

Spectral Learning of Hidden Markov Models with Group Persistence, Appendix

1. Proofs:

Lemma 3: K eigenvalues of an SHMM global transition matrix are same as the eigenvalues of its corresponding regime transition matrix B .

Proof: Let us consider the product $A^\top(b \otimes \mathbf{1}_M)$, where \otimes denotes the Kronecker product, and $b \in \mathbb{R}^M$ is an eigenvector of B^\top . Let λ_b denote the corresponding eigenvalue.

$$\begin{aligned} A^\top(b \otimes \mathbf{1}_M) &= \begin{bmatrix} \vdots \\ \sum_k B_{k,1} b_k A_{k,1}^\top \mathbf{1}_M \\ \vdots \end{bmatrix} = \begin{bmatrix} \vdots \\ \lambda_b b_k \mathbf{1}_M \\ \vdots \end{bmatrix} \\ &= \lambda_b (b \otimes \mathbf{1}_M). \end{aligned}$$

Note that each block $A_{i,j}$ is stochastic matrix with column sums equal to 1, so every $A_{i,j}^\top$ has an eigenvector equal to with eigenvalue one. And, since A and A^\top (and B, B^\top) have the same eigenvalues, we conclude that λ_b is an eigenvector of A . This argument applies to all eigenvectors b of B^\top . \square

Theorem 2:

$$\begin{aligned} \|\mathcal{P}(\Lambda) - \hat{\Lambda}\|_F &\leq c_3 \left(c_1 \frac{1 + \log(1/\delta)}{\sqrt{T}} + c_2 \|O^\dagger - \hat{O}^\dagger\|_F \right. \\ &\quad \left. + c_3 \|\xi^{-1} - \hat{\xi}^{-1}\|_F \right), \text{ with probability } 1 - \delta. \end{aligned}$$

where, $c_1 = \|O^\dagger\|^2 \|\xi^{-1}\|$, $c_2 = \|\widehat{S}_{1,2}\| \|\xi^{-1}\| \|(O^\top)^\dagger\| + \|\hat{O}^\dagger\| \|\widehat{S}_{1,2}\| \|\hat{\xi}^{-1}\|$, $c_3 = \|\hat{O}^\dagger\| \|\widehat{S}_{1,2}\| \|(O^\top)^\dagger\|$, and $c_4 = \sqrt{\kappa(A)\kappa(\hat{A})}$ and, T is the number data items used for $\widehat{S}_{1,2}$. We denote $\text{diag}(1/\xi)$ and $\text{diag}(1/\hat{\xi})$ respectively by, ξ^{-1} and $\hat{\xi}^{-1}$ to save space.

Proof: We know from (Bhatia et al., 1997) that for two diagonalizable matrices $A = V\Lambda V^{-1}$ and $\hat{A} = \widehat{V}\widehat{\Lambda}\widehat{V}^{-1}$,

$$\|\mathcal{P}(\Lambda) - \hat{\Lambda}\|_F \leq \sqrt{\kappa(A)\kappa(\hat{A})} \|A - \hat{A}\|_F. \quad (1)$$

So, proving the bound amounts to upper bounding the deviation in the estimation $\|A - \hat{A}\|_F$. We use the estimator

given in Equation 5 (in the paper, pseudo inverse estimator) for the sake of analysis.

$$\begin{aligned} \|A - \hat{A}\|_F & \quad (2) \\ &= \|O^\dagger S_{1,2} \xi^{-1} (O^\top)^\dagger - \hat{O}^\dagger \widehat{S}_{1,2} \widehat{\xi}^{-1} (\hat{O}^\top)^\dagger\|_F \\ &\leq \|O^\dagger S_{1,2} \xi^{-1} (O^\top)^\dagger - O^\dagger \widehat{S}_{1,2} \xi^{-1} (O^\top)^\dagger\|_F \\ &\quad + \|O^\dagger \widehat{S}_{1,2} \xi^{-1} (O^\top)^\dagger - \hat{O}^\dagger \widehat{S}_{1,2} \widehat{\xi}^{-1} (\hat{O}^\top)^\dagger\|_F \\ &\leq \|O^\dagger S_{1,2} \xi^{-1} (O^\top)^\dagger - O^\dagger \widehat{S}_{1,2} \xi^{-1} (O^\top)^\dagger\|_F \\ &\quad + \|O^\dagger \widehat{S}_{1,2} \xi^{-1} (O^\top)^\dagger - \hat{O}^\dagger \widehat{S}_{1,2} \xi^{-1} (O^\top)^\dagger\|_F \\ &\quad + \|\hat{O}^\dagger \widehat{S}_{1,2} \xi^{-1} (O^\top)^\dagger - \hat{O}^\dagger \widehat{S}_{1,2} \widehat{\xi}^{-1} (\hat{O}^\top)^\dagger\|_F, \\ &\leq \|O^\dagger S_{1,2} \xi^{-1} (O^\top)^\dagger - O^\dagger \widehat{S}_{1,2} \xi^{-1} (O^\top)^\dagger\|_F \\ &\quad + \|O^\dagger \widehat{S}_{1,2} \xi^{-1} (O^\top)^\dagger - \hat{O}^\dagger \widehat{S}_{1,2} \xi^{-1} (O^\top)^\dagger\|_F \\ &\quad + \|\hat{O}^\dagger \widehat{S}_{1,2} \xi^{-1} (O^\top)^\dagger - \hat{O}^\dagger \widehat{S}_{1,2} \widehat{\xi}^{-1} (\hat{O}^\top)^\dagger\|_F, \\ &= \|O^\dagger (S_{1,2} - \widehat{S}_{1,2}) \xi^{-1} (O^\top)^\dagger\|_F \\ &\quad + \|(O^\dagger - \hat{O}^\dagger) \widehat{S}_{1,2} \xi^{-1} (O^\top)^\dagger\|_F \\ &\quad + \|\hat{O}^\dagger \widehat{S}_{1,2} (\xi^{-1} - \widehat{\xi}^{-1}) (O^\top)^\dagger\|_F \\ &\quad + \|\hat{O}^\dagger \widehat{S}_{1,2} \widehat{\xi}^{-1} ((O^\top)^\dagger - (\hat{O}^\top)^\dagger)\|_F, \end{aligned}$$

where the inequalities are due to triangular inequality. The next step is to bound the individual terms:

$$\begin{aligned} \|O^\dagger (S_{1,2} - \widehat{S}_{1,2}) \xi^{-1} (O^\top)^\dagger\| & \quad (3) \\ &\leq \|O^\dagger\| \|S_{1,2} - \widehat{S}_{1,2}\| \|\xi^{-1}\| \|(O^\top)^\dagger\| \\ &\leq \|O^\dagger\|^2 \|\xi^{-1}\| \frac{1 + \log(1/\delta)}{\sqrt{T}}, \text{ with probability } 1 - \delta, \end{aligned}$$

where we omitted the subscript F to save space, used the property $\|AB\|_F \leq \|A\|_F \|B\|_F$ of the Frobenius norm, and the second inequality is from (Hsu et al., 2009). The three remaining terms are also handled using the same property of the Frobenius norm,

$$\begin{aligned} \|(O^\dagger - \hat{O}^\dagger) \widehat{S}_{1,2} \xi^{-1} (O^\top)^\dagger\| & \quad (4) \\ &\leq \|(O^\dagger - \hat{O}^\dagger)\| \|\widehat{S}_{1,2}\| \|\xi^{-1}\| \|(O^\top)^\dagger\|, \end{aligned}$$

055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107
108
109

$$\begin{aligned} & \|\widehat{O}^\dagger \widehat{S}_{1,2} (\xi^{-1} - \widehat{\xi}^{-1}) (O^\top)^\dagger\|_F \\ & \leq \|\widehat{O}^\dagger\| \|\widehat{S}_{1,2}\| \|\xi^{-1} - \widehat{\xi}^{-1}\| \|(O^\top)^\dagger\| \end{aligned} \quad (5)$$

$$\begin{aligned} & \|\widehat{O}^\dagger \widehat{S}_{1,2} \widehat{\xi}^{-1} ((O^\top)^\dagger - (\widehat{O}^\top)^\dagger)\| \\ & \leq \|\widehat{O}^\dagger\| \|\widehat{S}_{1,2}\| \|\widehat{\xi}^{-1}\| \|O^\dagger - \widehat{O}^\dagger\|. \end{aligned} \quad (6)$$

By arranging terms we see that $c_1 = \|O^\dagger\|^2 \|\xi^{-1}\|$, $c_2 = \|\widehat{S}_{1,2}\| \|\xi^{-1}\| \|(O^\top)^\dagger\| + \|\widehat{O}^\dagger\| \|\widehat{S}_{1,2}\| \|\widehat{\xi}^{-1}\|$, $c_3 = \|\widehat{O}^\dagger\| \|\widehat{S}_{1,2}\| \|(O^\top)^\dagger\|$, and $c_4 = \sqrt{\kappa(A)\kappa(\widehat{A})}$. \square

2. Additional Experiment

We do the depermutation experiment in Section 4.1 also for the HMM-M model. We could not include this result in the manuscript due to space limitations. It is given in Figure 1 of the appendix.

References

- Bhatia, Rajendra, Kittaneh, Fuad, and Li, Ren Cang. Some inequalities for commutators and an application to spectral variation. *Linear and Multilinear Algebra*, 1997.
- Hsu, Daniel, Kakade, Sham M., and Zhang, Tong. A spectral algorithm for learning hidden markov models. *Journal of Computer and System Sciences*, (1460-1480), 2009.

165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219

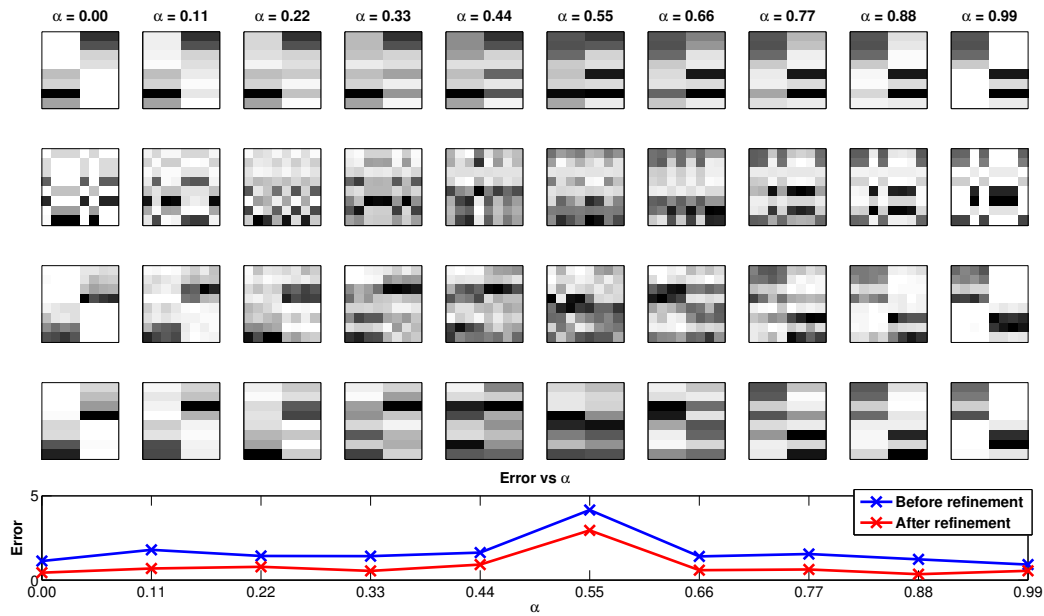


Figure 1. Depermuted estimated transition matrices on synthetic data, (First row) True Transition Matrices, (Second Row) Low Rank Reconstruction \hat{A}_r , (Third Row) $\tilde{\mathcal{P}}(A)$ before refinement, (Fourth Row) $\tilde{\mathcal{P}}(A)$ after refinement, (Fifth Row) The error $\sum_{i,j} \|\tilde{\mathcal{P}}(\hat{A}_{i,j}) - B_{\mathcal{P}_2(i,j)} \mathcal{P}_1(A_{\mathcal{P}_2(i,j)})\|_F$, before and after refinement.