



Finanziato  
dall'Unione europea  
NextGenerationEU



Ministero  
dell'Università  
e della Ricerca



# A REGULARIZED MAXIMUM LIKELIHOOD ESTIMATION FOR HIDDEN MARKOV MODELS WITH COVARIATES

---

LUCA BRUSA<sup>1</sup>, FULVIA PENNONI<sup>1</sup>,  
FRANCESCO BARTOLUCCI<sup>2</sup>, ROMINA PERUILH BAGOLINI<sup>2</sup>  
([luca.brusa@unimib.it](mailto:luca.brusa@unimib.it))

<sup>1</sup>University of Milano-Bicocca - Department of Statistics and Quantitative Methods

<sup>2</sup>University of Perugia - Department of Economics

# Outline

- 1 Hidden Markov models
- 2 Penalized maximum likelihood approach
- 3 Simulation study
- 4 Application
- 5 Conclusions
- 6 References

# Hidden Markov model: notation

- **Univariate binary response variables**  $\mathbf{Y}_i = (Y_i^{(1)}, \dots, Y_i^{(T)})$ , with

$$Y_i^{(t)} = \begin{cases} 1 & \text{if the event of interest is observed at time } t \text{ for unit } i \\ 0 & \text{otherwise} \end{cases}$$

- **Hidden process:**  $\mathbf{U}_i = (U_i^{(1)}, \dots, U_i^{(T)})$ , following a first-order Markov chain with state-space  $\{1, \dots, k\}$ , initial probabilities  $\pi_u$ , and transition probabilities  $\pi_{u|\bar{u}}$
- **Covariates:**  $\mathbf{x}_i = (\mathbf{x}_i^{(1)}, \dots, \mathbf{x}_i^{(T)})$ , with  $\mathbf{x}_i^{(t)}$  representing the vector of observed individual covariates for unit  $i$  at time  $t$ 
  - **lagged response:** allows for serial dependence between responses over time, thus relaxing the conditional independence of  $\mathbf{Y}$  given  $\mathbf{U}$  and  $\mathbf{x}$

# Model formulation

## 1 Measurement sub-model (affected by covariates):

- Conditional distribution of the response variable  $Y_i^{(t)}$  given the latent variable  $U_i^{(t)}$  and the covariates  $\mathbf{x}_i^{(t)}$
- **Conditional response probabilities:**

$$\phi_{u\mathbf{x}_i^{(t)}} = \mathbb{P} \left( Y_i^{(t)} = 1 \mid U_i^{(t)} = u, \mathbf{x}_i^{(t)} \right),$$

such that:

$$\log \frac{\phi_{u\mathbf{x}_i^{(t)}}}{1 - \phi_{u\mathbf{x}_i^{(t)}}} = \alpha_u + (\mathbf{x}_i^{(t)})' \beta$$

## 2 Latent sub-model (not affected by covariates):

- Non-parametric distribution of the latent process
- **Unobserved heterogeneity** between individuals, which remains when covariates in the measurement model cannot fully explain the variability

# Maximum likelihood estimation

- **Expectation-maximization** (EM) algorithm (Dempster et al., 1977) is often employed to perform **maximum likelihood estimation**
- It maximizes the observed-data log-likelihood function  $\ell(\theta)$  relying on the **complete-data log-likelihood** function  $\ell^*(\theta)$
- It alternates the following steps until convergence:
  - **E-step: compute the conditional expected value** of  $\ell^*(\theta)$  given the value of the parameters at the previous step and the observed data
  - **M-step: update the model parameters** by maximizing the expected value of  $\ell^*(\theta)$ :
    - explicit solutions are available for  $\pi_u$  and  $\pi_{u|\bar{u}}$
    - a Newton-Raphson algorithm is used for updating  $\alpha$  and  $\beta$

# Outline

- 1 Hidden Markov models
- 2 Penalized maximum likelihood approach
- 3 Simulation study
- 4 Application
- 5 Conclusions
- 6 References

# Motivation

- When the model is applied to **binary** and categorical response variables with a limited number of categories, the estimated support points  $\hat{\alpha}_u$  may be very large, leading to **widely separated latent states** (**separation problem**)
- This may results in:
  - **excessively higher relevance of one or more latent states** than others
  - **reduced importance of the available covariates** whose estimated effects may become negligible and insignificant
  - **instability of the estimates**

# Penalized maximum likelihood estimation

- Proposed **penalty**: to **reduce latent states separation**:

$$\mathcal{A} = \sum_{u=1}^k (\alpha_u - \bar{\alpha})^2,$$

where  $\bar{\alpha} = \frac{1}{k} \sum_{u=1}^k \alpha_u$

- Applied to** both the **observed-data log-likelihood**  $\ell(\boldsymbol{\theta})$  and the **complete-data log-likelihood**  $\ell^*(\boldsymbol{\theta})$ :

$$\tilde{\ell}(\boldsymbol{\theta}) = \ell(\boldsymbol{\theta}) - \lambda \mathcal{A} \quad \text{and} \quad \tilde{\ell}^*(\boldsymbol{\theta}) = \ell^*(\boldsymbol{\theta}) - \lambda \mathcal{A},$$

where  $\lambda \in \mathbb{R}^+$  is a **tuning parameter** controlling the penalization

- Penalized estimation through the **EM algorithm**:
  - the **E-step** remains unaltered
  - the **M-step** requires to **revise the Newton-Raphson** iteration for  $\alpha$

# Outline

- 1 Hidden Markov models
- 2 Penalized maximum likelihood approach
- 3 Simulation study**
- 4 Application
- 5 Conclusions
- 6 References

# Simulation study

- **Different scenarios** (40) to explore the performance of the proposal with  $k = 3$  hidden states: **sample size** ( $n = 250, 500$ ), **number of time occasions** ( $T = 10, 20$ ), hidden **state persistence** (high or low) and **separation** (five different behaviors:  $\alpha^j, j = 1, \dots, 5$ )
- **Four covariates** also including the **lagged response variable**; the corresponding vector of regression coefficients is  $\beta = (1, -1, 1, 1)'$
- Extensive **Monte Carlo simulation study**; for each scenario:
  - we randomly draw 50 samples
  - we estimate the HM model using both the standard approach and the penalized approach with  $\lambda = 0.01$  and  $\lambda = 0.05$

# Comparison criteria

- **Percentage variation** in the following quantities for both procedures:
  - **root mean squared relative error between true and estimated model parameters ( $\theta$ )** defined as:

$$\text{RMSRE} = \sqrt{\frac{1}{M} \sum_{m=1}^M \left( \frac{\hat{\theta}_m - \theta_m}{\theta_m} \right)^2}$$

- **standard errors of the covariate regression parameters  $\beta$** , obtained as minus the second derivative of the expected value of the complete-data log-likelihood
- **computational time**

# Results - RMSRE

- In the majority of cases, both values of  $\lambda$  ensure a **lower RMSRE when using the penalized estimation** method. This indicates that the estimated parameters are closer to the true values (**higher estimation precision**)
- The **fourth and fifth scenarios** (characterized by highly separated states) show the **greatest improvement**; the percentage decreases using penalized estimation are often exceeding 90%
- The penalization approach appears **less effective in the first scenario** (characterized by closely spaced states) and in cases with a high value of  $T$  (20 time occasions)
- **Only in two cases** the penalized estimation method exhibits **no improvement** with either value of  $\lambda$ , showing values of the RMSRE that are very similar to the standard estimation

# Other results

- **Standard errors:**

- In most cases, **penalization reduces the estimated standard errors**
- The proposed approach is **less effective under the first scenario**
- In all other cases, the **percentage decrease is significant**, often reaching very high values

- **Computational time:**

- Estimation with the **penalty** approach often **reduces the average computational time**
- Benefits are **particularly evident in the fourth and fifth scenarios** where hidden states are widely separated

# Outline

- 1 Hidden Markov models
- 2 Penalized maximum likelihood approach
- 3 Simulation study
- 4 Application**
- 5 Conclusions
- 6 References

# Hypotension during spinal anesthesia

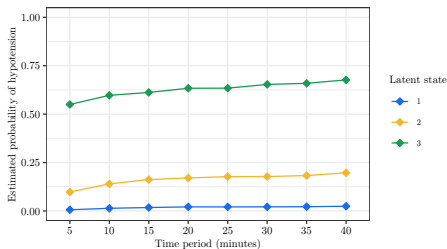
- Data<sup>1</sup> refer to **375 patients** undergoing spinal anesthesia during a surgery; they cover the period **from January 2008 to January 2011**
- Measurements are taken 8 times, at equally spaced intervals over a period of 40 minutes
- Variable  $Y_i^{(t)}$  indicates whether or not patient  $i$  has experienced **hypotension** (decrease in mean systolic blood pressure) at time  $t$
- Approximately **25% ( $n = 94$ ) of patients** recorded **at least one hypotensive episode**
- Time-fixed covariates, time-varying covariates, and lagged response
- Using a **cross-validation approach**, we select  **$k = 3$  hidden states** and a **penalization parameter of  $\lambda = 0.01$**  as optimal

---

<sup>1</sup>Data are freely available at <https://peerj.com/articles/648/>.

# Estimated conditional hypotension probability

- Patients in the **first hidden state** ( $\hat{\alpha}_1 = -0.827$ ) have an **almost negligible probability of hypotension** during the surgery
- Patients in the **second hidden state** ( $\hat{\alpha}_2 = 3.147$ ) experience a **low probability of hypotension**, ranging approximately from 0.10 to 0.20
- Patients in the **third hidden state** ( $\hat{\alpha}_3 = 7.359$ ) have a **high probability of hypotension** during surgery, ranging from 0.54 to 0.68



# Estimated regression coefficients

Covariate	$\hat{\beta}$	$\hat{se}$	$p$ -value
<b>Gender (Female)</b>	<b>1.541 *</b>	2.032	0.042
Position (Supine)	0.813	1.597	0.110
Operation (Urogoly)	0.363	0.338	0.735
Operation (General surgery)	-0.061	-0.068	0.946
ECG (Normal)	-0.878	-0.978	0.328
Age (year)	0.037 *	1.967	0.049
DBP	-0.167 **	-7.642	0.000
Pulse rate	-0.002	-0.216	0.829
Marcaïn-heavy	-0.049	-1.235	0.217
Midazolam	0.243 *	2.104	0.035
Chirocaine	-0.049	-1.346	0.178
Fentanyl	1.215	0.515	0.607
Hypotension ( $t - 1$ )	2.675 **	10.289	0.000

- Gender** (female) has a **significant positive effect** on the response variable, indicating that the conditional **probability of experimenting hypotension** given the latent state is **higher for females**

# Estimated regression coefficients

Covariate	$\hat{\beta}$	$\hat{se}$	p-value
Gender (Female)	1.541 *	2.032	0.042
Position (Supine)	0.813	1.597	0.110
Operation (Urogoly)	0.363	0.338	0.735
Operation (General surgery)	-0.061	-0.068	0.946
ECG (Normal)	-0.878	-0.978	0.328
Age (year)	0.037 *	1.967	0.049
DBP	-0.167 **	-7.642	0.000
Pulse rate	-0.002	-0.216	0.829
Marcaïn-heavy	-0.049	-1.235	0.217
Midazolam	0.243 *	2.104	0.035
Chirocaine	-0.049	-1.346	0.178
Fentanyl	1.215	0.515	0.607
Hypotension ( $t - 1$ )	2.675 **	10.289	0.000

- Older individuals exhibit higher log-odds of being diagnosed with hypotension compared to younger individuals

# Estimated regression coefficients

Covariate	$\hat{\beta}$	se	p-value
Gender (Female)	1.541 *	2.032	0.042
Position (Supine)	0.813	1.597	0.110
Operation (Urogoly)	0.363	0.338	0.735
Operation (General surgery)	-0.061	-0.068	0.946
ECG (Normal)	-0.878	-0.978	0.328
Age (year)	0.037 *	1.967	0.049
<b>DBP</b>	<b>-0.167 **</b>	-7.642	0.000
Pulse rate	-0.002	-0.216	0.829
Marcaïn-heavy	-0.049	-1.235	0.217
Midazolam	0.243 *	2.104	0.035
Chirocaine	-0.049	-1.346	0.178
Fentanyl	1.215	0.515	0.607
Hypotension ( $t - 1$ )	2.675 **	10.289	0.000

- **Diastolic blood pressure** has a significant **negative effect** on the log-odds of hypotension: lower pressure is associated with higher probabilities of experiencing hypotension

# Estimated regression coefficients

Covariate	$\hat{\beta}$	$\hat{se}$	$p$ -value
Gender (Female)	1.541 *	2.032	0.042
Position (Supine)	0.813	1.597	0.110
Operation (Urogoly)	0.363	0.338	0.735
Operation (General surgery)	-0.061	-0.068	0.946
ECG (Normal)	-0.878	-0.978	0.328
Age (year)	0.037 *	1.967	0.049
DBP	-0.167 **	-7.642	0.000
Pulse rate	-0.002	-0.216	0.829
Marcaïn-heavy	-0.049	-1.235	0.217
<b>Midazolam</b>	<b>0.243 *</b>	2.104	0.035
Chirocaine	-0.049	-1.346	0.178
Fentanyl	1.215	0.515	0.607
Hypotension ( $t - 1$ )	2.675 **	10.289	0.000

- **Midazolam** has a significant **positive effect**, indicating that higher concentration of this drug in the blood is associated with increased odds of experiencing hypotension during surgery. For the other drugs, the estimated coefficients are not significant

# Estimated regression coefficients

Covariate	$\hat{\beta}$	$\hat{se}$	p-value
Gender (Female)	1.541 *	2.032	0.042
Position (Supine)	0.813	1.597	0.110
Operation (Urogoly)	0.363	0.338	0.735
Operation (General surgery)	-0.061	-0.068	0.946
ECG (Normal)	-0.878	-0.978	0.328
Age (year)	0.037 *	1.967	0.049
DBP	-0.167 **	-7.642	0.000
Pulse rate	-0.002	-0.216	0.829
Marcaïn-heavy	-0.049	-1.235	0.217
Midazolam	0.243 *	2.104	0.035
Chirocaine	-0.049	-1.346	0.178
Fentanyl	1.215	0.515	0.607
<b>Hypotension (<math>t - 1</math>)</b>	<b>2.675 **</b>	10.289	0.000

- The **lagged response** has a significant **positive effect on hypotension** indicating serial correlation

# Outline

- 1 Hidden Markov models
- 2 Penalized maximum likelihood approach
- 3 Simulation study
- 4 Application
- 5 Conclusions**
- 6 References

## Limitations and future works

- **Implementation in C++ to improve computational speed**, particularly for datasets with a large number of repeated measurements and/or a large sample size
- **Evaluation of the model's predictive performance** and comparison with machine learning methods, also in connection with the use of HM models as early warning systems
- **Development of feature selection techniques** to identify relevant covariates, especially when many are available
- **Investigation of methods to achieve scalability**, such as parallel computation and dimension reduction, essential for handling large datasets efficiently

# Outline

- 1 Hidden Markov models
- 2 Penalized maximum likelihood approach
- 3 Simulation study
- 4 Application
- 5 Conclusions
- 6 References

# References I

- AKTAS SAMUR, A., COSKUNFIRAT, N., AND SAKA, O. (2014). Comparison of predictor approaches for longitudinal binary outcomes: Application to anesthesiology data. *PeerJ*, **2**, e648.
- BARTOLUCCI, F. AND FARCOMENI, A. (2009). A multivariate extension of the dynamic logit model for longitudinal data based on a latent markov heterogeneity structure. *J. Am. Stat. Assoc.*, **104**, 816–831.
- BARTOLUCCI, F., FARCOMENI, A., AND PENNONI, F. (2013). *Latent Markov models for longitudinal data*. Chapman & Hall/CRC, Boca Raton.
- BARTOLUCCI, F., FARCOMENI, A., AND PENNONI, F. (2014). Latent Markov models: a review of a general framework for the analysis of longitudinal data with covariates. *Test*, **23**, 433–465.

## References II

- BATES, D., HASTIE, T., AND TIBSHIRANI, R. (2023). Cross-validation: What does it estimate and how well does it do it? *J. Mach. Learn. Res.*, **24**, 1234–1256.
- DEMPSTER, A., LAIRD, N., AND RUBIN, D. (1977). Maximum likelihood from incomplete data via the EM algorithm (with discussion). *J. R. Stat. Soc. Series B Stat. Methodol.*, **39**, 1–38.
- SMYTH, P. (2000). Model selection for probabilistic clustering using cross-validated likelihood. *Stat. Comput.*, **9**, 63–72.
- WELCH, L. (2003). Hidden Markov models and the Baum-Welch algorithm. *IEEE Inform. Theory Soc. Newsl.*, **50**, 10–13.

# Acknowledgments

**Acknowledgment:** The authors acknowledge the financial support from the grant “Hidden Markov Models for Early Warning Systems” of Ministero dell’Università e della Ricerca (PRIN 2022TZEXKF) funded by European Union - Next Generation EU, Mission 4, Component 2, CUP J53D23004990006.

# State separation: a motivating example

- Data generation:
  - one binary response variable for  $n = 250$  units over  $T = 10$  time points
  - four covariates (including the lagged response), with  $\beta = [1, -1, 1, 1]'$
  - $k = 3$  latent states with support points  $\alpha = [-20, -5, 5]'$
  - 1<sup>st</sup> time point: 88 units in the 1<sup>st</sup> state, 71 in the 2<sup>nd</sup>, and 91 in the 3<sup>rd</sup>
- Model estimation using standard maximum likelihood:

	$\hat{\alpha}_1$	$\hat{\alpha}_2$	$\hat{\alpha}_3$	$\hat{\beta}_1$	$\hat{\beta}_2$	$\hat{\beta}_3$	$\hat{\beta}_4$
Estimate	-250.13	146.80	457.96	24.14	-80.21	56.69	103.48
Standard error	-	-	-	23.28	69.43	50.91	90.14
<i>p</i> -value	-	-	-	0.30	0.25	0.27	0.25

- poor accuracy in parameter estimation for both  $\alpha_u$  and  $\beta_j$
- 1<sup>st</sup> time point: 159 units in the 1<sup>st</sup> state and 91 in the 3<sup>rd</sup>
- conditional response probabilities equal to 0 (1<sup>st</sup> state) or 1 (3<sup>rd</sup> state)

# Cross-validated log-likelihood

- A **cross-validation** approach is employed to **jointly select** the **penalization parameter**  $\lambda$  and the **number of states**  $k$  of the hidden chain
- We consider  $M$  **partitions of the data**  $D: (D \setminus S_m, S_m)_{m=1, \dots, M}$
- For the  $m$ -th partition:
  - the model is estimated on the data subset  $D \setminus S_m$ , providing parameters estimates  $\hat{\theta}^{(k, \lambda)}(D \setminus S_m)$
  - $\ell\left(\hat{\theta}^{(k, \lambda)}(D \setminus S_m) \mid S_m\right)$  denotes the (possibly penalized) log-likelihood function where the model parameters are estimated on the training data  $D \setminus S_m$  but the log-likelihood is evaluated on the test data  $S_m$
  - the **cross-validated likelihood** is defined as

$$\ell_{\text{CV}} = \frac{1}{M} \sum_{m=1}^M \ell\left(\theta^{(k, \lambda)}(D \setminus S_m) \mid S_m\right).$$