

Asymptotic Theory for Linear-Chain Conditional Random Fields

Mathieu Sinn, Pascal Poupart
University of Waterloo

**14th International Conference on Artificial
Intelligence and Statistics, Ft. Lauderdale, Florida**

April 11, 2011



Motivation

An **L-CRF** models the distribution of \mathbf{Y}_n **conditional on** \mathbf{X}_n where

- ▶ $\mathbf{X}_n = (X_0, \dots, X_n)$ are the **observations** (ranging in \mathcal{X})
- ▶ $\mathbf{Y}_n = (Y_0, \dots, Y_n)$ are the **labels** (ranging in \mathcal{Y})

Basic assumption: \mathbf{Y}_n conditional on \mathbf{X}_n forms a **Markov Chain**.

Here we assume that

$$P_{\lambda}(\mathbf{Y}_n = \mathbf{y}_n \mid \mathbf{X}_n = \mathbf{x}_n) = \frac{1}{Z_{\lambda}(\mathbf{x}_n)} \prod_{t=1}^n \exp(\lambda^T \mathbf{f}(x_t, y_{t-1}, y_t))$$

for $\mathbf{x}_n = (x_0, \dots, x_n)$, $\mathbf{y}_n = (y_0, \dots, y_n)$.

Motivation

Suppose we observe \mathbf{X}_n and \mathbf{Y}_n where

- ▶ the distribution of \mathbf{X}_n is **unknown**
- ▶ the distribution of \mathbf{Y}_n conditional on \mathbf{X}_n follows an L-CRF with **known** features \mathbf{f} and **unknown** weights λ^* .

How can we **estimate** λ^* from the data?

Common approach is to maximize **conditional log-likelihood**, i.e.,

$$\text{compute } \hat{\lambda}_n = \arg \max_{\lambda \in \Theta} \mathcal{L}_n(\lambda)$$

$$\text{where } \mathcal{L}_n(\lambda) = \frac{1}{n} \left(\sum_{t=1}^n \lambda^T \mathbf{f}(X_t, Y_{t-1}, Y_t) - \log Z_\lambda(\mathbf{X}_n) \right).$$

Motivation

Intuitively: the **larger** n , the **closer** $\hat{\lambda}_n$ will be to λ^* .

We would like to establish conditions for **consistency** of the MLEs, i.e., under which $\hat{\lambda}_n$ **converges to** λ^* as $n \rightarrow \infty$.

Problem: We don't have a model for infinite sequences!

Solution:

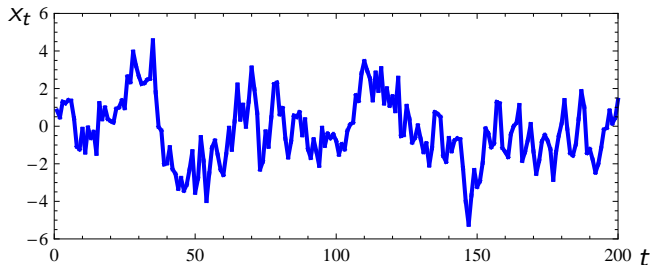
1. Define L-CRFs for double-sided **infinite sequences** $\mathbf{X} = (X_t)_{t \in \mathbb{Z}}$ and $\mathbf{Y} = (Y_t)_{t \in \mathbb{Z}}$.
2. Consider \mathbf{X}_n and \mathbf{Y}_n to be **finite subsequences** of \mathbf{X} and \mathbf{Y} .

Motivation

Example: Suppose that $\mathcal{X} = \mathbb{R}$ and $\mathcal{Y} = \{1, 2\}$. Let

$$\mathbf{f}(x_t, y_{t-1}, y_t) = \begin{bmatrix} \mathbf{1}(y_t = 1) x_t \\ \mathbf{1}(y_t = 2) x_t \\ \mathbf{1}(y_{t-1} = 1, y_t = 1) \\ \mathbf{1}(y_{t-1} = 2, y_t = 2) \end{bmatrix} \quad \text{and} \quad \boldsymbol{\lambda}^* = \begin{bmatrix} 1 \\ -1 \\ 2 \\ 2 \end{bmatrix}.$$

Here \mathbf{X} is an AR(1) process: $X_t = 0.75 X_{t-1} + \epsilon_t$ with $\epsilon_t \sim N(0, 1)$.

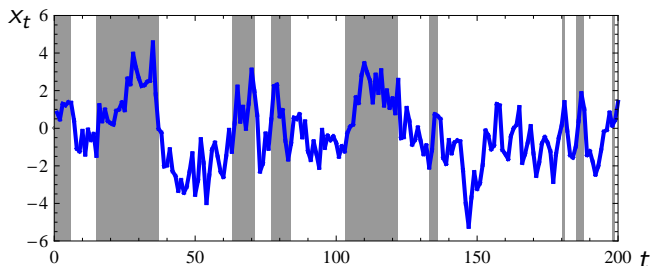


Motivation

Example: Suppose that $\mathcal{X} = \mathbb{R}$ and $\mathcal{Y} = \{1, 2\}$. Let

$$\mathbf{f}(x_t, y_{t-1}, y_t) = \begin{bmatrix} \mathbf{1}(y_t = 1) x_t \\ \mathbf{1}(y_t = 2) x_t \\ \mathbf{1}(y_{t-1} = 1, y_t = 1) \\ \mathbf{1}(y_{t-1} = 2, y_t = 2) \end{bmatrix} \quad \text{and} \quad \boldsymbol{\lambda}^* = \begin{bmatrix} 1 \\ -1 \\ 2 \\ 2 \end{bmatrix}.$$

\mathbf{Y} conditional on \mathbf{X} follows an L-CRF for infinite sequences.

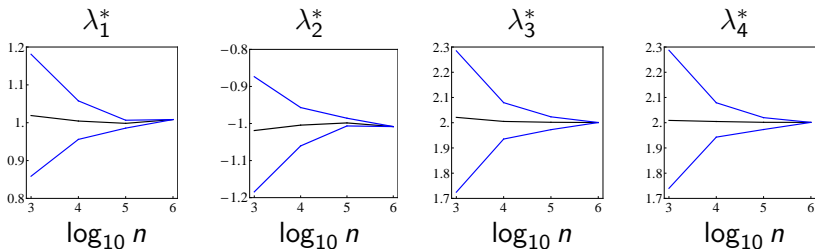


Motivation

Based on this **subsequence**, we obtain the MLE

$$\hat{\lambda}_n = \begin{bmatrix} 0.7 \\ -0.7 \\ 2.1 \\ 2.2 \end{bmatrix} \quad \text{of} \quad \lambda^* = \begin{bmatrix} 1 \\ -1 \\ 2 \\ 2 \end{bmatrix}.$$

As $n \rightarrow \infty$, the MLEs converge to λ^* :



Motivation

What if \mathbf{f} is **unknown**?

- ▶ Use **prior knowledge** to obtain a reasonable “**guess**”
- ▶ Use **parametric approximations**, e.g., linear CRFs

What can we say about the limit of $\hat{\lambda}_n$ in those cases?

Our present work

- ▶ confirms that, if \mathbf{f} is known, “**everything works fine**”
- ▶ provides a **framework** to study the asymptotics of $\hat{\lambda}_n$ in more difficult cases
- ▶ establishes connections between L-CRFs and the theories of **weak ergodicity** and Markov chains in **random environments**

L-CRFs for Infinite Sequences

L-CRFs for Infinite Sequences

Goal Given a sequence $\mathbf{x} = (x_t)_{t \in \mathbb{Z}}$ in \mathcal{X} , define the distribution of \mathbf{Y} conditional on $\mathbf{X} = \mathbf{x}$.

Approach Consider the matrix of **potentials**,

$$\mathbf{M}(\mathbf{x}) = \left[\exp(\boldsymbol{\lambda}^T \mathbf{f}(x, i, j)) \right]_{i, j \in \mathcal{Y}}$$

For $s, t \in \mathbb{Z}$ with $s \leq t$, let

$$\begin{aligned}\alpha_s^t(\mathbf{x}) &= \mathbf{M}(x_t)^T \dots \mathbf{M}(x_s)^T (1, 1, \dots, 1)^T \\ \beta_s^t(\mathbf{x}) &= \mathbf{M}(x_{s+1}) \dots \mathbf{M}(x_t) (1, 1, \dots, 1)^T.\end{aligned}$$

Proposition 1 [Lafferty et al. (2001), Wallach (2004)]

$$\begin{aligned}P_{\boldsymbol{\lambda}}^{(0, n)}(Y_t = y_t, \dots, Y_{t+k} = y_{t+k} \mid \mathbf{X} = \mathbf{x}) \\ = \frac{\alpha_1^t(\mathbf{x}, y_t) \beta_{t+k}^n(\mathbf{x}, y_{t+k})}{\alpha_1^t(\mathbf{x})^T \beta_t^n(\mathbf{x})} \prod_{i=1}^k \exp(\boldsymbol{\lambda}^T \mathbf{f}(x_{t+i}, y_{t+i-1}, y_{t+i})).\end{aligned}$$



L-CRFs for Infinite Sequences

Take into account the **observational context** (X_{-n}, \dots, X_n) :

$$P_{\lambda}^{(-n,n)}(Y_t = y_t, \dots, Y_{t+k} = y_{t+k} \mid \mathbf{X} = \mathbf{x}) \\ := \frac{\alpha_{-n}^t(\mathbf{x}, y_t) \beta_{t+k}^n(\mathbf{x}, y_{t+k})}{\alpha_{-n}^t(\mathbf{x})^T \beta_t^n(\mathbf{x})} \prod_{i=1}^k \exp(\lambda^T \mathbf{f}(x_{t+i}, y_{t+i-1}, y_{t+i})).$$

Theorem 1 Suppose the features \mathbf{f} are **bounded** and the weights λ are **finite**. Then the following limit is well-defined:

$$P_{\lambda}(Y_t = y_t, \dots, Y_{t+k} = y_{t+k} \mid \mathbf{X} = \mathbf{x}) \\ := \lim_{n \rightarrow \infty} P_{\lambda}^{(-n,n)}(Y_t = y_t, \dots, Y_{t+k} = y_{t+k} \mid \mathbf{X} = \mathbf{x}).$$

Moreover, there exists a $\kappa \in (0, 1)$ (not depending on \mathbf{x}) such that

$$\left| P_{\lambda}^{(-n,n)}(Y_t = y_t, \dots, Y_{t+k} = y_{t+k} \mid \mathbf{X} = \mathbf{x}) \right. \\ \left. - P_{\lambda}(Y_t = y_t, \dots, Y_{t+k} = y_{t+k} \mid \mathbf{X} = \mathbf{x}) \right| = O(\kappa^n).$$



L-CRFs for Infinite Sequences

Proof of Theorem 1 Recall the definition

$$P_{\lambda}^{(-n,n)}(Y_t = y_t, \dots, Y_{t+k} = y_{t+k} \mid \mathbf{X} = \mathbf{x}) \\ := \frac{\alpha_{-n}^t(\mathbf{x}, y_t) \beta_{t+k}^n(\mathbf{x}, y_{t+k})}{\alpha_{-n}^t(\mathbf{x})^T \beta_t^n(\mathbf{x})} \prod_{i=1}^k \exp(\lambda^T \mathbf{f}(x_{t+i}, y_{t+i-1}, y_{t+i})).$$

The rows of the matrices

$$\mathbf{M}(x_t)^T \dots \mathbf{M}(x_{-n})^T \quad \text{and} \quad \mathbf{M}(x_{t+1}) \dots \mathbf{M}(x_n)$$

tend to **proportionality** as $n \rightarrow \infty$ [Seneta, 2006]. Hence

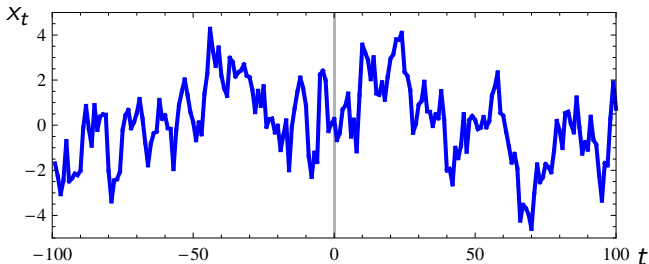
$$\lim_{n \rightarrow \infty} \frac{\alpha_{-n}^t(\mathbf{x}, i)}{\alpha_{-n}^t(\mathbf{x}, j)} = r_{ij} \quad \text{and} \quad \lim_{n \rightarrow \infty} \frac{\beta_{t+k}^n(\mathbf{x}, i)}{\beta_{t+k}^n(\mathbf{x}, j)} = s_{ij}.$$

Kolmogorov's extension theorem: The collection of all marginal distributions specifies the distribution of \mathbf{Y} conditional on $\mathbf{X} = \mathbf{x}$.

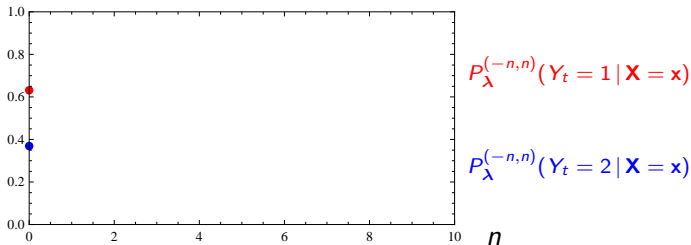


L-CRFs for Infinite Sequences

Example: Consider the distribution of Y_t at $t = 0$:

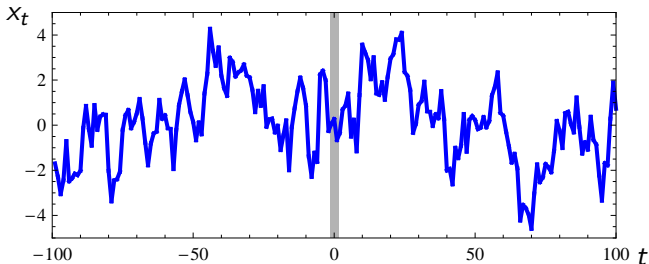


As n increases, the distribution **converges**:

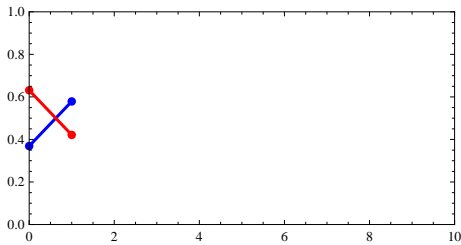


L-CRFs for Infinite Sequences

Example: Consider the distribution of Y_t at $t = 0$:



As n increases, the distribution **converges**:

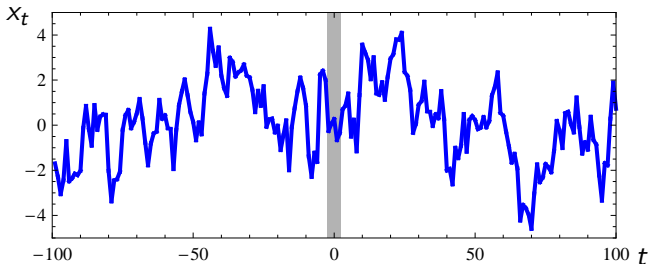


$$P_{\lambda}^{(-n,n)}(Y_t = 1 | \mathbf{X} = \mathbf{x})$$

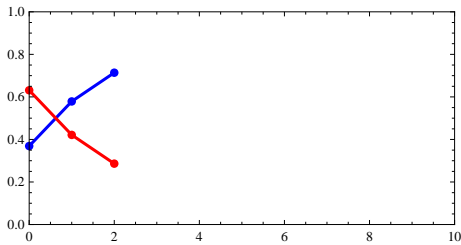
$$P_{\lambda}^{(-n,n)}(Y_t = 2 | \mathbf{X} = \mathbf{x})$$

L-CRFs for Infinite Sequences

Example: Consider the distribution of Y_t at $t = 0$:



As n increases, the distribution **converges**:

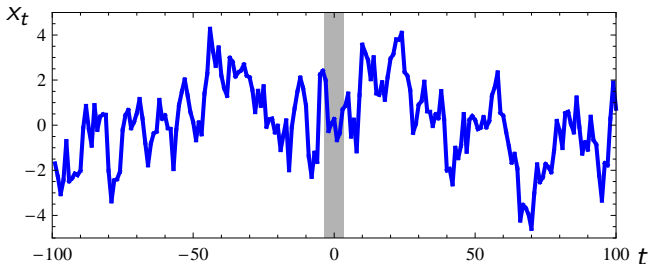


$$P_{\lambda}^{(-n,n)}(Y_t = 1 | \mathbf{X} = \mathbf{x})$$

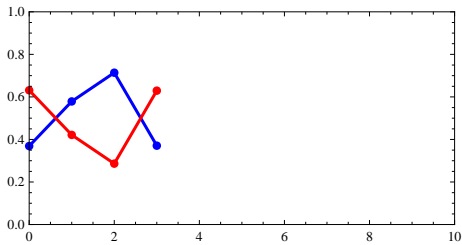
$$P_{\lambda}^{(-n,n)}(Y_t = 2 | \mathbf{X} = \mathbf{x})$$

L-CRFs for Infinite Sequences

Example: Consider the distribution of Y_t at $t = 0$:



As n increases, the distribution **converges**:

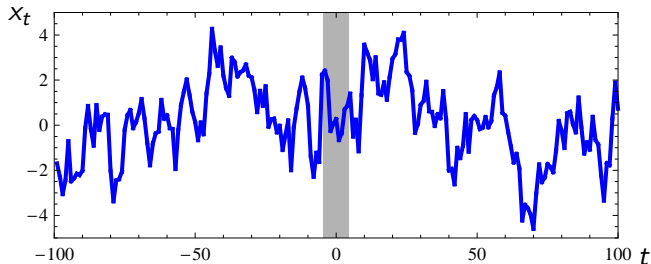


$$P_{\lambda}^{(-n,n)}(Y_t = 1 | \mathbf{X} = \mathbf{x})$$

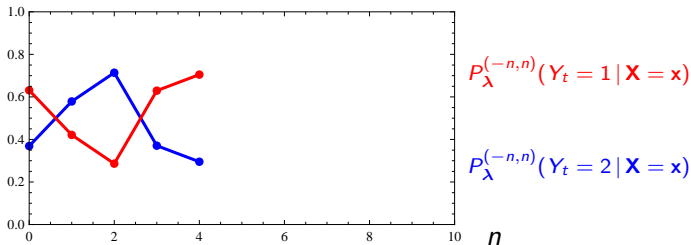
$$P_{\lambda}^{(-n,n)}(Y_t = 2 | \mathbf{X} = \mathbf{x})$$

L-CRFs for Infinite Sequences

Example: Consider the distribution of Y_t at $t = 0$:

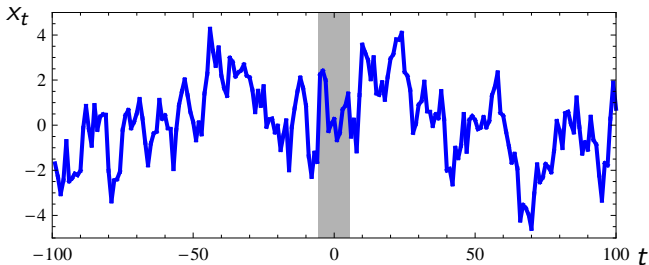


As n increases, the distribution **converges**:

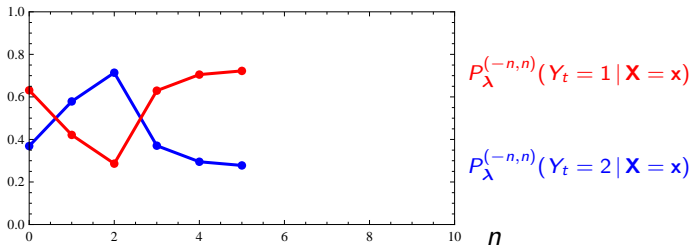


L-CRFs for Infinite Sequences

Example: Consider the distribution of Y_t at $t = 0$:

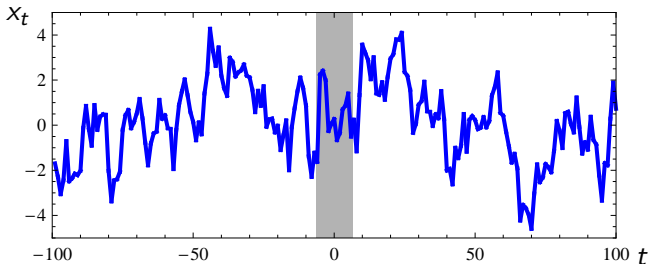


As n increases, the distribution **converges**:

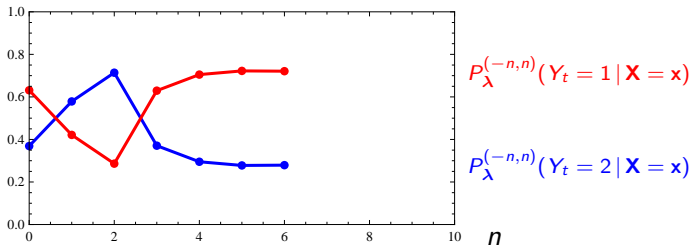


L-CRFs for Infinite Sequences

Example: Consider the distribution of Y_t at $t = 0$:

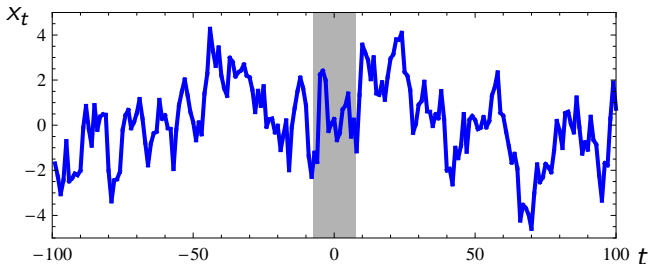


As n increases, the distribution **converges**:

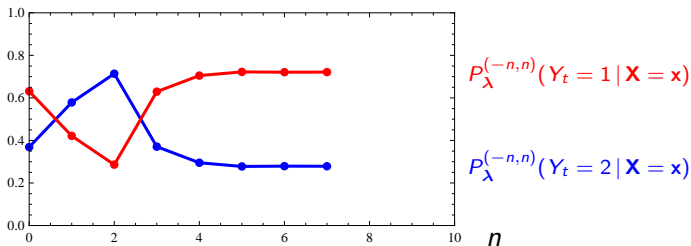


L-CRFs for Infinite Sequences

Example: Consider the distribution of Y_t at $t = 0$:

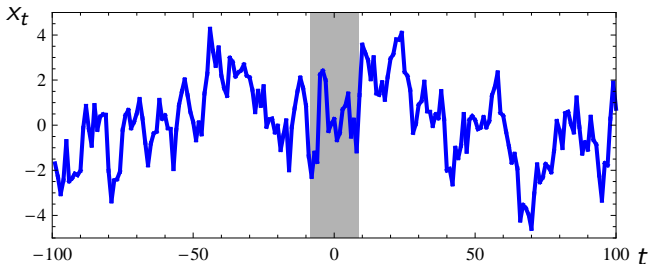


As n increases, the distribution **converges**:

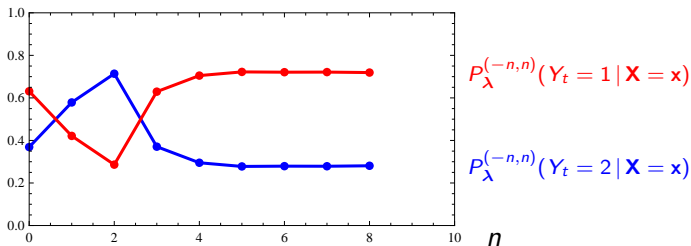


L-CRFs for Infinite Sequences

Example: Consider the distribution of Y_t at $t = 0$:

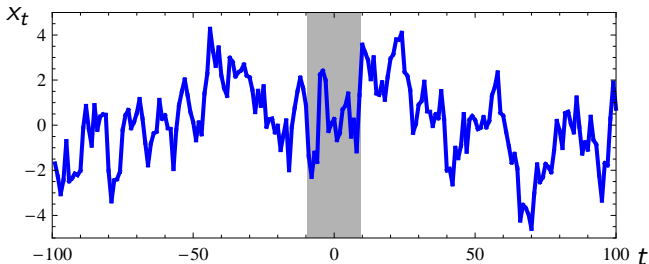


As n increases, the distribution **converges**:

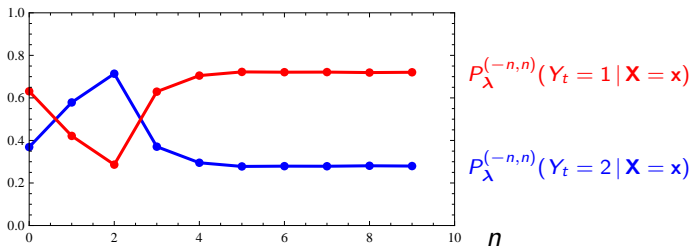


L-CRFs for Infinite Sequences

Example: Consider the distribution of Y_t at $t = 0$:

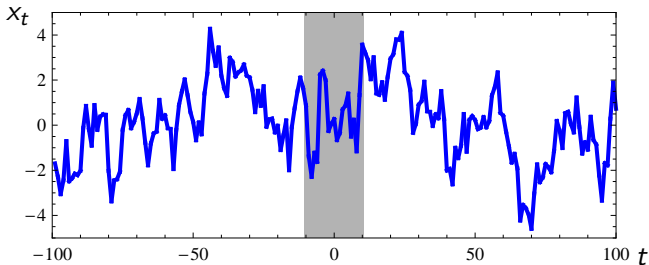


As n increases, the distribution **converges**:

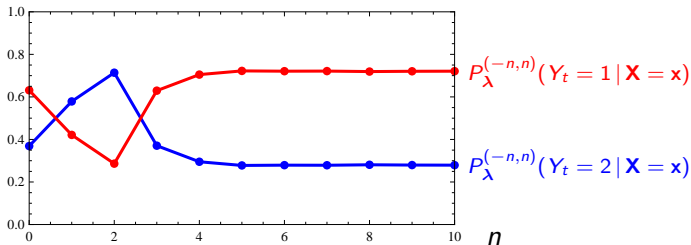


L-CRFs for Infinite Sequences

Example: Consider the distribution of Y_t at $t = 0$:



As n increases, the distribution **converges**:



L-CRFs for Infinite Sequences

Corollary 1 Suppose the assumptions of Theorem 1 hold.

(i) \mathbf{Y} conditional on $\mathbf{X} = \mathbf{x}$ is a **Markov chain** with

$$\begin{aligned} P_{\lambda}(Y_{t+1} = y_{t+1} \mid Y_t = y_t, \mathbf{X} = \mathbf{x}) \\ = \exp(\boldsymbol{\lambda}^T \mathbf{f}(x_{t+1}, y_t, y_{t+1})) \lim_{n \rightarrow \infty} \frac{\beta_{t+1}^n(\mathbf{x}, y_{t+1})}{\beta_t^n(\mathbf{x}, y_t)}. \end{aligned}$$

In particular, $P_{\lambda}(Y_{t+1} = y_{t+1} \mid Y_t = y_t, \mathbf{X} = \mathbf{x}) \geq c$ for some $c > 0$ not depending on \mathbf{x} .

- (ii) If \mathbf{X} is stationary, then the joint sequence (\mathbf{X}, \mathbf{Y}) is **stationary**.
- (iii) If \mathbf{X} is ergodic, then (\mathbf{X}, \mathbf{Y}) is **ergodic**. In particular, for any $g : \mathcal{X} \times \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ satisfying $E_{\lambda}|g(X_t, Y_{t-1}, Y_t)| < \infty$,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n g(X_t, Y_{t-1}, Y_t) = E_{\lambda}[g(X_t, Y_{t-1}, Y_t)] \quad P_{\lambda}\text{-a.s.}$$



Consistency of the MLEs

Consistency of the MLEs

Suppose we observe the **finite subsequences** $\mathbf{X}_n = (X_0, \dots, X_n)$ and $\mathbf{Y}_n = (Y_0, \dots, Y_n)$ of \mathbf{X} and \mathbf{Y} where

- ▶ the distribution of \mathbf{X} is **unknown**
- ▶ the distribution of \mathbf{Y} conditional on \mathbf{X} follows an L-CRF for infinite sequences with **known** features \mathbf{f} and **unknown** weights λ^* .

In order to **estimate** λ^* ,

$$\text{compute } \hat{\lambda}_n = \arg \max_{\lambda \in \Theta} \mathcal{L}_n(\lambda)$$

$$\text{where } \mathcal{L}_n(\lambda) = \frac{1}{n} \left(\sum_{t=1}^n \lambda^T \mathbf{f}(X_t, Y_{t-1}, Y_t) - \log Z_\lambda(\mathbf{X}_n) \right).$$

Under which conditions $\lim_{n \rightarrow \infty} \hat{\lambda}_n = \lambda^*$ P_{λ^*} -a.s.?

Consistency of the MLEs

Lemma 1 Suppose that (A1) the features \mathbf{f} are **bounded** and the weights $\boldsymbol{\lambda}$ are **finite**, (A2) the process \mathbf{X} is **ergodic**, (A3) the parameter space Θ is **compact**. Then the following holds:

(i) There exists a function $\mathcal{L}(\boldsymbol{\lambda})$ such that, for every $\boldsymbol{\lambda} \in \Theta$,

$$\lim_{n \rightarrow \infty} \mathcal{L}_n(\boldsymbol{\lambda}) = \mathcal{L}(\boldsymbol{\lambda}) \quad P_{\boldsymbol{\lambda}^*}\text{-a.s.}$$

(ii) The convergence of $\mathcal{L}_n(\boldsymbol{\lambda})$ to $\mathcal{L}(\boldsymbol{\lambda})$ is **uniform** on Θ .

(iii) The limit of $\nabla^2 \mathcal{L}_n(\boldsymbol{\lambda})$ is **finite**:

$$\lim_{n \rightarrow \infty} \nabla^2 \mathcal{L}_n(\boldsymbol{\lambda}) = - \left(\gamma_{\boldsymbol{\lambda}}(0) + 2 \sum_{k=1}^{\infty} \gamma_{\boldsymbol{\lambda}}(k) \right) \quad P_{\boldsymbol{\lambda}^*}\text{-a.s.}$$

where $\gamma_{\boldsymbol{\lambda}}(k) = \text{Cov}_{\boldsymbol{\lambda}}[\mathbf{f}(X_t, Y_{t-1}, Y_t), \mathbf{f}(X_{t+k}, Y_{t+k-1}, Y_{t+k})]$.



Consistency of the MLEs

Theorem 2 Suppose that (A1)-(A3) hold. Moreover, suppose that the limit of $\nabla^2 \mathcal{L}_n(\boldsymbol{\lambda})$ is **non-singular** for every $\boldsymbol{\lambda} \in \Theta$. Then $\mathcal{L}(\boldsymbol{\lambda})$ is **strictly concave** and hence

$$\lim_{n \rightarrow \infty} \hat{\boldsymbol{\lambda}}_n = \boldsymbol{\lambda}^* \quad P_{\boldsymbol{\lambda}^*}\text{-a.s.}$$

What if the limit of $\nabla^2 \mathcal{L}_n(\boldsymbol{\lambda})$ is **singular**?

Example: Suppose that \mathbf{f} contains two **identical features** f_i, f_j . Then $\mathcal{L}(\boldsymbol{\lambda}) = \mathcal{L}(\boldsymbol{\lambda}^*)$ for any $\boldsymbol{\lambda} \in \Theta$ with $\lambda_i + \lambda_j = \lambda_i^* + \lambda_j^*$ and $\lambda_k = \lambda_k^*$ for all other k .

Discussion

Discussion

Open questions for **future research**:

- ▶ What are **sufficient conditions** for non-singularity of the limit of $\nabla^2 \mathcal{L}(\lambda)$ and hence **identifiability** of λ^* ?
- ▶ Can we weaken the assumption of **boundedness** of \mathbf{f} ?
- ▶ Straightforward: generalization to L-CRFs of **higher orders**.
- ▶ What can we say about the quality of the **approximations**, e.g., by linear CRFs, when the true features are **unknown**?
- ▶ In the context of **human activity recognition**: what features should we use to obtain models that are **robust** towards **individual differences**?

References

- [1] R. Cogburn (1984). The ergodic theory of Markov chains in random environments. *Z. Wahrscheinlichkeitstheorie verw. Gebiete* 66:109-128.
- [2] O. Hernández-Lerma and J.B. Lasserre (2003). *Markov chains and invariant probabilities*. Basel, Switzerland: Birkhäuser.
- [3] J. Lafferty, A. McCallum and F.C.N. Pereira (2001). Conditional random fields: probabilistic models for segmenting and labeling sequence data. *Proceedings of the IEEE International Conference on Machine Learning (ICML)*.
- [4] E. Seneta (2006). *Non-negative matrices and Markov chains. Revised Edition*. New York, NY: Springer.
- [5] H. Wallach (2004). Conditional random fields: an introduction. *Technical Report MS-CIS-04-21*. University of Pennsylvania, Philadelphia.