(h)}$ parameters for each of the $m$ dynamic regimes and let $x_{jt}$ be the latent $m$-dimensional categorical state that governs the dynamics for channel $j$ at time $t$. The model assumes that $y_{jt}$ follows a locally stationary $\mathbf{AR(h)}$ model with $m$ state dynamics:

$$
\text{(SM1.1)} \qquad y_{jt} = \sum_{l=1}^{h} a_k^{x_{jt}} y_{j(t-l)} + e_{jt},
$$

where the lag $l$ AR dynamics at time $t$, $\mathbf{a}^{x_{jt}}$, are indexed by the latent state, $x_{jt}$, and $e_{jt}$ is mean zero Gaussian noise independent across series, $E(e_{jt}) = 0$ and $E(e_{jt}e_{j't'}) = 0$ for all $(j, t) \neq (j', t')$. The transitions between dynamic regimes are assumed to evolve independently

[‡]The Voleon Group, Berkeley, CA (alextank@uw.edu)

[§]Department of Biostatistics, University of Washington, Seattle WA (xiudil@uw.edu)

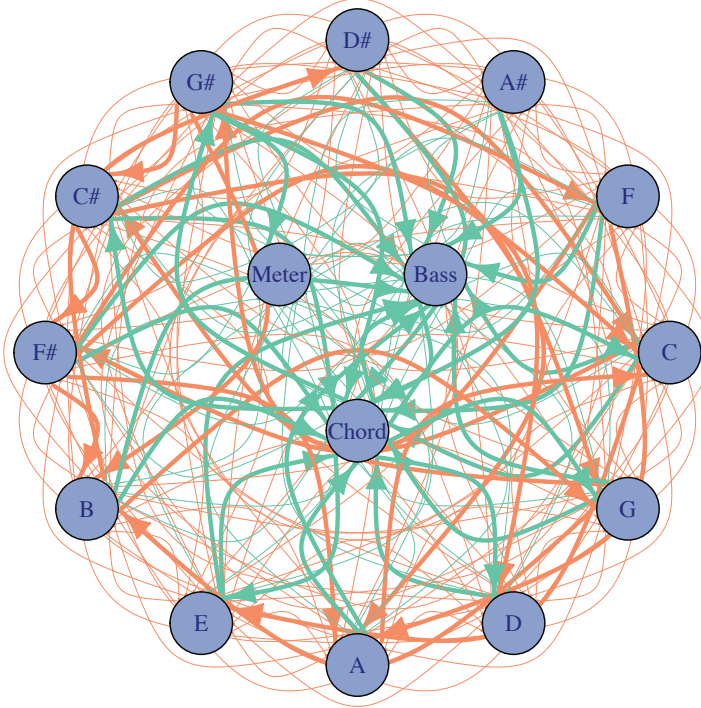[¶]Departments of Computer Science & Engineering and Statistics, University of Washington, Seattle WA (ebfox@uw.edu)

[‖]Department of Biostatistics, University of Washington, Seattle WA (ashojaie@uw.edu)

30  between series according to a hidden Markov model. See [SM11] for more details on the model.
31  Due to the long length of the series, we use a stochastic gradient MCMC algorithm [SM7] to
32  fit the model with $m = 5$ categorical states. We display the segmentation of a single channel
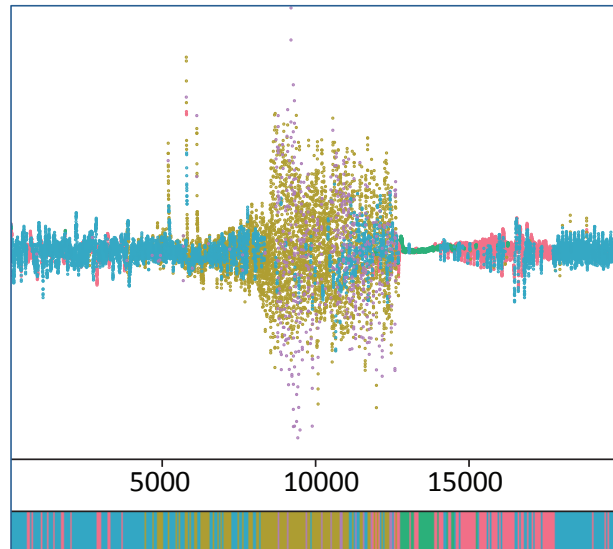    using this method in Figure SM2.

## mLTD Graph



**Figure SM1.** *The Granger causality graph for the 'Bach Choral Harmony' data set using the mLTD method. The harmony notes are displayed around the edge in a circle corresponding to the circle of fifths. Orange links display directed interactions between the harmony notes while green links display interactions to and from the 'bass', 'chord', and 'meter' variables.*
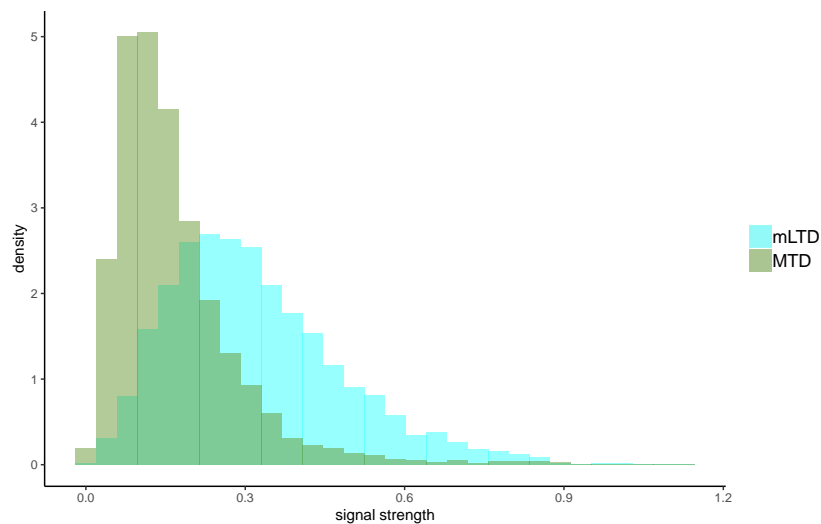
33

34  **SM1.3. Additional Simulation Results.** Figure SM3 compares the signal strengths in the
35  mLTD and MTD models for the case where each series has $m = 4$ possible states and $d = 15$.
36  To capture the effect of time series $j$ on time series $i$, we unfold the transition probability
37  tensor $p(x_{it}|x_{1(t-1)}, \ldots, x_{d(t-1)})$ along the mode defined by $x_{j(t-1)}$, and obtain an $m \times m^d$
38  matrix. We then compute the $l_2$ distances between any two rows of the resulting matrix. For
39  the MTD model, this is equivalent (up to scaling) to the $l_2$ distance between columns of $\mathbf{Z}^{ij}$,
40  since the effect is additive. We repeat this procedure for all $(i, j)$ pairs and aggregate the
41  results over 20 replications. Figure SM3 shows a histogram of nonzero signals in the MTD
42  and mLTD models.
43  We observe that, in our simulation settings, the difference among transition probabilities
44  in the mLTD model is larger than that in the MTD model, leading to stronger connections.
45  Next, we present median ROC curves over 20 replications for the proposed methods, under

**Figure SM2.** *Colored segmentation with $m = 5$ states of a single iEEG channel during a seizure using the Markov switching autoregressive model.*



**Figure SM3.** *Signal strengths in the mLTD and MTD models.*

46  different simulation settings. The results displayed in Figures SM4-SM5, Figures SM6-SM7
47  and Figures SM8-SM9, correspond to data generated by MTD, mLTD and latent VAR models,
48  respectively. We observe that for all three methods, the performance improves with increasing
49  sample size $T$ and worsens with increasing dimension $d$.
50     We also show the points on the ROC curves that correspond to tuning parameter values

51  chosen by BIC and cross-validation. In general, cross-validation tends to over-select Granger
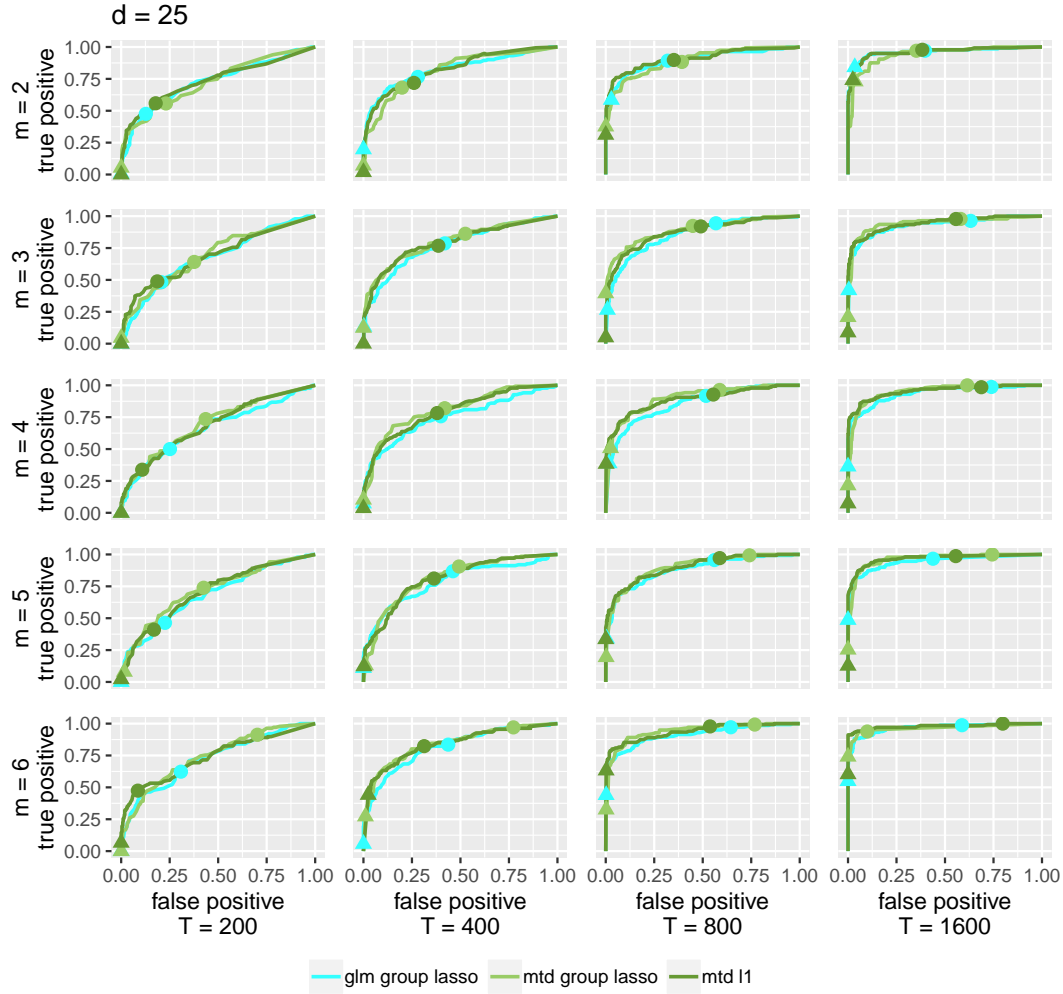52  causality relationships.  This highlights the importance of thresholding when using cross-
53  validation in practice. In contrast, BIC generally gives an overly sparse model when sample
54  size is small; but it performs much better with large sample sizes.



**Figure SM4.** *Median ROC curves over 20 simulation runs, for data generated by a sparse MTD process with d = 15.  Triangles correspond to tuning parameter values chosen by BIC; while dots correspond to the values chosen by cross-validation.*

55  Finally, in Figure SM10, we show the average run time of the three proposed methods
56  under different sample size $T$ and number of time series $d$, where each time series has 4
57  categories. We observe that in general mLTD group lasso runs faster than MTD with either
58  group lasso or lasso penalty. This is due to the constraints on the parameter set in the MTD
59  model, which requires additional projection steps. For all three methods, the run time scales
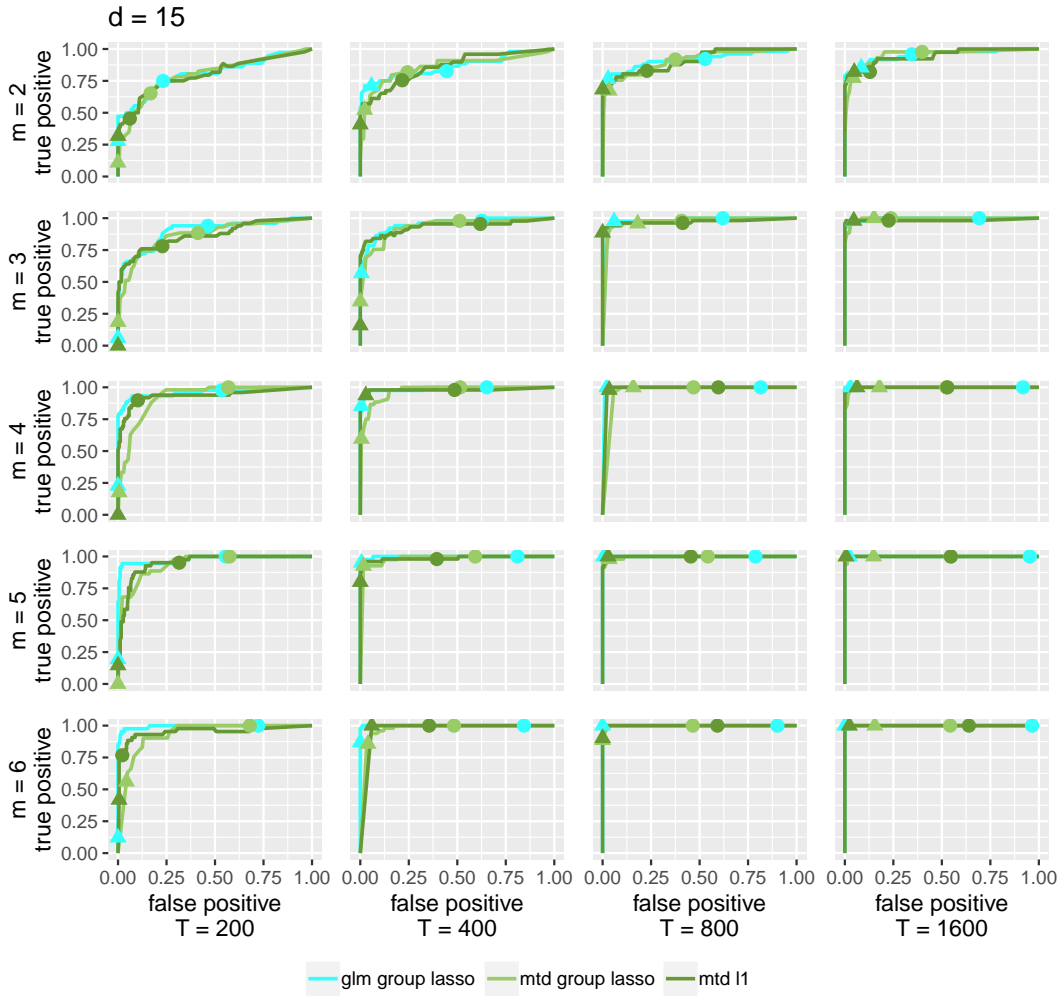
**Figure SM5.** *Median ROC curves over 20 simulation runs, for data generated by a sparse MTD process with $d = 25$. Triangles correspond to tuning parameter values chosen by BIC; while dots correspond to the values chosen by cross-validation.*

60  nearly linearly in sample size.

61  ## SM2. Proofs of Results in Section 3.

62  *Proof of Proposition 3.3.* If the columns of $\mathbf{Z}^j$ are all equal, then for all fixed values of
63  $x_{\setminus j(t-1)}$ the conditional distribution is the same for all values of $x_{j(t-1)}$. If one column is
64  different, then the conditional distribution for all values of $x_{\setminus j(t-1)}$ will depend on $x_{j(t-1)}$.

65  To prove the second claim, we let $\mathbf{Z}$ and $\tilde{\mathbf{Z}}$ be two parameterizations for the same MTD
66  model. Suppose that they give different causality conclusions. Then, there exists some $j \in$
67  $\{1, \ldots, d\}$ such that the columns of $\mathbf{Z}^j$ are all equal, while the columns of $\tilde{\mathbf{Z}}^j$ are not, or the
68  other way around. There must thus exist a row where at least two columns differ in this row.

**Figure SM6.** *Median ROC curves over 20 simulation runs, for data generated by a sparse mLTD process with d = 15. Triangles correspond to tuning parameter values chosen by BIC; while dots correspond to the values chosen by cross-validation.*

Without loss of generality, we assume that $\mathbf{Z}_{11}^1 \neq \mathbf{Z}_{12}^1$ but $\tilde{\mathbf{Z}}_{11}^1 = \tilde{\mathbf{Z}}_{12}^1$. Then under $\mathbf{Z}$, we have that

$$P\left(x_{it} = 1 | x_{1(t-1)} = 1, x_{2(t-1)}, \ldots, x_{d(t-1)}\right) \neq P\left(x_{it} = 1 | x_{1(t-1)} = 2, x_{2(t-1)}, \ldots, x_{d(t-1)}\right).$$

However, under $\tilde{\mathbf{Z}}$ we have that

$$P\left(x_{it} = 1 | x_{1(t-1)} = 1, x_{2(t-1)}, \ldots, x_{d(t-1)}\right) = P\left(x_{it} = 1 | x_{1(t-1)} = 2, x_{2(t-1)}, \ldots, x_{d(t-1)}\right).$$

This is a clear contradiction, as $\tilde{\mathbf{Z}}$ and $\mathbf{Z}$ are different parameterizations of the same model, and hence all conditional probabilities should be the same.                                          ∎

**Figure SM7.** *Median ROC curves over 20 simulation runs, for data generated by a sparse mLTD process with $d = 25$. Triangles correspond to tuning parameter values chosen by BIC; while dots correspond to the values chosen by cross-validation.*
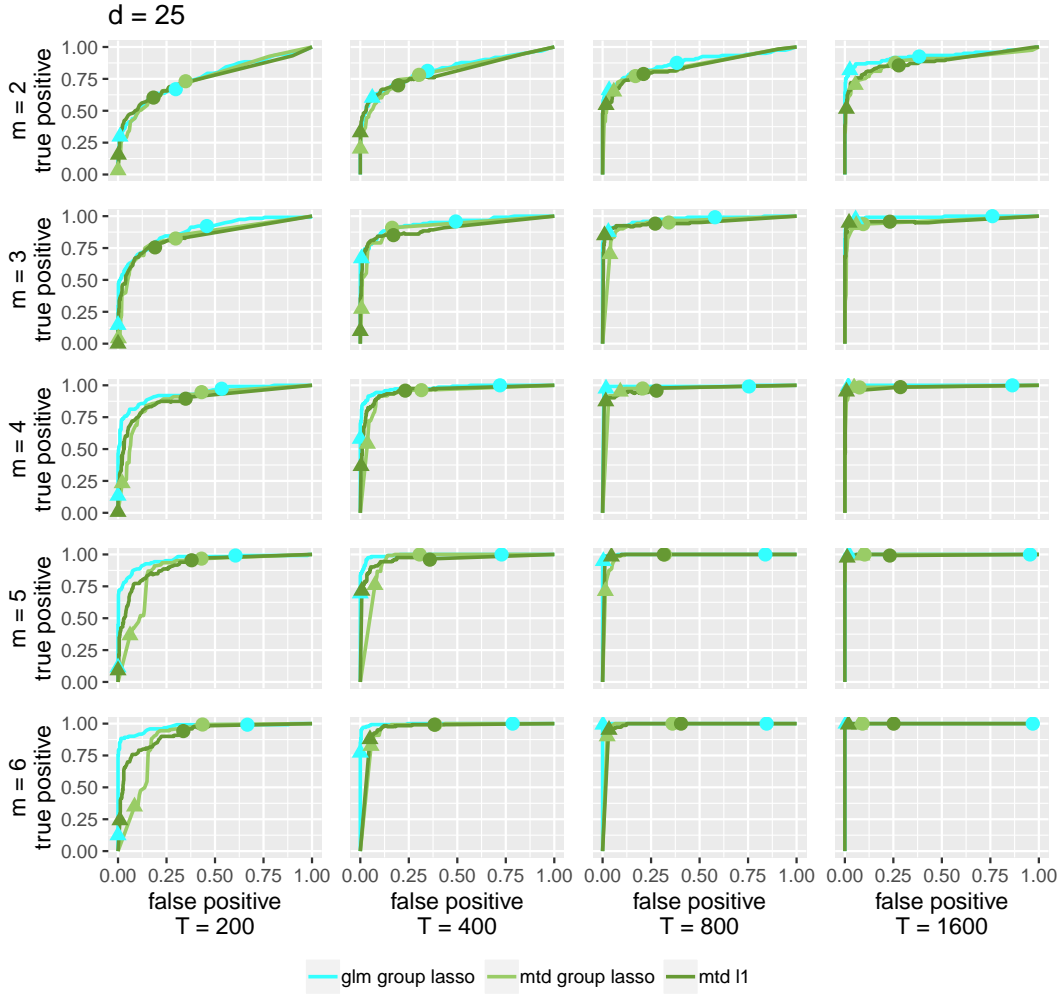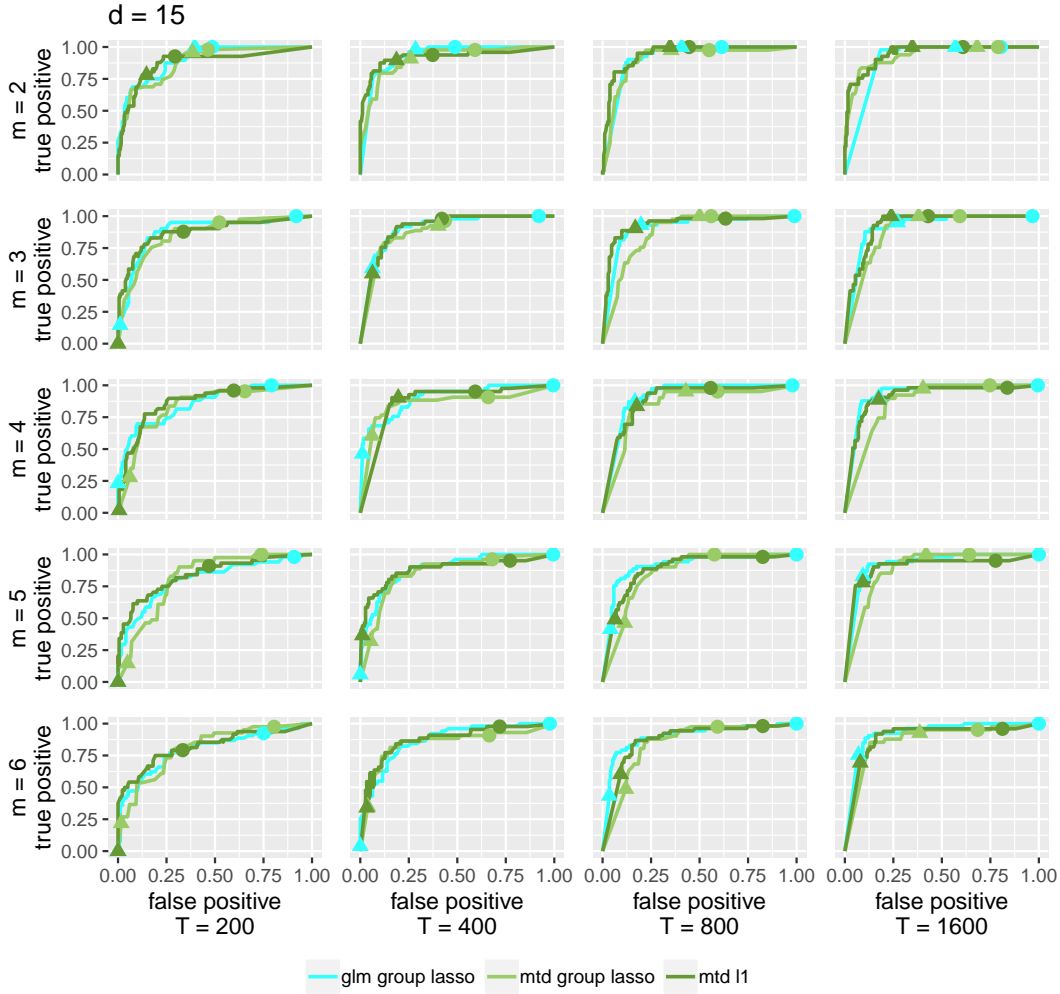
71  *Proof of Theorem 1.* First we show that any parameter set $\mathbf{Z}$ can be converted to another
72  set $\tilde{\mathbf{Z}}$ that contains at least one 0 element in each row of each matrix; and that $\tilde{\mathbf{Z}}$ satisfies the
73  constraints of the MTD model. Let $\mathbf{Z}$ be the parameter set for an MTD model. For each $\mathbf{Z}^j$
74  let the vector $\alpha^j$ be the minimal element in each row, $\alpha^j_k = \min \mathbf{Z}^j_{k:}$. Let $\tilde{\mathbf{Z}}^j = \mathbf{Z}^j - \alpha_j$ and
75  $\tilde{\mathbf{z}}^0 = \mathbf{z}^0 + \sum_{j=1}^d \alpha_j$. This $\tilde{\mathbf{Z}}$ gives the same MTD distribution as $\mathbf{Z}$. Furthermore, this $\tilde{\mathbf{Z}}$ has
76  a zero element in each row of each $\tilde{\mathbf{Z}}^j$ by construction.
77      The non-negativity constraint is trivially satisfied by $\tilde{\mathbf{Z}}$ as we subtract the minimum in each
78  row. For all $j$, we have that $\mathbf{1}^T \mathbf{Z}^j = \gamma_j \mathbf{1}^T$. Then $\mathbf{1}^T \tilde{\mathbf{Z}}^j = \mathbf{1}^T (\mathbf{Z}^j - \alpha^j \mathbf{1}^T) = (\gamma_j - \mathbf{1}^T \alpha^j) \mathbf{1}^T = $
79  $\tilde{\gamma}_j \mathbf{1}^T$, where we define $\tilde{\gamma}_j = \gamma_j - \mathbf{1}^T \alpha^j$. We note that $\tilde{\gamma}_j \geq 0$ as we subtract the row minimum.
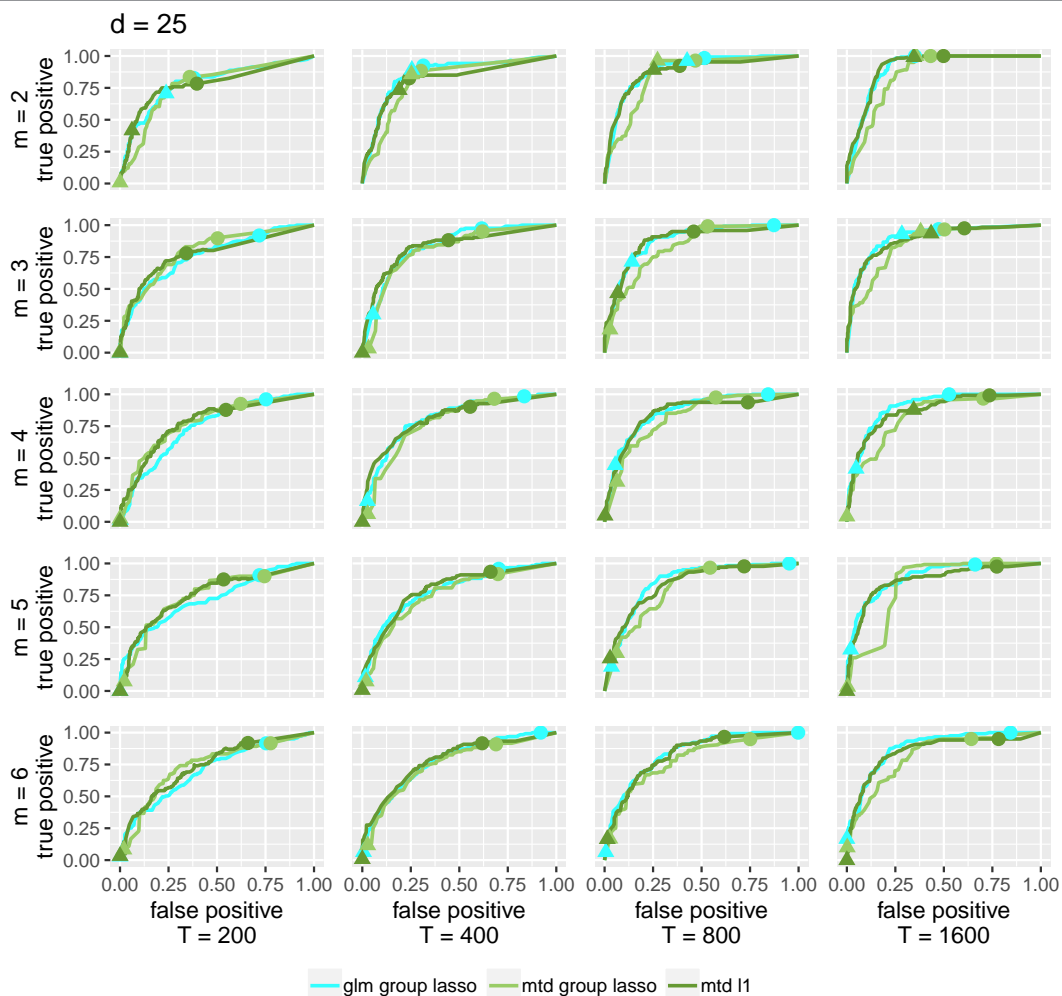
**Figure SM8.** *Median ROC curves over 20 simulation runs, for data generated by a sparse latent VAR process with $d = 15$. Triangles correspond to tuning parameter values chosen by BIC; while dots correspond to the values chosen by cross-validation.*

80  Hence within each $\tilde{\mathbf{Z}}^j$, the column sums are all equal. Finally, we have that $\tilde{\gamma}_0 = \gamma_0 +$
81  $\sum_{j=1}^{d} \mathbf{1}^T \alpha^j$ and $\sum_{j=0}^{d} \gamma_j = 1$, so $\sum_{j=0}^{d} \tilde{\gamma}_j = \gamma_0 + \sum_{j=1}^{d} \mathbf{1}^T \alpha^j + \sum_{j=1}^{d} \left(\gamma_j - \mathbf{1}^T \alpha^j\right) = \sum_{j=0}^{d} \gamma_j =$
82  1. Hence $\tilde{\gamma}_j$'s sum up to 1.
83      Next, we show that this new parameter set is uniquely determined. Suppose two parameter
84  sets $\mathbf{X}$ and $\mathbf{Y}$ provide the same MTD distribution. Let $\tilde{\mathbf{X}}$ be as above for $\mathbf{X}$ and $\tilde{\mathbf{Y}}$ of $\mathbf{Y}$.
85      We use a proof by contradiction. Suppose that $\tilde{\mathbf{Y}} \neq \tilde{\mathbf{X}}$. There must exist some $j$ and some
86  row $k$ such that $\tilde{\mathbf{X}}_{k:}^j \neq \tilde{\mathbf{Y}}_{k:}^j$. Let $l_X$ be the index of the zero element for $\mathbf{X}^j$, i.e., such that
87  $\tilde{\mathbf{X}}_{kl}^j = 0$, and likewise for $l_Y$. If there are more than one zero elements, pick any. Furthermore,
88  if $\tilde{\mathbf{X}}_{k:}^j$ and $\tilde{\mathbf{Y}}_{k:}^j$ share a zero in the same location (if there are one or more zero elements in

**Figure SM9.** *Median ROC curves over 20 simulation runs, for data generated by a sparse latent VAR process with $d = 25$. Triangles correspond to tuning parameter values chosen by BIC; while dots correspond to the values chosen by cross-validation.*

89    each), then let $l_X$ and $l_Y$ be that index so that $l_X = l_Y$.

90       If $l_X = l_Y$, let $l'$ be an index such that $\tilde{\mathbf{X}}^j_{kl'} \neq \tilde{\mathbf{Y}}^j_{kl'}$. This index must exist by construction.

91    Let the categories of other series (not for series $j$), $x_{\setminus j(t-1)}$, be fixed arbitrarily. The difference

**Figure SM10.** *Average run time of three proposed methods over 10 replications, with $m = 4$ and $\lambda = 100$ for MTD group lasso, MTD $L_1$ and $\lambda = 12.5$ for mLTD group lasso.*

92  between the conditional distributions for $\mathbf{X}$ are

$$
\tilde{\mathbf{X}}^j_{kl'} = \tilde{\mathbf{X}}^j_{kl'} - \tilde{\mathbf{X}}^j_{kl_X}
$$
$$
= \left( \tilde{\mathbf{X}}^j_{kl'} + \alpha_{jk} \right) - \left( \tilde{\mathbf{X}}^j_{kl_X} + \alpha_{jk} \right)
$$
93
$$
= \mathbf{X}^j_{kl'} - \mathbf{X}^j_{kl_X}
$$
$$
= \left( \mathbf{x}^0_k + \sum_{i \in \backslash j} \mathbf{X}^i_{kx_{i(t-1)}} + \mathbf{X}^j_{kl'} \right) - \left( \mathbf{x}^0_k + \sum_{i \in \backslash j} \mathbf{X}^i_{kx_{i(t-1)}} + \mathbf{X}^j_{kl_X} \right)
$$
94
$$
= p_X \left( x_t = k | x_{\backslash j(t-1)}, x_{j(t-1)} = l' \right) - p_X \left( x_t = k | x_{\backslash j(t-1)}, x_{j(t-1)} = l_X \right).
$$

95  A similar calculation for $\mathbf{Y}$ shows that

96  
97
$$
\tilde{\mathbf{Y}}^j_{kl'} = p_Y \left( x_t = k | x_{\backslash j(t-1)}, x_{j(t-1)} = l' \right) - p_Y \left( x_t = k | x_{\backslash j(t-1)}, x_{j(t-1)} = l_Y \right).
$$

98  However, $\tilde{\mathbf{Y}}^j_{kl'} \neq \tilde{\mathbf{X}}^j_{kl'}$, thus showing that

99  
100
$$
p_Y \left( x_t = k | x_{\backslash j(t-1)}, x_{j(t-1)} = l' \right) - p_Y \left( x_t = k | x_{\backslash j(t-1)}, x_{j(t-1)} = l_Y \right) \neq
$$
$$
p_X \left( x_t = k | x_{\backslash j(t-1)}, x_{j(t-1)} = l' \right) - p_X \left( x_t = k | x_{\backslash j(t-1)}, x_{j(t-1)} = l_X \right).
$$

101  This inequality contradicts our assumption that the MTD distributions parametrized by $\mathbf{X}$
102  and $\mathbf{Y}$ are the same since $l_X = l_Y$.
103      If $l_X \neq l_Y$, then

104  
105
$$
p_X \left( x_t = k | x_{\backslash j(t-1)}, x_{j(t-1)} = l_Y \right) - p_X \left( x_t = k | x_{\backslash j(t-1)}, x_{j(t-1)} = l_X \right) = \tilde{\mathbf{X}}^j_{kl_Y},
$$

106  and

107  
108
$$
p_Y \left( x_t = k | x_{\backslash j(t-1)}, x_{j(t-1)} = l_Y \right) - p_Y \left( x_t = k | x_{\backslash j(t-1)}, x_{j(t-1)} = l_X \right) = -\tilde{\mathbf{Y}}^j_{kl_X}.
$$

However, $-\tilde{\mathbf{Y}}^j_{kl_X} \neq \tilde{\mathbf{X}}^j_{kl_Y}$ since at least one of $\tilde{\mathbf{Y}}^j_{kl_X}$ and $\tilde{\mathbf{X}}^j_{kl_Y}$ are nonzero and both are nonnegative. Again, this shows that

$$p_Y\left(x_t = k|x_{\setminus j(t-1)}, x_{j(t-1)} = l_Y\right) - p_Y\left(x_t = k|x_{\setminus j(t-1)}, x_{j(t-1)} = l_X\right) \neq$$
$$p_X\left(x_t = k|x_{\setminus j(t-1)}, x_{j(t-1)} = l_Y\right) - p_X\left(x_t = k|x_{\setminus j(t-1)}, x_{j(t-1)} = l_X\right),$$

which contradicts our assumption that the MTD distributions parametrized by $\mathbf{X}$ and $\mathbf{Y}$ are the same.

The same argument shows that the reduction is unique. ∎

*Proof of Proposition 3.1.* First we check the parameter set satisfies the constraints of MTD model. Since $\mathbf{Z}$ and $\tilde{\mathbf{Z}}$ are valid MTD parameter sets, we have that $\forall j, \mathbf{1}^T\mathbf{Z}^j = \gamma_j\mathbf{1}^T, \mathbf{Z}^j \geq 0; \mathbf{1}^T\tilde{\mathbf{Z}}^j = \tilde{\gamma}_j\mathbf{1}^T, \tilde{\mathbf{Z}}^j \geq 0$, and $\mathbf{1}^T\gamma = 1, \gamma \geq 0; \mathbf{1}^T\tilde{\gamma} = 1, \tilde{\gamma} \geq 0$. Consider the new parameter set $\alpha\mathbf{Z} + (1-\alpha)\tilde{\mathbf{Z}}$; we have that for all $j$,

$$\mathbf{1}^T(\alpha\mathbf{Z}^j + (1-\alpha)\tilde{\mathbf{Z}}^j)$$
$$= \alpha(\mathbf{1}^T\mathbf{Z}^j) + (1-\alpha)(\mathbf{1}^T\tilde{\mathbf{Z}}^j)$$
$$= (\alpha\gamma_j + (1-\alpha)\tilde{\gamma}_j)\mathbf{1}^T$$
$$= \bar{\gamma}_j\mathbf{1}^T,$$

where we define $\bar{\gamma}_j = \alpha\gamma_j + (1-\alpha)\tilde{\gamma}_j$ for all $j$. Then

(SM2.1) $$\mathbf{1}^T\bar{\gamma} = \mathbf{1}^T(\alpha\gamma + (1-\alpha)\tilde{\gamma}) = \alpha + (1-\alpha) = 1.$$

Finally since $\mathbf{Z}^j, \tilde{\mathbf{Z}}^j, \gamma$ and $\tilde{\gamma}$ are all non-negative, we have that $\alpha\mathbf{Z}^j + (1-\alpha)\tilde{\mathbf{Z}}^j \geq 0 \ \forall j$ and $\bar{\gamma} \geq 0$.

Next we demonstrate that the probability tensor given by this new parameter set is the same as those given by $\mathbf{Z}$ and $\tilde{\mathbf{Z}}$. For any two MTD factorizations $\mathbf{Z}$ and $\tilde{\mathbf{Z}}$ that have the same conditional distribution $p\left(x_{kt}|x_{t-1}\right)$ for all $x_{kt}$ and $x_{t-1}$, then for any $0 < \alpha < 1$, the probability tensor of the MTD model for the parameter set $\alpha\mathbf{Z} + (1-\alpha)\tilde{\mathbf{Z}}$ is given by

$$\alpha\mathbf{z}^0_{x_{kt}} + (1-\alpha)\tilde{\mathbf{z}}^0_{x_{kt}} + \sum_{j=1}^d\left(\alpha\mathbf{Z}^j_{x_{kt}x_{j(t-1)}} + (1-\alpha)\tilde{\mathbf{Z}}^j_{x_{kt}x_{j(t-1)}}\right)$$

$$= \alpha\left(\mathbf{z}^0_{x_{kt}} + \sum_{j=1}^d\mathbf{Z}^j_{x_{kt}x_{j(t-1)}}\right) + (1-\alpha)\left(\tilde{\mathbf{z}}^0_{x_{kt}} + \sum_{i=1}^d\tilde{\mathbf{Z}}^j_{x_{kt}x_{j(t-1)}}\right)$$

$$= \alpha p\left(x_{kt}|x_{(t-1)}\right) + (1-\alpha)p\left(x_{kt}|x_{(t-1)}\right)$$

$$= p\left(x_{kt}|x_{(t-1)}\right).$$

This shows that $\alpha\mathbf{Z} + (1-\alpha)\tilde{\mathbf{Z}}$ has the same distribution as both $\mathbf{Z}$ and $\tilde{\mathbf{Z}}$, so that the set of parameters with the same distribution is a convex set. ∎

*Proof of Theorem 2.* First, we note that a solution always exists since the log likelihood $L(\mathbf{Z}) = -\sum_{t=1}^T \log\left(\mathbf{z}^0_{x_{it}} + \sum_{j=1}^d \mathbf{Z}^j_{x_{it}x_{j(t-1)}}\right)$ and penalty are both bounded below by zero and

143   the feasible set is closed and bounded. Suppose an optimal solution is $\mathbf{Z}$ for which there exists
144   some $j$ such that one row, call it $k$, of $\mathbf{Z}^j$ does not have a zero element. Let $\alpha = \min\left(\mathbf{Z}^j_{k:}\right)$
145   be the minimum value in row $k$ and let $\tilde{\mathbf{Z}}^j$ be equal to $\mathbf{Z}^j \; \forall j$ except that $\tilde{\mathbf{Z}}^j_{k:} = \mathbf{Z}^j_{k:} - \alpha$ and
146   $\tilde{z}^0_k = z^0_k + \alpha$. Due to the nonidentifiability of the MTD model $L(\tilde{\mathbf{Z}}) = L(\mathbf{Z})$, while we have
147   that $\Omega\left(\tilde{\mathbf{Z}}^j\right) < \Omega\left(\mathbf{Z}^j\right)$, implying for $\lambda > 0$

148
149
$$L(\tilde{\mathbf{Z}}) + \lambda\Omega(\tilde{\mathbf{Z}}) < L(\mathbf{Z}) + \lambda\Omega(\mathbf{Z}),$$

150   showing that $\mathbf{Z}$ cannot be an optima.                                                                   ∎

151       **SM3. Proof of Estimation Consistency.** First, we re-introduce some of our notations.
152   Recall that we define a covariate vector $W \in \mathbb{R}^{m+dm^2}$ as follows: $W_t = (W^T_{t0}, W^T_{t1}, \ldots, W^T_{td})^T$;
153   $W_{t0} = \left(W^1_{t0}, \ldots, W^m_{t0}\right)^T \in \mathbb{R}^m$ where $W^l_{t0} = I\{x_{it} = l\}$; and $W_{tj} = \left((W^1_{tj})^T, \ldots, (W^m_{tj})^T\right)^T \in$
154   $\mathbb{R}^{m^2}$, for $j \in \{1, \ldots, d\}$, where $W^l_{tj} = (W^{l1}_{tj}, \ldots, W^{lm}_{tj})^T$ and $W^{lk}_{tj} = I\left\{x_{it} = l, x_{j(t-1)} = k\right\}$.
155   Let $\mathcal{A}_t$ denote the sub $\sigma$-algebra generated by $x_1, \ldots, x_t$. Then $\{W_t\}$ is adapted to $\{\mathcal{A}_t\}$. For
156   a general MTD parameter set, we collect the parameters in a vector form $\beta \in \mathbb{R}^{m+dm^2}$ where
157   $\beta = \left(\beta^T_0, \beta^T_1, \ldots, \beta^T_d\right)^T$, $\beta_0 = \mathbf{z}^0$ and $\beta_j = \text{vec}(\mathbf{Z}^j)$ for $j \in \{1, \ldots, d\}$. The MTD model can be
158   written as

159   (SM3.1)                                    $$p(x_{it}|x_{t-1}) = W^T_t \beta.$$

160   For a general $\beta$, we define $R_n$ and $R$ to be the empirical and conditional expected negative
161   log-likelihood risks, respectively,

162   (SM3.2)          $$R_n(\beta) = -\frac{1}{T}\sum_{t=1}^{T}\log(W^T_t\beta); \quad R(\beta) = -\frac{1}{T}\sum_{t=1}^{T}\mathbb{E}\left[\log(W^T_t\beta)|\mathcal{A}_{t-1}\right].$$

163   Denote the group lasso penalty by $\Omega(\beta) = \sum_{j=1}^{d}\|\beta_j\|_2 = \sum_{j=1}^{d}\|\mathbf{Z}^j\|_F$. In the remainder of
164   this section, we will use the superscript 0 to denote the true parameter value.
165       We now turn to the proofs of the estimation consistency results.

166       **SM3.1. Proof of Lemma 6.2.** By definition, we have
167
168   (SM3.3)   $R_n(\beta) - R(\beta) - (R_n(\beta^0) - R(\beta^0))$

169                $$= -\frac{1}{T}\sum_{t=1}^{T}\left\{(\log(W^T_t\beta) - \log(W^T_t\beta^0)) - \mathbb{E}\left[(\log(W^T_t\beta) - \log(W^T_t\beta^0))|\mathcal{A}_{t-1}\right]\right\}.$$
170

171   For simplicity, we define $\tilde{\Omega}(\beta) = \|\beta_0\|_1 + \Omega(\beta)$. We will consider the following empirical process
172   indexed by $f$,

173   (SM3.4)                $$M_n(f) = \frac{1}{T}\sum_{t=1}^{T}(f(W_t) - \mathbb{E}[f(W_t)|\mathcal{A}_{t-1}]), \quad f \in \mathcal{F},$$

174  where the function class $\mathcal{F}$ is defined as

175  (SM3.5) $$\mathcal{F} = \left\{ f : f(W_t) = \log(W_t^T \beta) - \log(W_t^T \beta^0), \tilde{\Omega}(\beta - \beta^0) \leq M \right\}.$$

176  In the following, we will consider expectation of the supremum of this empirical process. Since
177  $W^T \beta^0$ is the transition probability, values of $W$ such that $W^T \beta^0 = 0$ will not contribute to
178  the expectation as these types of transition occur with probability 0.

179      Take $M_{\max} = c(T,d)/2$. If $\tilde{\Omega}(\beta - \beta^0) \leq M_{\max}$, $\left| W_t^T(\beta - \beta^0) \right| \leq M_{\max}$. Then by As-
180  sumption 2, we can regard $\mathcal{F}$ as a class of $[\log(c(T,d)/2), -\log(c(T,d)/2)]$-valued functions
181  for some function $c$ that only depends on the sample size $T$ and the number of time series
182  $d$. Hence we rescale it by multiplying $c(T,d)/2$, and denote the new class by $\tilde{\mathcal{F}}$ so that $\tilde{\mathcal{F}}$ is
183  bounded by 1 and is Lipschitz-continuous with Lipschitz constant 1.

184      We use the notion of sequential Rademacher complexity and covering number developed
185  in [SM9], which generalizes the definition of Rademacher complexity and covering number to
186  the setting of dependent samples. For a general function class $\mathcal{G}$ mapping from $\mathcal{Z}$ to $\mathbb{R}$, its
187  sequential Rademacher complexity is defined as

188  (SM3.6) $$\mathcal{R}_n = \sup_{\mathbf{z}} \mathcal{R}_n(\mathcal{G}, \mathbf{z}), \quad \text{where } \mathcal{R}_n(\mathcal{G}, \mathbf{z}) = \mathbb{E}\left[ \sup_{g \in \mathcal{G}} \frac{1}{T} \sum_{t=1}^{T} \epsilon_t g(\mathbf{z}_t(\epsilon)) \right],$$

189  where $(\epsilon_t)_{t=1}^T$ is a sequence of independent Rademacher random variables, i.e., Uniform $\{-1,1\}$
190  and $\mathbf{z}$ is a $\mathcal{Z}$-valued tree of depth $T$. Further, define
     (SM3.7)

191  $$\mathcal{D}_n(\mathcal{G}) = \sup_{\mathbf{z}} \mathcal{D}_n(\mathcal{G}, \mathbf{z}), \quad \text{where } \mathcal{D}_n(\mathcal{G}, \mathbf{z}) = \inf_{\alpha} \left\{ 4\alpha + 12/\sqrt{T} \int_{\alpha}^{1} \sqrt{\log \mathcal{N}_2(\delta, \mathcal{G}, \mathbf{z})} d\delta \right\},$$

192  and $\mathcal{N}_2(\cdot, \mathcal{G}, \mathbf{z})$ is the $l_2$ covering number of $\mathcal{G}$ over a tree $\mathbf{z}$ of depth $T$. See [SM9] for a
193  complete introduction to sequential Rademacher complexities and covering numbers.

194      By Theorem 2 and Theorem 4 in [SM9] we can bound the expectation by the sequential
195  Rademacher complexity and a Dudley-type entropy integral,

196  (SM3.8) $$\mathbb{E}\left[ \sup_{f \in \tilde{\mathcal{F}}} |M_n(f)| \right] = \mathbb{E}\left[ \sup_{f \in \tilde{\mathcal{F}} \cup -\tilde{\mathcal{F}}} M_n(f) \right] \leq 2\mathcal{R}_n(\tilde{\mathcal{F}} \cup -\tilde{\mathcal{F}}) \leq 2\mathcal{D}_n(\tilde{\mathcal{F}} \cup -\tilde{\mathcal{F}}).$$

197      We note that since $\beta^0$ is fixed, the covering number of $\tilde{\mathcal{F}}$ is the same as that of $\mathcal{G} = \{g(\cdot) :$
198  $g(W_t) = \log(W_t^T \beta), \tilde{\Omega}(\beta - \beta^0) \leq M\}$. Using the same arguments as in Lemma 13 of [SM9],
199  we can show that

200  (SM3.9) $$\log \mathcal{N}_2(\delta, \tilde{\mathcal{F}}, \mathbf{z}) = \log \mathcal{N}_2(\delta, \mathcal{G}, \mathbf{z}) \leq \log \mathcal{N}_\infty(\delta, \mathcal{H}, \mathbf{z}),$$

201  where $\mathcal{H} = \{h : h(W_t) = W_t^T \beta - W_t^T \beta^0, \tilde{\Omega}(\beta - \beta^0) \leq M\}$. Hence we have that

202  $$\mathcal{D}_n(\tilde{\mathcal{F}} \cup -\tilde{\mathcal{F}}) = \sup_{\mathbf{z}} \inf_{\alpha} \left\{ 4\alpha + 12/\sqrt{T} \int_{\alpha}^{1} \sqrt{\log \mathcal{N}_2\left(\delta, \tilde{\mathcal{F}} \cup -\tilde{\mathcal{F}}, \mathbf{z}\right)} d\delta \right\}$$

203  $$\leq \sup_{\mathbf{z}} \inf_{\alpha} \left\{ 4\alpha + 12/\sqrt{T} \int_{\alpha}^{1} \sqrt{\log \mathcal{N}_\infty\left(\delta, \mathcal{H} \cup -\mathcal{H}, \mathbf{z}\right)} d\delta \right\}$$

204
205  (SM3.10) $$= \mathcal{D}_n^\infty(\mathcal{H} \cup -\mathcal{H}).$$

Applying Lemma 9 in [SM9], we then get

(SM3.11) $$\mathcal{D}_n^\infty(\mathcal{H} \cup -\mathcal{H}) \leq 8\mathcal{R}_n(\mathcal{H} \cup -\mathcal{H})\left(1 + 4\sqrt{2}\log^{3/2}\left(eT^2\right)\right).$$

Our last step is to bound the Rademacher complexity of the class $\mathcal{H} \cup -\mathcal{H}$. Note that by definition,

$$\mathcal{R}_n(\mathcal{H} \cup -\mathcal{H}) = \sup_{\mathbf{w}} \mathbb{E}\left[\sup_{h \in \mathcal{H}}\left|\frac{1}{T}\sum_{t=1}^{T}\epsilon_t h(\mathbf{w}_t(\epsilon))\right|\right]$$

$$= \sup_{\mathbf{w}} \mathbb{E}\left[\sup_{\beta:\tilde{\Omega}(\beta-\beta^0)\leq M}\left|\frac{1}{T}\sum_{t=1}^{T}\epsilon_t \mathbf{w}_t(\epsilon)^T(\beta-\beta^0)\right|\right]$$

$$\leq \sup_{\mathbf{w}} \mathbb{E}\left[\left\|\frac{1}{T}\sum_{t=1}^{T}\epsilon_t \mathbf{w}_t(\epsilon)\right\|_\infty\right]\sup_{\beta:\tilde{\Omega}(\beta-\beta^0)\leq M}\left\|\beta-\beta^0\right\|_1$$

$$\leq mM \sup_{\mathbf{w}} \mathbb{E}\left[\max_{j\in\{1,\ldots,m+dm^2\}}\left|\frac{1}{T}\sum_{t=1}^{T}\epsilon_t \mathbf{w}_{tj}(\epsilon)\right|\right]$$

(SM3.12) $$\leq mM\sqrt{\frac{2\log(2(m+dm^2))}{T}}$$

where the fourth line follows from Lemma SM3.1 and the fifth line follows by applying the finite class lemma in the dependent setting [SM9] and a union bound.

Finally, combining (SM3.8), (SM3.10), (SM3.11) and (SM3.12), we have that

(SM3.13) $$\mathbb{E}\left[\sup_{f\in\mathcal{F}}|M_n(f)|\right] \leq \frac{32}{c(T,d)}mM\sqrt{\frac{2\log(2(m+dm^2))}{T}}\left(1+4\sqrt{2}\log^{3/2}\left(eT^2\right)\right).$$

Thus, by Markov inequality, we can take

(SM3.14) $$\lambda_\epsilon = O_p\left(\frac{1}{c(T,d)}\sqrt{\frac{\log(d)\log^3(T)}{T}}\right).$$

Finally, we need

(SM3.15) $$\frac{32\lambda_\epsilon(1+\delta)^2|S|}{\delta^2\phi^2(1/(1-\delta),S,\tau)} \leq \frac{1}{2}c(T,d),$$

which holds with probability tending to 1 by Assumption 2 and Assumption 3.

**SM3.2. Useful lemmas.** Before proving our main theorem, we first establish several lemmas which will be useful later in the proof.

The first lemma establishes a margin condition for the negative loglikelihood loss.

230  **Lemma SM3.1.** *(Margin condition) For all $\beta$ satisfying the MTD model constraints, $R(\beta) -$*
231  $R(\beta^0) \geq \frac{1}{2}\tilde{\tau}^2(\beta - \beta^0)$, *where $\tilde{\tau}(\beta)$ is a semi-norm defined as*

232  (SM3.16)
$$\tilde{\tau}(\beta) = \sqrt{\frac{1}{T}\beta^T \left(\sum_{t=1}^{T} \mathbb{E}\left[W_t W_t^T | \mathcal{A}_{t-1}\right]\right)\beta}.$$

233  *Proof.* As $\beta^0$ is the true parameter in the conditional distribution specified by MTD model,
234  it maximizes $\mathbb{E}[\log(W_t^T \beta)|\mathcal{A}_{t-1}]$ for all $t$, and hence minimizes $R(\beta)$. (The minimizer is not
235  unique, as in general the MTD model is not identifiable. But restricting each row to have at
236  least one zero can make the solution unique.)
237  Let $H(\beta) = 0$ denote the set of equality constraints on a valid MTD parameter set. Then,
238  consider the Lagrangian form of the MTD optimization,

239  (SM3.17)
$$R(\beta) + \lambda_1^T H(\beta) + \lambda_2^T(-\beta),$$

240  where $\lambda_1$ and $\lambda_2$ are the Lagrange multipliers associated with the equality and inequality
241  constraints respectively. Then $\beta^0$ satisfies the following KKT conditions:

242  (SM3.18)
$$\frac{\partial R(\beta)}{\partial \beta}\Big|_{\beta^0} + (\lambda_1^0)^T \frac{\partial H(\beta)}{\partial \beta}\Big|_{\beta^0} - \lambda_2^0 = 0;$$

243  (SM3.19)     $\qquad H(\beta^0) = 0;$

244  (SM3.20)     $\qquad (\lambda_2^0)^T \beta^0 = 0;$

245  (SM3.21)     $\qquad \lambda_2^0 \geq 0, \beta^0 \geq 0.$
246

247  We define a new function

248  (SM3.22)
$$\tilde{R}(\beta) = R(\beta) + (\lambda_1^0)^T H(\beta) + (\lambda_2^0)^T(-\beta).$$

249  Note that for all $\beta$ satisfying the MTD model constraints, $H(\beta) = 0$. Thus,

250
$$\tilde{R}(\beta) - \tilde{R}(\beta^0) = R(\beta) - R(\beta^0) + (\lambda_1^0)^T(H(\beta) - H(\beta^0)) + (\lambda_2^0)^T(\beta^0 - \beta)$$

251  (SM3.23)
$$= R(\beta) - R(\beta^0) + (\lambda_2^0)^T(\beta^0 - \beta)$$

252  (SM3.24)
$$= R(\beta) - R(\beta^0) - (\lambda_2^0)^T\beta,$$
253

254  where the last line follows from the KKT conditions. At the same time, using a first order
255  Taylor expansion and noting that the derivative of $\tilde{R}(\beta)$ at $\beta^0$ is 0, we get

256  (SM3.25)
$$\tilde{R}(\beta) - \tilde{R}(\beta^0) = (\beta - \beta^0)^T \frac{\partial^2 \tilde{R}}{\partial \beta^2}\Big|_{\beta^*}(\beta - \beta^0)/2,$$

257  for some $\beta^*$ between $\beta$ and $\beta^0$. Then, we have

258  (SM3.26)
$$R(\beta) - R(\beta^0) = (\lambda_2^0)^T\beta + (\beta - \beta^0)^T \frac{\partial^2 \tilde{R}}{\partial \beta^2}\Big|_{\beta^*}(\beta - \beta^0)/2.$$

259  Since the equality and inequality constraints are both linear, $\partial^2 \tilde{R}/\partial \beta^2 = \partial^2 R/\partial \beta^2$ and we
260  have

261  (SM3.27)
$$\frac{\partial^2 R}{\partial \beta^2} = \frac{1}{T}\sum_{t=1}^{T}\mathbb{E}\left[\frac{1}{(W_t^T\beta)^2}W_tW_t^T|\mathcal{A}_{t-1}\right].$$

262  Here, $W_t^T\beta$ models conditional probability, and is bounded between 0 and 1. Hence the above
263  expression is lower bounded by $\sum_{t=1}^{T}\mathbb{E}[W_tW_t^T|\mathcal{A}_{t-1}]/T$. Also, we have that $(\lambda_2^0)^T\beta \geq 0$.
264  Together, we have

265  (SM3.28)
$$R(\beta) - R(\beta^0) \geq \frac{1}{2}(\beta-\beta^0)^T\frac{\sum_{t=1}^{T}\mathbb{E}[W_tW_t^T|\mathcal{A}_{t-1}]}{T}(\beta-\beta^0). \qquad \blacksquare$$

266      Recall that $S$ denotes the active set of $\beta^0$, i.e., $S = \{j : j > 0, \ \beta_j^0 \neq \mathbf{0}\}$ and $S^c$ denotes its
267  complement in $\{1,\ldots,d\}$. We define $\Omega^+(\beta) = \sum_{j\in S}\|\beta_j\|_1$ and $\Omega^-(\beta) = \sum_{j\in S^c}\|\beta_j\|_1$. The
268  next lemma shows some basic properties of the penalty $\Omega(\cdot)$.

269      **Lemma SM3.2.** *(Properties of the penalty) The penalty $\Omega(\cdot)$ satisfies the following for any*
270  $\beta$*:*
271      1. $\|\beta\|_1 \leq \|\beta_0\|_1 + m\Omega(\beta)$.
272      2. $\Omega(\beta^0) - \Omega(\beta) \leq \Omega^+(\beta-\beta^0) - \Omega^-(\beta-\beta^0)$.

273      *Proof.*      1. $\|\beta\|_1 = \sum_{j=0}^{d}\|\beta_j\|_1$. For $j \neq 0, \beta_j \in \mathbb{R}^{m^2}$. By Lyapunov inequality
274      $\frac{1}{m^2}\|\beta_j\|_1 \leq \sqrt{\frac{1}{m^2}\|\beta_j\|_2^2}$, and hence $\|\beta_j\|_1 \leq m\|\beta_j\|_2$. Invoking the definition of $\Omega(\beta)$
275      completes the proof.
276      2. We note that $\Omega(\beta) = \Omega^+(\beta)+\Omega^-(\beta)$. By the triangle inequality, $\|\beta_j^0\|_1 \leq \|\beta_j^0 - \beta_j\|_1 +$
277      $\|\beta_j\|_1$. Summing over $j \in S$ we have $\Omega^+(\beta^0) - \Omega^+(\beta) \leq \Omega^+(\beta^0 - \beta)$. By definition
278      $\Omega^-(\beta^0) = 0$ and $\beta_j - \beta_j^0 = \beta_j$ for $j \in S^c$, which implies that $\Omega^-(\beta - \beta^0) = \Omega^-(\beta)$.
279      Thus,

280
$$\Omega(\beta^0) - \Omega(\beta) = \Omega^+(\beta^0) - \Omega^+(\beta) - \Omega^-(\beta)$$
281  (SM3.29)
282
$$\leq \Omega^+(\beta-\beta^0) - \Omega^-(\beta) = \Omega^+(\beta-\beta^0) - \Omega^-(\beta-\beta^0). \qquad \blacksquare$$

283      Recall that we have defined a semi-norm $\tilde{\tau}(\beta) = \sqrt{\beta^T\sum_{t=1}^{T}\mathbb{E}[W_tW_t^T|\mathcal{A}_{t-1}]\beta/T}$. However,
284  this semi-norm itself is random as we condition on the past. The next lemma shows that it is
285  close to a deterministic semi-norm $\tau(\cdot)$, and the compatibility constants defined with $\tilde{\tau}$ and $\tau$
286  are close. To this end, we will use concentration inequalities for Markov chains developed in
287  [SM8].

288      **Lemma SM3.3.** *Under Assumption 1 and Assumption 4, with probability at least $1 - 1/T$,*

289  (SM3.30)   $\dfrac{\phi^2(L,S,\tilde{\tau})}{\phi^2(L,S,\tau)} \geq 1 - (1+(1+L)m)^2 C'\sqrt{\dfrac{\log(2(m+dm^2)^2) + \log(T)}{T\gamma_{ps}}}|S|/\phi^2(L,S,\tau).$

290  *Thus, under Assumptions 1, 3 and 4, for $T$ sufficiently large, $\phi^2(L,S,\tilde{\tau})/\phi^2(L,S,\tau) > 1/2$*
291  *with probability at least $1 - 1/T$.*

*Proof.* For any $j, k \in \{1, \ldots, m + dm^2\}$, $W_j W_k$ is bounded between 0 and 1. For simplicity, we will assume, for now, that $x_0 \sim \pi$, i.e., the chain starts in the stationary distribution. We will relax this assumption later. Applying Theorem 3.11 in [SM8],

(SM3.31)
$$\mathbb{P}\left(\left|\frac{1}{T}\sum_{t=1}^{T}\mathbb{E}[W_{tj}W_{tk}|\mathcal{A}_{t-1}] - \mathbb{E}_{\pi}[W_{1j}W_{1k}]\right| \geq t\right) \leq 2\exp\left(-\frac{T^2t^2\gamma_{ps}}{8(T + 1/\gamma_{ps}) + 20Tt}\right).$$

And, using a union bound,

(SM3.32) $\quad \mathbb{P}\left(\sup_{j,k}\left|\frac{1}{T}\sum_{t=1}^{T}\mathbb{E}[W_{tj}W_{tk}|\mathcal{A}_{t-1}] - \mathbb{E}[W_{1j}W_{1k}]\right| \geq t\right) \leq$

$$2(m + dm^2)^2\exp\left(-\frac{T^2t^2\gamma_{ps}}{8(T + 1/\gamma_{ps}) + 20Tt}\right).$$

In order to obtain a concentration bound, we will choose $t = o(1)$ and consider large $T$. Hence, the right-hand-side is of the same order as $2(m + dm^2)^2\exp(-CTt^2\gamma_{ps})$, provided that $1/\gamma_{ps} = o(T)$. Now setting $t = \sqrt{\log(2(m + dm^2)^2/\alpha)/CT\gamma_{ps}}$,

(SM3.33) $\quad \mathbb{P}\left(\max_{j,k}\left|\frac{1}{T}\sum_{t=1}^{T}\mathbb{E}[W_{tj}W_{tk}|\mathcal{A}_{t-1}] - \mathbb{E}[W_{1j}W_{1k}]\right| \geq \sqrt{\frac{\log(2(m + dm^2)^2/\alpha)}{CT\gamma_{ps}}}\right) \leq \alpha,$

for $T$ sufficiently large.

Then, for all $\beta$

$$\left|\tau^2(\beta) - \tilde{\tau}^2(\beta)\right| = \left|\beta^T\left(\frac{1}{T}\sum_{t=1}^{T}\mathbb{E}[W_t W_t^T|\mathcal{A}_{t-1}] - \mathbb{E}_{\pi}[W_1 W_1^T]\right)\beta\right|$$

$$\leq \|\beta\|_1^2 \left\|\frac{1}{T}\sum_{t=1}^{T}\mathbb{E}[W_t W_t^T|\mathcal{A}_{t-1}] - \mathbb{E}_{\pi}[W_1 W_1^T]\right\|_{\infty}$$

(SM3.34)
$$\leq \|\beta\|_1^2 C'\sqrt{\frac{\log(2(m + dm^2)^2/\alpha)}{T\gamma_{ps}}},$$

where by (SM3.33) the last line holds with probability at least $1 - \alpha$.

Recall the definition of $\Gamma$ and compatibility constant $\phi$,

(SM3.35) $\qquad \Gamma_{\Omega}(L, S, \tau) = \left(\min\left\{\tau(\beta) : \|\beta_0\|_1 + \Omega^+(\beta) = 1, \Omega^-(\beta) \leq L\right\}\right)^{-1}$

(SM3.36) $\qquad \phi^2(L, S, \tau) = \Gamma_{\Omega}^{-2}(L, S, \tau)|S|.$

Thus,

$$\frac{\phi^2(L, S, \tilde{\tau})}{\phi^2(L, S, \tau)} = \frac{\Gamma_{\Omega}^2(L, S, \tau)}{\Gamma_{\Omega}^2(L, S, \tilde{\tau})} = \frac{\min \tilde{\tau}^2(\beta)}{\min \tau^2(\beta)} \geq 1 + \frac{\min \tilde{\tau}^2(\beta) - \tau^2(\beta)}{\min \tau^2(\beta)}$$

(SM3.37) $\qquad \geq 1 - (1 + (1 + L)m)^2 C'\sqrt{\frac{\log(2(m + dm^2)^2/\alpha)}{T\gamma_{ps}}}|S|/\phi^2(L, S, \tau),$

with probability at least $1 - \alpha$. Setting $\alpha = 1/T$, we see that with probability approaching 1, the ratio is greater than $\frac{1}{2}$ for sufficiently large $T$, provided that $|S|\sqrt{\log(d)/T}\gamma_{ps} = o(1)$ and $\phi^2(L, S, \tau)$ is bounded away from 0.

If the chain does not start in stationary distribution, a result similar to (SM3.31) can be established, provided that the distribution of $x_0$ is not too far away from $\pi$. In the rest of this subsection, we use $\mathbb{P}_q$ to denote the probability under the case $x_0 \sim q$. Define

(SM3.38)
$$N_q = \begin{cases} \mathbb{E}_\pi \left[ \left( \frac{q(x)}{\pi(x)} \right)^2 \right] & \text{if } q \text{ is absolutely continuous with respect to } \pi, \\ +\infty & \text{otherwise.} \end{cases}$$

Applying Proposition 3.15 in [SM8], we get

$$\mathbb{P}_q \left( \left| \frac{1}{T} \sum_{t=1}^T \mathbb{E}[W_{tj}W_{tk}|\mathcal{A}_{t-1}] - \mathbb{E}_\pi[W_{1j}W_{1k}] \right| \geq t \right)$$

$$\leq N_q^{1/2} \left[ \mathbb{P} \left( \left| \frac{1}{T} \sum_{t=1}^T \mathbb{E}[W_{tj}W_{tk}|\mathcal{A}_{t-1}] - \mathbb{E}_\pi[W_{1j}W_{1k}] \right| \geq t \right) \right]^{1/2}$$

(SM3.39)
$$\leq 2N_q^{1/2} \exp \left( - \frac{T^2 t^2 \gamma_{ps}}{16(T + 1/\gamma_{ps}) + 40Tt} \right).$$

This bound is essentially the same as in (SM3.31), except that we are working with different constants. The rest of the proof follows. ∎

**SM3.3. Proof of Theorem 6.1.** Next we prove our main theorem, which is a modification of the proof of Theorem 7.2 in [SM10]. The difference is that we handle the unpenalized intercept as in [SM2] and we have time dependence in the data. For notational convenience, define
$$M = \frac{4\lambda(1 + \delta)^2 |S|}{\delta \phi^2(1/(1 - \delta), S, \tau)}, \text{ and } t = \frac{M}{M + \Omega(\hat{\beta} - \beta^0) + \|\hat{\beta}_0 - \beta_0^0\|_1}.$$

Define $\tilde{\beta} = t\hat{\beta} + (1 - t)\beta^0$. With this construction, $\|\tilde{\beta}_0 - \beta_0^0\|_1 + \Omega(\tilde{\beta} - \beta^0) \leq M$.

We note that although in general $\tilde{\beta}$ may not have a zero in each row of the corresponding $\mathbf{Z}^j$ matrices, and hence may not be identifiable, it does satisfy the equality and inequality constraints of the MTD model. By the convexity of $R_n + \lambda\Omega$, we have that

$$R_n(\tilde{\beta}) + \lambda\Omega(\tilde{\beta}) \leq tR_n(\hat{\beta}) + t\lambda\Omega(\hat{\beta}) + (1 - t)R_n(\beta^0) + (1 - t)\lambda\Omega(\beta^0)$$

(SM3.40)
$$\leq R_n(\beta^0) + \lambda\Omega(\beta^0).$$

We rewrite this and apply Lemma 6.2 and Lemma SM3.2,

$$0 \leq R(\tilde{\beta}) - R(\beta^0) \leq - \left[ [R_n(\tilde{\beta}) - R(\tilde{\beta})] - [R_n(\beta^0) - R(\beta^0)] \right] + \lambda\Omega(\beta^0) - \lambda\Omega(\tilde{\beta})$$

$$\leq \lambda_\epsilon M + \lambda\Omega(\beta^0) - \lambda\Omega(\tilde{\beta})$$

(SM3.41)
$$\leq \lambda_\epsilon M + \lambda\Omega^+(\tilde{\beta} - \beta^0) - \lambda\Omega^-(\tilde{\beta} - \beta^0).$$

We consider two cases.

- Case 1: If $\lambda\|\tilde{\beta}_0 - \beta_0^0\|_1 + \lambda\Omega^+(\tilde{\beta} - \beta^0) \leq (1-\delta)\lambda_\epsilon M/\delta$, we have that

$$(\text{SM3.42}) \qquad \delta\lambda\left(\|\tilde{\beta}_0 - \beta_0^0\|_1 + \Omega^+(\tilde{\beta} - \beta^0)\right) \leq \lambda_\epsilon M,$$

and

$$(\text{SM3.43}) \qquad \delta\lambda\Omega^-(\tilde{\beta} - \beta^0) \leq \lambda_\epsilon M.$$

Hence,

$$(\text{SM3.44}) \qquad \delta\lambda\left(\|\tilde{\beta}_0 - \beta_0^0\|_1 + \Omega(\tilde{\beta} - \beta^0)\right) \leq 2\lambda_\epsilon M.$$

- Case 2: If instead $\lambda\|\tilde{\beta}_0 - \beta_0^0\|_1 + \lambda\Omega^+(\tilde{\beta} - \beta^0) \geq (1-\delta)\lambda_\epsilon M/\delta$, then by (SM3.41)

$$R(\tilde{\beta}) - R(\beta^0) + \lambda\Omega^-(\tilde{\beta} - \beta^0) \leq \lambda\Omega^+(\tilde{\beta} - \beta^0) + \frac{\delta}{(1-\delta)}\lambda\left(\Omega^+(\tilde{\beta} - \beta^0) + \|\tilde{\beta}_0 - \beta_0^0\|_1\right)$$

$$(\text{SM3.45}) \qquad \qquad \leq \lambda\left(\Omega^+(\tilde{\beta} - \beta^0) + \|\tilde{\beta}_0 - \beta_0^0\|_1\right)/(1-\delta),$$

where the second inequality holds because $0 < \delta < 1$. Since $R(\tilde{\beta}) - R(\beta^0) \geq 0$,

$$(\text{SM3.46}) \qquad \Omega^-(\tilde{\beta} - \beta^0) \leq \left(\Omega^+(\tilde{\beta} - \beta^0) + \|\tilde{\beta}_0 - \beta_0^0\|_1\right)/(1-\delta),$$

which allows us to use the compatibility condition later. Again from (SM3.41),

$$R(\tilde{\beta}) - R(\beta) + \lambda\Omega^-(\tilde{\beta} - \beta^0) + \delta\lambda\left(\Omega^+(\tilde{\beta} - \beta^0) + \|\tilde{\beta}_0 - \beta_0^0\|_1\right)$$

$$\leq \lambda_\epsilon M + (1+\delta)\lambda\left(\Omega^+(\tilde{\beta} - \beta^0) + \|\tilde{\beta}_0 - \beta_0^0\|_1\right)$$

$$\leq \lambda(1+\delta)\tilde{\tau}(\tilde{\beta} - \beta^0)\Gamma_\Omega(1/(1-\delta), S, \tilde{\tau}) + \lambda_\epsilon M$$

$$\leq \frac{1}{2}(\lambda^2(1+\delta)^2\Gamma_\Omega^2(1/(1-\delta), S, \tilde{\tau})) + \frac{1}{2}\tilde{\tau}^2(\tilde{\beta} - \beta^0) + \lambda_\epsilon M$$

$$(\text{SM3.47}) \qquad \leq \frac{1}{2}(\lambda^2(1+\delta)^2\Gamma_\Omega^2(1/(1-\delta), S, \tilde{\tau})) + R(\tilde{\beta}) - R(\beta) + \lambda_\epsilon M,$$

where the second inequality follows by applying Assumption 3 with stretching factor $1/(1-\delta)$, and the fourth inequality follows from Lemma SM3.1. It follows that

$$\delta\lambda\left(\Omega(\tilde{\beta} - \beta^0) + \|\tilde{\beta}_0 - \beta_0^0\|_1\right) \leq \frac{1}{2}(\lambda(1+\delta)\Gamma_\Omega(1/(1-\delta), S, \tilde{\tau}))^2 + \lambda_\epsilon M$$

$$= \frac{1}{2}(\lambda(1+\delta))^2\frac{|S|}{\phi^2(L, S, \tilde{\tau})} + \lambda_\epsilon M$$

$$(\text{SM3.48}) \qquad \leq \frac{(\lambda(1+\delta))^2|S|}{\phi^2(L, S, \tau)} + \lambda_\epsilon M,$$

with probability approaching 1.

374      Hence, in both cases we have that with probability going to 1,

$$\delta\lambda\left(\Omega(\tilde{\beta}-\beta^0)+\|\tilde{\beta}_0-\beta_0^0\|_1\right)\leq 2\lambda_\epsilon M+(\lambda(1+\delta)\Gamma_\Omega(1/(1-\delta),S,\tau))^2$$

(SM3.49)                                        $$= \delta\lambda M/4+2\lambda_\epsilon M\leq\delta\lambda M/2,$$

where the inequality follows from the fact that $\lambda\geq 8\lambda_\epsilon/\delta$ and the equality follows from the definition of $M$. Finally, this implies that

(SM3.50)                           $$\Omega(\tilde{\beta}-\beta^0)+\|\tilde{\beta}_0-\beta_0^0\|_1\leq M/2,$$

which in turn, by the construction of $\tilde{\beta}$, implies that

(SM3.51)                           $$\Omega(\hat{\beta}-\beta^0)+\|\hat{\beta}_0-\beta_0^0\|_1\leq M.$$

**SM4. Optimization Algorithms.** In the main text, we presented a projected gradient algorithm for optimization. Here, we present some alternative methods for optimization of the MTD objective and discuss in what contexts they might be applicable.

**SM5. Frank-Wolfe.** In very high-dimensional settings, with large state spaces, the projection step in the MTD projected gradient algorithm presented in the main text becomes increasingly more computationally intensive. Frank-Wolfe algorithms, on the other hand, are projection free algorithms for solving constrained convex optimization problems and have recently gained popularity due to their simplicity and scalability in sparse, high-dimensional regression and machine learning [SM4]. Fortunately, the Frank-Wolfe algorithm for MTD also takes a simple form that allows updating only a small number of parameters at a time. In very sparse, high dimensional problems with large state spaces, where most entries are zero, this is typically advantageous [SM4]. We develop the algorithm and provide a timing comparison to the projected gradient algorithm in the main text. We leave the development of Frank-Wolfe using various variants [SM5] for future work.

**SM5.1. Frank-Wolfe MTD.** Let $\mathbf{Z}^{(0)}$ be the initial MTD model. Let $L(\mathbf{Z})=L_{MTD}(\mathbf{Z})+\lambda\Omega(\mathbf{Z})$. The Frank-Wolfe algorithm iterates between the following steps starting with $k=0$:

  1. Find a direction $\hat{\mathbf{D}}$ that maximizes the dot product with the gradient while staying in the constraint set:

(SM5.1)        $$\hat{\mathbf{D}}=\operatorname*{argmin}_{\mathbf{D}}\left(\mathbf{z}^0\right)^T\nabla_{\mathbf{z}^0}L\left(\mathbf{Z}^{(k)}\right)+\sum_{j=1}\operatorname{trace}\left(\left(\mathbf{D}^j\right)^T\nabla_{\mathbf{Z}^j}L\left(\mathbf{Z}^{(k)}\right)\right)$$

$$\text{subject to }\mathbf{1}^T\mathbf{D}^j=\gamma_j\mathbf{1}^T,\ \mathbf{D}^j\geq 0\ \forall j,\quad \mathbf{1}^T\gamma=1,\gamma\geq 0.$$

  2. Choose $\theta$ by line search or set $\theta=\frac{2}{2+k}$.
  3. Set $\mathbf{Z}^{(k+1)}=\theta\hat{\mathbf{D}}+(1-\theta)\mathbf{Z}^{(k)}$.

Step 1 involves solving a linear programming problem. Since the solution to Step 1 stays in the constraint set, any step taken in Step 2 for $\theta\in(0,1)$ remains in the constraint set. Fortunately, the linear program in Step 1 has a simple, closed form solution with linear complexity in the number of parameters, $O\left(m^2 d+m\right)$.

**Proposition SM5.1.** *First let* $\mathbf{F}^j = \nabla_{\mathbf{Z}^j} L\left(\mathbf{Z}^{(k)}\right)$. *Let* $q_k^j$ *be the row index of the minimal element in column* $k$ *of* $\mathbf{F}_{:k}^j$ *and let* $s^j$ *be the sum of the minimal elements in each column:* $s^j = \sum_{k=1}^m \mathbf{Z}_{q_k^j k}^j$. *Furthermore, let* $j^*$ *be the index of the minimum* $s^j$ : $j^* = \underset{j}{\operatorname{argmin}}\left(s^j\right)$. *Then* $\mathbf{D}^*$ *is given by*

$$\hat{\mathbf{D}}^j = 0 \ \forall j \neq j^*,$$

$$\hat{\mathbf{D}}_{lk}^{j*} = \begin{cases} 1 & \textit{if } l = q_k^j \\ 0 & \textit{if } l \neq q_k^j \end{cases}.$$

Intuitively, to stay in the MTD constraint set any feasible step must place equal mass on each column of a $\mathbf{Z}^j$, and that the minima is attained by only taking steps in the direction of $\mathbf{Z}^j$ with a minimal sum of columnwise minima.

Proposition SM5.1 implies that if the model is initialized with $\left(\mathbf{Z}^j\right)^{(0)} = 0$ for all $j$, then at step $k$ at most only $km$ entries in $\mathbf{Z}^{(k)}$ will be nonzero, and typically less in high-dimensional sparse settings since certain entries with strong signal will be updated repeatedly. The final Frank-Wolfe algorithm for MTD is shown in Algorithm SM5.1.

*Proof of Proposition SM5.1.* We study the KKT conditions. The Lagrangian is given by:

$$\sum_j \sum_l \sum_k \mathbf{D}_{lk}^j \mathbf{F}_{lk}^j + \sum_j \sum_k \lambda_k^j \left(\left(\sum_l \mathbf{D}_{lk}^j\right) - \gamma_j\right) + \nu\left(1^T \gamma - 1\right) + \sum_j \sum_k \sum_l \phi_{lk}^j \mathbf{D}_{lk}^j.$$

So that the KKT conditions for an optima are given by:

(SM5.2)           $$\mathbf{F}_{lk}^j = \lambda_k^j + \gamma_{lk}^j;$$

(SM5.3)           $$\sum_k^{m_j} \lambda_k^j = \nu \quad \forall j;$$

(SM5.4)           $$\phi_{lk}^j \geq 0 \ (\text{dual feasibility}) ;$$

(SM5.5)           $$\phi_{lk}^j \hat{\mathbf{D}}_{lk}^j = 0 \quad (\text{complimentary slackness}) .$$

We show that for the primary feasible solution given in Proposition SM5.1, there exists a set of dual variables that obey the KKT conditions, showing that the solution in Proposition SM5.1 is indeed the global optima.

For the primal solution given in Proposition SM5.1, let the dual variables for $j^*$ be

$$\lambda_k^{j*} = \mathbf{F}_{q_k^{j*} k}^{j*} \text{ and } \phi_{q_k^{j*} k}^{j*} = 0 \ \forall k \in (1, \ldots, m_j),$$

which obeys (SM5.2) and the complimentary slackness in (SM5.5) since $\hat{\mathbf{D}}_{q_k^{j*}}^{j*} = 1$. For all other entries of $\hat{\mathbf{D}}^{j*}$, $\phi_{lk}^{j*} = F_{lk}^j - \lambda_k^j = F_{lk}^j - F_{q_k^{j*} k}^{j*}$, so that all entries in $\phi_{lk}^{j*}$ and $\lambda_k^{j*}$ obey the KKT conditions for all $l, k$ in (SM5.4). The complimentary slackness holds in (SM5.5) since for these $l, k$ $\hat{\mathbf{D}}_{lk}^{j*} = 0$. Finally, set $\nu = \sum_k^{m_{j*}} \lambda_k^{j*} = \sum_k^{m_{j*}} F_{q_k^{j*} k}^{j*}$ which by construction satisfies condition (SM5.3).

445      For $j \neq j^*$, let $\lambda_k^j = F_{q_k^j k}^j - \frac{\tilde{\nu}^j - \nu}{m_j}$ where $\tilde{\nu}^j = \sum_k^{m_j} F_{q_k^j k}^j$. By construction, $\sum_j^{m_j} \lambda_k^j = \nu$

446 satisfying (SM5.3). Furthermore, letting $\phi_{lk}^j = F_{lk}^j - \lambda_k^j$, we have that $\phi_{lk}^j > 0$ since $F_{lk}^j >$

447 $F_{q_k^j k}^j > F_{q_k^j k}^j - \frac{\tilde{\nu}^j - \nu}{m_j} = \lambda_k^j$ and $\tilde{\nu}^j - \nu = \sum_k^{m_j} F_{q_k^j k}^j - \sum_k^{m_{j*}} F_{q_k^{j*} k}^{j*} > 0$ satisfying (SM5.4). For all

448 these entries the complimentary slackness condition holds since $\hat{\mathbf{D}}_{lk}^j = 0$, satisfying (SM5.5).

449      Taken together, we have found a set of dual feasible points that obey the KKT conditions

450 for the solution in Proposition SM5.1, showing that the solution is the optima.      ■

---

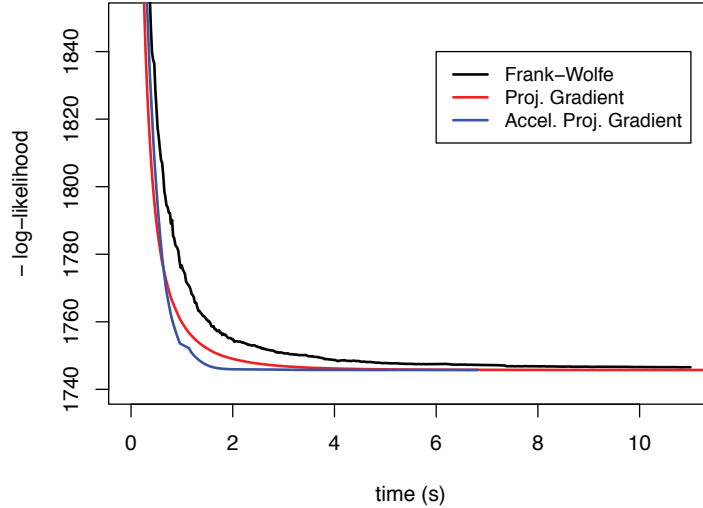**Algorithm SM5.1** Projection free Frank-Wolfe algorithm for MTD.

---

    Initialize $\left(\mathbf{Z}^j\right)^{(0)} = 0 \quad \forall j, \left(\mathbf{z}^0\right)^{(0)} = \frac{1}{m}$
    **for** $k = 0, 1, 2, \ldots$ **do**
       compute $\nabla L \left(\mathbf{Z}^{(k)}\right)$
       determine $\hat{\mathbf{D}}$ according to Proposition SM5.1
       determine $\theta$ by line search or $\theta = \frac{2}{2+k}$
       $\mathbf{Z}^{(k)} = (1 - \theta)\mathbf{Z}^{(k+1)} + \theta\hat{\mathbf{D}}$
    **end for**

---

451      **SM5.2. Run time comparison between Frank-Wolfe and Projected Gradient.** We com-
452 pare the Frank-Wolfe algorithm for MTD to the projected gradient algorithm in the main text.
453 In Figure SM11 we show the value of the objective as a function of time for Frank-Wolfe, pro-
454 jected gradient descent, and accelerated projected gradient descent on a synthetic data set.
455 For Frank-Wolfe, we use the step size of $\theta = \frac{2}{2+k}$. In this case, the Frank-Wolfe algorithm is
456 slower to converge than the projected or accelerated projected gradient algorithm. We suspect
457 that the gains of Frank-Wolfe over projected gradient will be in very high-dimensional settings
with large state spaces, but we leave that exploration for future work.



**Figure SM11.** *Run time comparison between Frank-Wolfe, projected gradient, and accelerated projected gradient on a $d = 25$, $T = 400$, and $m = 5$ synthetic data set.*

**SM5.3. Majorization-Minimization.** Here we use the convex formulation of MTD in the main text to derive a majorization-minimization (MM) algorithm [SM3]. The closed form updates are only given when there is no penalty function $\Omega(\mathbf{Z})$, so that this algorithm is not as generally applicable as the projected gradient algorithm presented in the main text. Interestingly, we find that the MM updates of the convex formulation correspond exactly to the MTD EM algorithm of [SM6] for the non-convex parameterization. This proves that the EM algorithm for MTD converges to a global optima even though the log-likelihood is non-convex.

We derive the MM algorithm for the convex MTD formulation with no penalty term (and no intercept):

(SM5.6)
$$\begin{aligned} \underset{\mathbf{Z},\gamma}{\text{minimize}}\ & L_{\text{MTD}}(\mathbf{Z}) \\ \text{subject to}\ & \mathbf{1}^T\mathbf{Z}^j = \gamma_j\mathbf{1}^T, \mathbf{Z}^j \geq 0\ \forall j, \quad \mathbf{1}^T\gamma = 1, \gamma \geq 0. \end{aligned}$$

To derive the MM algorithm, we first form the surrogate function

$$Q\left(\mathbf{Z},\mathbf{Z}^{(n)}\right) = \sum_{t=1}^{T}\sum_{j=1}^{d} p_{jt}\log\frac{Z^j_{x_{it}x_{j(t-1)}}}{p_{jt}},$$

where $p_{jt} = \frac{Z^{j(n)}_{x_{it}x_{j(t-1)}}}{\sum_{l=1}^{d} Z^{l(n)}_{x_{it}x_{l(t-1)}}}$. Now, $Q\left(\mathbf{Z},\mathbf{Z}^{(n)}\right)$ satisfies the MM algorithm conditions that $Q\left(\mathbf{Z},\mathbf{Z}^{(n)}\right) \geq L_{\text{MTD}}(\mathbf{Z})$ and $Q(\mathbf{Z},\mathbf{Z}) = L_{\text{MTD}}(\mathbf{Z})$. This implies we may iteratively minimize $Q\left(\mathbf{Z},\mathbf{Z}^{(n)}\right)$:

$$\mathbf{Z}^{(n+1)} = \underset{\mathbf{Z},\gamma}{\text{argmin}}\, Q\left(\mathbf{Z},\mathbf{Z}^{(n)}\right),$$

and that this sequence of $\mathbf{Z}^{(n+1)}$ converges to a global optima since Problem (SM5.6) is convex.

**Proposition SM5.2.** *The solution to Problem (SM5.6) under the MTD constraints is given in closed form:*

(SM5.7)
$$\mathbf{Z}^{j(n+1)}_{lk} = \left(\frac{\tilde{p}^j_{lk}}{\sum_l \tilde{p}^j_{lk}}\right)\left(\frac{\sum_{lk}\tilde{p}^j_{lk}}{\sum_j\sum_{lk}\tilde{p}^j_{lk}}\right),$$

*where $\tilde{p}^j_{lk} = \sum_{t=1} p_{jt}\mathbf{1}_{\left(x_{it}=l,x_{j(t-1)}=k\right)}$.*

**Corollary SM5.3.** *The EM algorithm for the unpenalized MTD model in the original $(\gamma,\mathbf{P})$ parameterization converges to a global optima of the non-convex log-likelihood.*

*Proof of Proposition SM5.2 and Corollary SM5.3.* The optimization problem for the MM update in Problem (SM5.6) is given by

(SM5.8)
$$\underset{\mathbf{Z},\gamma}{\text{minimize}} -\sum_{t=1}^{T}\sum_{j=1}^{d} p_{jt}\log\frac{Z^j_{x_{it}x_{j(t-1)}}}{p_{jt}}$$

487

subject to $\mathbf{1}^T \mathbf{Z}^j = \gamma_j \mathbf{1}^T \ \forall j, \quad \mathbf{1}^T \gamma = 1,$

where we have removed the non-negativity constraints because these are automatically enforced in the log terms of the $Q\left(Z, Z^{(n)}\right)$ objective. We may first rewrite the objective in (SM5.8) equivalently as

(SM5.9)
$$\underset{\mathbf{Z},\gamma}{\text{minimize}} -\sum_{j=1}^{d}\sum_{l=1}^{m}\sum_{k=1}^{m} \tilde{p}_{lk}^{j} \log Z_{lk}^{j}$$

subject to $\mathbf{1}^T \mathbf{Z}^j = \gamma_j \mathbf{1}^T \ \forall j, \mathbf{1}^T \gamma = 1,$

where $\tilde{p}_{lk}^{j} = \sum_{t=1} p_{jt} \mathbf{1}_{\left(x_{it}=l, x_{j(t-1)}=k\right)}$. We derive the solution by solving the KKT conditions. The Lagrangian of (SM5.9) is given by

$$\sum_{j=1}^{d}\sum_{l=1}^{m}\sum_{k=1}^{m} \tilde{p}_{lk}^{j} \log Z_{lk}^{j} + \sum_{j}\sum_{k} \lambda_{k}^{j} \left( \left( \sum_{l} Z_{lk}^{j} \right) - \gamma_j \right) + \nu \left( \mathbf{1}^T \gamma - 1 \right),$$

where $\lambda_j^k$ and $\nu$ are Lagrange multipliers. The solution must satisfy the KKT conditions: [SM1]

(SM5.10)
$$Z_{lk}^{j} = \frac{\tilde{p}_{lk}^{j}}{\lambda_{k}^{j}} \ \forall j, l, k,$$

(SM5.11)
$$\nu = \sum_{k} \lambda_{k}^{j} \ \forall j,$$

(SM5.12)
$$\mathbf{1}^T \mathbf{Z}^j = \gamma_j \mathbf{1}^T \ \forall j, \quad \mathbf{1}^T \gamma = 1.$$

Summing over Equation (SM5.10) for all rows $l$ gives

$$\gamma_j = \frac{\sum_{l} \tilde{p}_{lk}^{j}}{\lambda_{k}^{j}}.$$

Re-arranging and summing over $k$ gives

$$\frac{\sum_{lk} \tilde{p}_{lk}^{j}}{\gamma_j} = \sum_{k} \lambda_{k}^{j} = \nu,$$

and finally re-arranging once more and summing over $j$ gives

$$\frac{\sum_{j} \sum_{lk} \tilde{p}_{lk}^{j}}{\nu} = \sum_{j} \gamma_j = 1.$$

Plugging these results back into those above implies that $\nu = \sum_j \sum_{lk} \tilde{p}^j_{lk}$, $\gamma_j = \frac{\sum_{lk} \tilde{p}^j_{lk}}{\sum_j \sum_{lk} \tilde{p}^j_{lk}}$, $\lambda^j_k = \frac{(\sum_l \tilde{p}^j_{lk})(\sum_j \sum_{lk} \tilde{p}^j_{lk})}{\sum_{lk} \tilde{p}^j_{lk}}$. Plugging into Equation (SM5.10) gives the final update for $\mathbf{Z}^{(n+1)}$ as

(SM5.13)
$$Z^{j(n+1)}_{lk} = \left(\frac{\tilde{p}^j_{lk}}{\sum_l \tilde{p}^j_{lk}}\right)\left(\frac{\sum_{lk} \tilde{p}^j_{lk}}{\sum_j \sum_{lk} \tilde{p}^j_{lk}}\right)$$

(SM5.14)
$$= P^{j(n+1)}_{lk} \gamma^{(n+1)}_j,$$

where $P^{j(n+1)}_{lk} = \left(\frac{\sum_{lk} \tilde{p}^j_{lk}}{\sum_j \sum_{lk} \tilde{p}^j_{lk}}\right)$ and $\gamma^{(n+1)}_j = \left(\frac{\tilde{p}^j_{lk}}{\sum_l \tilde{p}^j_{lk}}\right)$.

This update for $P^{j(n+1)}_{lk}$ and $\gamma^{(n+1)}_j$ is identical to the updates for the EM algorithm in the original $(\mathbf{P}, \gamma)$ parameterization [SM6]. Since the MM algorithm on a convex problem converges to a global optima, it follows that the EM algorithm for the original non-convex MTD parameterization also converges to a global optima. ∎

## REFERENCES

[1] S. BOYD AND L. VANDENBERGHE, *Convex optimization*, Cambridge university press, 2004.

[2] A. HARIS, N. SIMON, AND A. SHOJAIE, *Generalized sparse additive models*, arXiv preprint arXiv:1903.04641, (2019).

[3] D. R. HUNTER AND K. LANGE, *A tutorial on mm algorithms*, The American Statistician, 58 (2004), pp. 30–37, https://doi.org/10.1198/0003130042836, https://doi.org/10.1198/0003130042836, https://arxiv.org/abs/https://doi.org/10.1198/0003130042836.

[4] M. JAGGI, *Revisiting frank-wolfe: Projection-free sparse convex optimization.*, in ICML (1), 2013, pp. 427–435.

[5] S. LACOSTE-JULIEN AND M. JAGGI, *On the global linear convergence of frank-wolfe optimization variants*, in Advances in Neural Information Processing Systems 28, C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, eds., Curran Associates, Inc., 2015, pp. 496–504, http://papers.nips.cc/paper/5925-on-the-global-linear-convergence-of-frank-wolfe-optimization-variants.pdf.

[6] S. LÈBRE AND P.-Y. BOURGUIGNON, *An em algorithm for estimation in the mixture transition distribution model*, Journal of Statistical Computation and Simulation, 78 (2008), pp. 713–729, https://doi.org/10.1080/00949650701266666, https://doi.org/10.1080/00949650701266666, https://arxiv.org/abs/https://doi.org/10.1080/00949650701266666.

[7] Y.-A. MA, N. J. FOTI, AND E. B. FOX, *Stochastic gradient mcmc methods for hidden markov models*, in Proceedings of the 34th International Conference on Machine Learning-Volume 70, JMLR. org, 2017, pp. 2265–2274.

[8] D. PAULIN ET AL., *Concentration inequalities for markov chains by marton couplings and spectral methods*, Electronic Journal of Probability, 20 (2015).

[9] A. RAKHLIN, K. SRIDHARAN, AND A. TEWARI, *Sequential complexities and uniform martingale laws of large numbers*, Probability Theory and Related Fields, 161 (2015), pp. 111–153.

[10] S. VAN DE GEER, *Estimation and testing under sparsity*, Lecture notes in mathematics, 2159 (2016).

[11] D. WULSIN, E. FOX, AND B. LITT, *Parsing epileptic events using a markov switching process model for correlated time series*, in International Conference on Machine Learning, 2013, pp. 356–364.