
Asymptotic Theory for Linear-Chain Conditional Random Fields

Mathieu Sinn
University of Waterloo

Pascal Poupart
University of Waterloo

Abstract

In this theoretical paper we develop an asymptotic theory for Linear-Chain Conditional Random Fields (L-CRFs) and apply it to derive conditions under which the Maximum Likelihood Estimates (MLEs) of the model weights are strongly consistent. We first define L-CRFs for infinite sequences and analyze some of their basic properties. Then we establish conditions under which ergodicity of the observations implies ergodicity of the joint sequence of observations and labels. This result is the key ingredient to derive conditions for strong consistency of the MLEs. Interesting findings are that the consistency crucially depends on the limit behavior of the Hessian of the likelihood function and that, asymptotically, the state feature functions do not matter.

1 INTRODUCTION

Conditional Random Fields (CRFs) are a widely popular model to describe the statistical dependence between sequences of “observations” and “labels” (Lafferty et al., 2001). Applications include natural language processing (Sutton and McCallum, 2006), the analysis of genome data (Culotta et al., 2005), and human activity recognition (Omar et al., 2010). Extensions of CRFs are hierarchical CRFs (Liao et al. 2007b), relational CRFs (Taskar et al. 2002), and semi-Markov CRFs (Sarawagi and Cohen 2004). The key idea of CRFs is to represent the distribution of the labels conditional on the observations by a Markov random field. The simplest non-trivial class of models is that of Linear-Chain CRFs (L-CRFs) where the random field forms a plain Markov chain.

Appearing in Proceedings of the 14th International Conference on Artificial Intelligence and Statistics (AISTATS) 2011, Fort Lauderdale, FL, USA. Volume 15 of JMLR: W&CP 15. Copyright 2011 by the authors.

In this paper, we study asymptotical properties of the Maximum Likelihood Estimates (MLEs) for the model weights. More specifically, we assume that we are given a sequence of observations and labels where the distribution of the labels conditional on the observations follows an L-CRF with known feature functions but unknown weights. In this setting, we investigate conditions under which the MLEs converge to the true weights as the length of the sequences goes to infinity. Note that, to state and to analyze this problem, a definition of L-CRFs for infinite sequences is required.

Our research is motivated by the following questions: How can the weights and feature functions be jointly estimated in the case where both are unknown? How robust is the training and inference towards a sampling bias (that is, when training and test data come from different distributions)? How well is the model identifiable in the presence of noisy data? To tackle these problems of great practical importance, the present paper aims to achieve a better understanding of the simplest case, namely, when the feature functions are known and a sampling bias or noisy data is absent. Furthermore, it provides a theoretical framework and useful techniques to study the more complicated cases.

This paper is structured as follows: In Sec. 2 we introduce some notation and review the definition of L-CRFs for finite sequences. In Sec. 3 we define L-CRFs for infinite sequences and derive some of their basic properties. Sec. 4 establishes conditions under which ergodicity of the sequence of observations implies ergodicity of the joint sequence of observations and labels. In Sec. 5 we apply the previous results to derive conditions under which the MLEs are strongly consistent. Sec. 6 concludes the paper.

2 PRELIMINARIES

Throughout this paper, let \mathbb{N} , \mathbb{Z} , \mathbb{R} denote the sets of natural numbers, integers and real numbers, respectively. Let \mathcal{X} be a metric space with the Borel sigma-field \mathcal{A} . To fix ideas, think of \mathcal{X} as \mathbb{R}^d equipped with the Euclidean norm. Furthermore, consider a finite set \mathcal{Y} and let \mathcal{B} denote the power set of \mathcal{Y} . In the following

we suppose that $\mathcal{Y} = \{1, \dots, \ell\}$ for some $\ell \in \mathbb{N}$. Now let (Ω, \mathcal{F}) be a measurable space, and let $\mathbf{X} = (X_t)_{t \in \mathbb{Z}}$ and $\mathbf{Y} = (Y_t)_{t \in \mathbb{Z}}$ be sequences of measurable mappings from (Ω, \mathcal{F}) into \mathcal{X} and \mathcal{Y} , respectively. We refer to \mathbf{X} and \mathbf{Y} as the sequences of *observations* and *labels*, respectively.

A conventional Linear-Chain Conditional Random Field (L-CRF) specifies the conditional distribution of $\mathbf{Y}_n = (Y_0, \dots, Y_n)$ given $\mathbf{X}_n = (X_0, \dots, X_n)$. The distributions are parameterized as follows: Let $\mathbf{f}^{\text{state}}$ and $\mathbf{f}^{\text{trans}}$ be vectors of \mathbb{R} -valued functions defined on $\mathcal{X} \times \mathcal{Y}$ and $\mathcal{X} \times \mathcal{Y}^2$, respectively, and $\boldsymbol{\mu}, \boldsymbol{\nu}$ be \mathbb{R} -valued vectors of corresponding lengths. We call $\mathbf{f}^{\text{state}}$ and $\mathbf{f}^{\text{trans}}$ the *state* and *transition* feature functions, and $\boldsymbol{\mu}, \boldsymbol{\nu}$ the *model weights*. Write $\boldsymbol{\lambda}$ for the stacked vector $(\boldsymbol{\mu}, \boldsymbol{\nu})$, and let

$$\begin{aligned} \boldsymbol{\lambda}^T \mathbf{f}(\mathbf{x}_n, \mathbf{y}_n) &= \sum_{t=0}^n \boldsymbol{\mu}^T \mathbf{f}^{\text{state}}(x_t, y_t) \\ &+ \sum_{t=1}^n \boldsymbol{\nu}^T \mathbf{f}^{\text{trans}}(x_t, y_{t-1}, y_t) \end{aligned} \quad (1)$$

for $\mathbf{x}_n = (x_0, \dots, x_n) \in \mathcal{X}^{n+1}$ and $\mathbf{y}_n = (y_0, \dots, y_n) \in \mathcal{Y}^{n+1}$. We make three remarks at this point: First, the feature functions could also depend on t ; here we restrict ourselves to models where the dependencies between observations and labels do not vary with time. Second, the assumption that the feature functions only depend on the current observation x_t is without loss of generality; if they depended on, say, x_{t-l}, \dots, x_{t+l} , then simply consider the modified observations $\tilde{x}_t = (x_{t-l}, \dots, x_{t+l})$ instead of x_t . Third, equation (1) suggests that there might be some overlap between the state and transition feature functions; this is indeed the case and will play an important role later.

Now let $\mathbf{x} = (x_t)_{t \in \mathbb{Z}}$ be a sequence in \mathcal{X} and write \mathbf{x}_n to denote the projection of \mathbf{x} onto the components (x_0, \dots, x_n) . For any $\mathbf{y}_n = (y_0, \dots, y_n) \in \mathcal{Y}^{n+1}$, define the conditional probability

$$\begin{aligned} P_{\boldsymbol{\lambda}}^{(0:n)}(\mathbf{Y}_n = \mathbf{y}_n \mid \mathbf{X} = \mathbf{x}) \\ = \frac{1}{Z_{\boldsymbol{\lambda}}(\mathbf{x}_n)} \exp(\boldsymbol{\lambda}^T \mathbf{f}(\mathbf{x}_n, \mathbf{y}_n)) \end{aligned} \quad (2)$$

where $Z_{\boldsymbol{\lambda}}(\mathbf{x}_n)$ is the normalizing partition function

$$Z_{\boldsymbol{\lambda}}(\mathbf{x}_n) = \sum_{\mathbf{y}'_n \in \mathcal{Y}^{n+1}} \exp(\boldsymbol{\lambda}^T \mathbf{f}(\mathbf{x}_n, \mathbf{y}'_n)). \quad (3)$$

Note that the conditional probability in (2) only depends on the components x_0, \dots, x_n of \mathbf{x} . Thus, $P_{\boldsymbol{\lambda}}^{(0:n)}$ specifies a conditional distribution for \mathbf{Y}_n given \mathbf{X}_n . In the next section we will define L-CRFs for infinite sequences by considering conditional distributions which depend on the entire sequence \mathbf{x} .

The key step in applying L-CRFs is the specification of the feature functions and the model weights. A common approach is to select the feature functions manually according to some expert domain knowledge; alternative ways are to iteratively select those functions from a set of candidates which give the largest increase of the conditional log-likelihood (McCallum, 2003), or to use virtual evidence boosting (Liao et al. 2007a). Given a fixed set of feature functions, the model weights are usually learned from labeled training data by maximizing the conditional log-likelihood. In Sec. 5 we suppose that the “true” feature functions are known, and consider asymptotic properties of the Maximum Likelihood estimates.

In the remaining part of this section we introduce an alternative representation for the right hand side of (2), which is fundamental to all that follows. Consider the function $\alpha_{\boldsymbol{\lambda}} : \mathcal{X} \rightarrow \mathbb{R}^{\ell}$ with the i th component

$$\alpha_{\boldsymbol{\lambda}}(x, i) = \exp(\boldsymbol{\mu}^T \mathbf{f}^{\text{state}}(x, i)).$$

Furthermore, let $M_{\boldsymbol{\lambda}}(x)$ be the $\ell \times \ell$ -matrix with the (i, j) -th component

$$m_{\boldsymbol{\lambda}}(x, i, j) = \exp(\boldsymbol{\mu}^T \mathbf{f}^{\text{state}}(x, j) + \boldsymbol{\nu}^T \mathbf{f}^{\text{trans}}(x, i, j)).$$

In the terminology of Markov Random Fields, $\alpha_{\boldsymbol{\lambda}}(x_t, y_t)$ is the *potential* of the event $X_t = x_t$ and $Y_t = y_t$ (with respect to the model weights $\boldsymbol{\lambda}$), while $m_{\boldsymbol{\lambda}}(x_t, y_{t-1}, y_t)$ is the potential of the event $X_t = x_t$, $Y_{t-1} = y_{t-1}$ and $Y_t = y_t$. For a sequence $\mathbf{x} = (x_t)_{t \in \mathbb{Z}}$ in \mathcal{X} and $s, t \in \mathbb{Z}$ with $s \leq t$, define

$$\begin{aligned} \boldsymbol{\alpha}_s^t(\boldsymbol{\lambda}, \mathbf{x}) &= M_{\boldsymbol{\lambda}}(x_t)^T \dots M_{\boldsymbol{\lambda}}(x_s)^T \boldsymbol{\alpha}_{\boldsymbol{\lambda}}(x_{s-1}), \\ \boldsymbol{\beta}_s^t(\boldsymbol{\lambda}, \mathbf{x}) &= M_{\boldsymbol{\lambda}}(x_{s+1}) \dots M_{\boldsymbol{\lambda}}(x_t) (1, 1, \dots, 1)^T. \end{aligned}$$

Write $\alpha_s^t(\boldsymbol{\lambda}, \mathbf{x}, i)$ and $\beta_s^t(\boldsymbol{\lambda}, \mathbf{x}, j)$ to denote the i th and j th component of $\boldsymbol{\alpha}_s^t(\boldsymbol{\lambda}, \mathbf{x})$ and $\boldsymbol{\beta}_s^t(\boldsymbol{\lambda}, \mathbf{x})$. Note that $\alpha_s^t(\boldsymbol{\lambda}, \mathbf{x}, y_t)$ is the potential of the event $X_s = x_s, \dots, X_t = x_t$ and $Y_t = y_t$, while $\beta_s^t(\boldsymbol{\lambda}, \mathbf{x}, y_s)$ is the potential of the event $X_s = x_s, \dots, X_t = x_t$ and $Y_s = y_s$. The following is well-known (Wallach, 2004):

Proposition 1. *Let $t, k \in \mathbb{N}$ such that $t + k \leq n$. Then, for all $y_t, \dots, y_{t+k} \in \mathcal{Y}$,*

$$\begin{aligned} P_{\boldsymbol{\lambda}}^{(0:n)}(Y_t = y_t, \dots, Y_{t+k} = y_{t+k} \mid \mathbf{X} = \mathbf{x}) \\ = \frac{\alpha_1^t(\boldsymbol{\lambda}, \mathbf{x}, y_t) \beta_{t+k}^n(\boldsymbol{\lambda}, \mathbf{x}, y_{t+k})}{\boldsymbol{\alpha}_1^t(\boldsymbol{\lambda}, \mathbf{x})^T \boldsymbol{\beta}_t^n(\boldsymbol{\lambda}, \mathbf{x})} \\ \times \prod_{i=1}^k m_{\boldsymbol{\lambda}}(x_{t+i}, y_{t+i-1}, y_{t+i}) \end{aligned}$$

where, as usual, products over empty index sets are equal to 1.

3 L-CRFs FOR INFINITE SEQUENCES

In this section we define L-CRFs for infinite sequences and derive some of their basic properties. We will assume that the following condition is satisfied:

(A1) The feature functions $\mathbf{f}^{\text{state}}$ and $\mathbf{f}^{\text{trans}}$ are bounded, and the model weights $\boldsymbol{\lambda}$ are finite.

The next lemma states some obvious consequence of assumption (A1) which will be needed below. Let m_{inf} and m_{sup} denote the infimum and the supremum over all $m_{\boldsymbol{\lambda}}(x, i, j)$ with $x \in \mathcal{X}$ and $i, j \in \mathcal{Y}$.

Lemma 1. *Suppose that (A1) holds. Then $m_{\text{inf}} > 0$ and $m_{\text{sup}} < \infty$. In particular, the quantities φ^2 and ψ^2 defined by*

$$\begin{aligned} \varphi^2 &:= \inf \left\{ \min_{i,j,k,l \in \mathcal{Y}} \frac{m(x, k, i) m(x, l, j)}{m(x, k, j) m(x, l, i)} : x \in \mathcal{X} \right\}, \\ \psi^2 &:= \inf \left\{ \min_{i,j,k,l \in \mathcal{Y}} \frac{m(x, i, k) m(x, j, l)}{m(x, j, k) m(x, i, l)} : x \in \mathcal{X} \right\}, \end{aligned}$$

are both strictly greater than zero.

3.1 Definition

Let $\mathbf{x} = (x_t)_{t \in \mathbb{Z}}$ be fixed. Our goal is to define the distribution of the infinite sequence \mathbf{Y} conditional on $\mathbf{X} = \mathbf{x}$. Let $t \in \mathbb{Z}$ and $k \in \mathbb{N}$. We first specify the conditional marginal distribution of (Y_t, \dots, Y_{t+k}) . Then, by applying Kolmogorov's extension theorem, we obtain the conditional distribution for the entire sequence \mathbf{Y} . Let $y_t, \dots, y_{t+k} \in \mathcal{Y}$. For $n \in \mathbb{N}$ such that $-n \leq t$ and $n \leq t+k$, define

$$\begin{aligned} P_{\boldsymbol{\lambda}}^{(-n, n)}(Y_t = y_t, \dots, Y_{t+k} = y_{t+k} | \mathbf{X} = \mathbf{x}) \\ &:= \frac{\alpha_{-n}^t(\boldsymbol{\lambda}, \mathbf{x}, y_t) \beta_{t+k}^n(\boldsymbol{\lambda}, \mathbf{x}, y_{t+k})}{\alpha_{-n}^t(\boldsymbol{\lambda}, \mathbf{x})^T \beta_t^n(\boldsymbol{\lambda}, \mathbf{x})} \\ &\quad \times \prod_{i=1}^k m_{\boldsymbol{\lambda}}(x_{t+i}, y_{t+i-1}, y_{t+i}). \end{aligned}$$

Comparing this to Proposition 1 we see that $P_{\boldsymbol{\lambda}}^{(-n, n)}$ specifies the conditional distribution of (Y_t, \dots, Y_{t+k}) in the fashion of a conventional L-CRF considering the finite observational context (X_{-n}, \dots, X_n) . Theorem 2 below shows that this distribution converges as n tends to infinity. For the proof, we will need the following result.

Theorem 1. *Suppose that (A1) holds. For a fixed sequence $(x_t)_{t \in \mathbb{N}}$ in \mathcal{X} , let $\mathbf{M}_n = (m_n(i, j))_{i, j \in \mathcal{Y}}$ be given by $\mathbf{M}_n = \mathbf{M}_{\boldsymbol{\lambda}}(x_1) \dots \mathbf{M}_{\boldsymbol{\lambda}}(x_n)$. Then there exist positive numbers r_{ij} such that, for all $i, j, k \in \mathcal{Y}$,*

$$\lim_{n \rightarrow \infty} \frac{m_n(i, k)}{m_n(j, k)} = r_{ij},$$

i.e., the rows of \mathbf{M}_n tend to proportionality as n tends to infinity. Moreover,

$$\max_{k \in \mathcal{Y}} \left(\frac{m_n(i, k)}{m_n(j, k)} \right) \geq r_{ij} \geq \min_{k \in \mathcal{Y}} \left(\frac{m_n(i, k)}{m_n(j, k)} \right)$$

for all $n \in \mathbb{N}$, and

$$\min_{k \in \mathcal{Y}} \left(\frac{m_n(i, k)}{m_n(j, k)} \right) \geq \left[1 - 4 \left(\frac{1 - \psi}{1 + \psi} \right)^n \right] \max_{k \in \mathcal{Y}} \left(\frac{m_n(i, k)}{m_n(j, k)} \right)$$

with ψ given in Lemma 1. The same results apply, with φ instead of ψ , for $\mathbf{M}_n = \mathbf{M}_{\boldsymbol{\lambda}}(x_1)^T \dots \mathbf{M}_{\boldsymbol{\lambda}}(x_n)^T$.

Proof. The proof, which uses well-known results from the theory of weak ergodicity (Seneta, 2006), is included in the supplementary material. \square

Now we are prepared to establish the following result.

Theorem 2. *Suppose that (A1) holds. Then the following limit is well-defined:*

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

Moreover, there exist constants $c > 0$ and $0 < \kappa < 1$, which do not depend on \mathbf{x} , such that

$$\begin{aligned} |P_{\boldsymbol{\lambda}}^{(-n, n)}(Y_t = y_t, \dots, Y_{t+k} = y_{t+k} | \mathbf{X} = \mathbf{x}) \\ - P_{\boldsymbol{\lambda}}(Y_t = y_t, \dots, Y_{t+k} = y_{t+k} | \mathbf{X} = \mathbf{x})| \leq c\kappa^n. \end{aligned}$$

Proof. Define $\mathbf{G}_n := \mathbf{M}_{\boldsymbol{\lambda}}(x_t)^T \dots \mathbf{M}_{\boldsymbol{\lambda}}(x_{-n})^T$ and $\mathbf{H}_n := \mathbf{M}_{\boldsymbol{\lambda}}(x_{t+k+1}) \dots \mathbf{M}_{\boldsymbol{\lambda}}(x_n)$. Write $g_n(i, j)$ and $h_n(i, j)$ for the (i, j) -th components of \mathbf{G}_n and \mathbf{H}_n . Note that $\alpha_{-n}^t(\boldsymbol{\lambda}, \mathbf{x}) = \mathbf{G}_n \boldsymbol{\alpha}_{-n}(x_{-n-1})$ and $\beta_{t+k}^n(\boldsymbol{\lambda}, \mathbf{x}) = \mathbf{H}_n(1, 1, \dots, 1)^T$. Furthermore, with $\mathbf{F} := \mathbf{M}_{\boldsymbol{\lambda}}(x_{t+1}) \dots \mathbf{M}_{\boldsymbol{\lambda}}(x_{t+k})$, we have

$$\alpha_{-n}^t(\boldsymbol{\lambda}, \mathbf{x})^T \beta_t^n(\boldsymbol{\lambda}, \mathbf{x}) = \alpha_{-n}^t(\boldsymbol{\lambda}, \mathbf{x})^T \mathbf{F} \beta_{t+k}^n(\boldsymbol{\lambda}, \mathbf{x}).$$

According to Theorem 1, there exist numbers r_{ij} , s_{ij} such that

$$\lim_{n \rightarrow \infty} \frac{g_n(i, k)}{g_n(j, k)} = r_{ij} \quad \text{and} \quad \lim_{n \rightarrow \infty} \frac{h_n(i, k)}{h_n(j, k)} = s_{ij}$$

for all $k \in \mathcal{Y}$. Consequently, the ratio of $\alpha_{-n}^t(\boldsymbol{\lambda}, \mathbf{x}, i)$ to $\alpha_{-n}^t(\boldsymbol{\lambda}, \mathbf{x}, j)$ converges to r_{ij} , and the ratio of $\beta_{t+k}^n(\boldsymbol{\lambda}, \mathbf{x}, i)$ to $\beta_{t+k}^n(\boldsymbol{\lambda}, \mathbf{x}, j)$ converges to s_{ij} . Hence,

$$\lim_{n \rightarrow \infty} \frac{\alpha_{-n}^t(\boldsymbol{\lambda}, \mathbf{x}, i) \beta_{t+k}^n(\boldsymbol{\lambda}, \mathbf{x}, j)}{\alpha_{-n}^t(\boldsymbol{\lambda}, \mathbf{x})^T \mathbf{F} \beta_{t+k}^n(\boldsymbol{\lambda}, \mathbf{x})} = \frac{1}{r_i^T \mathbf{F} s_j}$$

where $\mathbf{r}_i = (r_{1i}, \dots, r_{\ell i})^T$ and $\mathbf{s}_j = (s_{1j}, \dots, s_{\ell j})^T$, which proves the first part of the theorem. For the second part, note that one can choose

$$\kappa = \frac{(1-\varphi)(1-\psi)}{(1+\varphi)(1+\psi)}$$

with φ^2 and ψ^2 as defined in Lemma 1, and

$$c = 16 \left(\frac{m_{\text{sup}}}{m_{\text{inf}}} \right)^2 \ell^{k+1} m_{\text{sup}}^{2k}.$$

See the supplementary material for details. \square

It is easy to verify that

$$P_{\lambda}(Y_t = y_t, \dots, Y_{t+k} = y_{t+k} | \mathbf{X} = \mathbf{x}) \geq 0$$

and

$$\sum_{y_t, \dots, y_{t+k} \in \mathcal{Y}} P_{\lambda}(Y_t = y_t, \dots, Y_{t+k} = y_{t+k} | \mathbf{X} = \mathbf{x}) = 1,$$

so the distribution of (Y_t, \dots, Y_{t+k}) conditional on $\mathbf{X} = \mathbf{x}$ is well-defined. Furthermore, the collection of all such marginal distributions with $t \in \mathbb{Z}$ and $k \in \mathbb{N}$ satisfies the consistency conditions of Kolmogorov's extension theorem and hence specifies a unique probability measure on the space of sequences $\mathbf{y} = (y_t)_{t \in \mathbb{Z}}$ with $y_t \in \mathcal{Y}$. In this way, we obtain the distribution of the infinite sequence \mathbf{Y} conditional on $\mathbf{X} = \mathbf{x}$.

Throughout the rest of this paper, let Θ be a set of model weights such that each $\lambda \in \Theta$ satisfies assumption (A1). For simplicity we assume that the distribution of \mathbf{X} does not depend on λ . More precisely, let \mathcal{X} denote the space of sequences $\mathbf{x} = (x_t)_{t \in \mathbb{Z}}$ in \mathcal{X} , and \mathcal{A} be the corresponding product sigma-field. We assume there exists a probability measure $P_{\mathbf{X}}$ on $(\mathcal{X}, \mathcal{A})$ such that $P_{\mathbf{X}}(A) = P_{\lambda}(\mathbf{X} \in A)$ for all $A \in \mathcal{A}$ and $\lambda \in \Theta$.

3.2 Basic Properties

The following corollary summarizes consequences of the previous definition. Statement (ii) says that the probability for the transition between any two labels is always bounded away from zero, regardless of the observational context. This result will be of particular importance in Sec. 4 where we investigate under which conditions ergodicity of \mathbf{X} implies ergodicity of the joint sequence (\mathbf{X}, \mathbf{Y}) .

Corollary 1. *Suppose that condition (A1) holds and hence the distribution of \mathbf{Y} conditional on $\mathbf{X} = \mathbf{x}$ is well-defined for every $\mathbf{x} = (x_t)_{t \in \mathbb{Z}}$.*

- (i) *In the definition of the distribution of \mathbf{Y} conditional on $\mathbf{X} = \mathbf{x}$ we may assume, without loss of generality, that the weighted state feature functions are equal to zero, $\boldsymbol{\mu}^T \tilde{\mathbf{f}}^{\text{state}} = 0$.*

- (ii) *\mathbf{Y} conditional on $\mathbf{X} = \mathbf{x}$ is a Markov process with*

$$\begin{aligned} P_{\lambda}(Y_{t+1} = y_{t+1} | Y_t = y_t, \mathbf{X} = \mathbf{x}) \\ = m_{\lambda}(x_t, y_t, y_{t+1}) \lim_{n \rightarrow \infty} \frac{\beta_{t+1}^n(\lambda, \mathbf{x}, y_{t+1})}{\beta_t^n(\lambda, \mathbf{x}, y_t)}. \end{aligned}$$

Furthermore, for all $y_t, y_{t+1} \in \mathcal{Y}$,

$$P_{\lambda}(Y_{t+1} = y_{t+1} | Y_t = y_t, \mathbf{X} = \mathbf{x}) \geq \frac{1}{\ell} \left(\frac{m_{\text{inf}}}{m_{\text{sup}}} \right)^2.$$

- (iii) *If \mathbf{X} is stationary, then the joint sequence (\mathbf{X}, \mathbf{Y}) is stationary.*

Proof. (i) Consider the state and transition feature functions $\tilde{\mathbf{f}}^{\text{state}}$, $\tilde{\mathbf{f}}^{\text{trans}}$ and the weights $\tilde{\boldsymbol{\mu}}$, $\tilde{\boldsymbol{\nu}}$ given in such a way that $\tilde{\boldsymbol{\mu}}^T \tilde{\mathbf{f}}^{\text{state}} = 0$ and

$$\begin{aligned} \tilde{\boldsymbol{\nu}}^T \tilde{\mathbf{f}}^{\text{trans}}(x_t, y_{t-1}, y_t) &= \boldsymbol{\mu}^T \mathbf{f}^{\text{state}}(x_t, y_t) \\ &+ \boldsymbol{\nu}^T \mathbf{f}^{\text{trans}}(x_t, y_{t-1}, y_t). \end{aligned}$$

It is easy to see that the resulting distributions for \mathbf{Y} conditional on $\mathbf{X} = \mathbf{x}$ are identical; in particular, the limit in Theorem 2 does not depend on the values of the vectors $\boldsymbol{\alpha}_{\lambda}(x_{t-n})$. (ii) The equality is directly obtained by the definition of the conditional distribution of \mathbf{Y} , and it is easily verified that it holds for any $P_{\lambda}(Y_{t+1} = y_{t+1} | Y_t = y_t, \dots, Y_{t-k} = y_{t-k}, \mathbf{X} = \mathbf{x})$ with $k \in \mathbb{N}$. To establish the lower bound note that, according to Theorem 1,

$$\frac{\beta_{t+1}^n(\lambda, \mathbf{x}, y_{t+1})}{\beta_{t+1}^n(\lambda, \mathbf{x}, y_t)} \geq \frac{m_{\text{inf}}}{m_{\text{sup}}}$$

for all $\mathbf{x} \in \mathcal{X}$ and $n \in \mathbb{N}$. Since $\beta_t^n(\lambda, \mathbf{x}) = M_{\lambda}(x_{t+1}) \beta_{t+1}^n(\lambda, \mathbf{x})$, the result follows. (iii) The stationarity of the joint sequence is obvious because

$$\begin{aligned} P_{\lambda}(\mathbf{X} \in \mathbf{A}, \mathbf{Y} \in \mathbf{B}) \\ = \int_{\mathbf{A}} P_{\lambda}(\mathbf{Y} \in \mathbf{B} | \mathbf{X} = \mathbf{x}) P_{\mathbf{X}}(d\mathbf{x}) \end{aligned}$$

for any measurable events \mathbf{A} and \mathbf{B} , and by the stationarity of \mathbf{X} the integral on the right hand side is invariant with respect to time shifts of \mathbf{X} and \mathbf{Y} in the integrand. \square

Throughout the rest of this paper, any expected value such as $E_{\lambda}[g(X_t, Y_t)]$ or $E_{\lambda}[g(X_t, Y_{t-1}, Y_t)]$ is with respect to the joint stationary distribution of \mathbf{X} and \mathbf{Y} , provided that \mathbf{X} is stationary.

4 JOINT ERGODICITY

In this section, we establish conditions under which the joint sequence (\mathbf{X}, \mathbf{Y}) is ergodic. The key step is

to embed this joint sequence in a Markov chain in a random environment. Recall that \mathcal{X} is the space of sequences $\mathbf{x} = (x_t)_{t \in \mathbb{Z}}$ in \mathcal{X} , and \mathcal{A} is the corresponding product sigma-field. Let τ denote the shift operator on \mathcal{X} , that is, $\tau^k \mathbf{x} = \mathbf{x}'$ with $\mathbf{x}' = (x'_t)_{t \in \mathbb{Z}}$ given by $x'_t = x_{t+k}$. As we show below, a sufficient condition for joint ergodicity is that besides (A1) the following condition is satisfied:

- (A2) \mathbf{X} is ergodic, i.e., the probability measure $P_{\mathbf{X}}$ on $(\mathcal{X}, \mathcal{A})$ satisfies $P_{\mathbf{X}}(A) = P_{\mathbf{X}}(\tau^{-1}A)$ for every $A \in \mathcal{A}$, and $P_{\mathbf{X}}(A) = 0$ or $P_{\mathbf{X}}(A) = 1$ for every set $A \in \mathcal{A}$ such that $A = \tau^{-1}A$.

A particular consequence of condition (A2) is that \mathbf{X} is stationary. Furthermore, for every measurable function $g : \mathcal{X} \rightarrow \mathbb{R}$ such that $E[g(\mathbf{X})] < \infty$,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n g(\tau^t \mathbf{X}) = E[g(\mathbf{X})]$$

P -almost surely (Cornfeld et al., 1982), where we omit the subscript λ in P_{λ} and E_{λ} as the distribution of \mathbf{X} does not depend on λ . Note that condition (A2) is not too restrictive. For example, any stationary ARMA process is ergodic, and a sufficient condition for stationary Gaussian processes to be ergodic is that the autocovariances go to zero as the lag goes to infinity.

4.1 Markov Chains in Random Environments

Let us demonstrate how the joint sequence (\mathbf{X}, \mathbf{Y}) can be embedded in a Markov chain. By $\tilde{\mathcal{X}}$ we denote the space of sequences $\tilde{\mathbf{x}} = (\mathbf{x}_t)_{t \in \mathbb{Z}}$ with $\mathbf{x}_t \in \mathcal{X}$ for $t \in \mathbb{Z}$, that is, each \mathbf{x}_t is a sequence in \mathcal{X} . Write $\tilde{\mathcal{A}}$ for the corresponding product sigma-field. Let $\tilde{\tau}$ denote the shift operator on $\tilde{\mathcal{X}}$, that is, $\tilde{\tau}^k \tilde{\mathbf{x}} = \tilde{\mathbf{x}'}$ with $\tilde{\mathbf{x}'} = (\mathbf{x}'_t)_{t \in \mathbb{Z}}$ given by $\mathbf{x}'_t = \mathbf{x}_{t+k}$. For a given probability measure π on $(\mathcal{X}, \mathcal{A})$, let the probability measure $\tilde{\pi}$ on $(\tilde{\mathcal{X}}, \tilde{\mathcal{A}})$ be defined as follows: For $\tilde{A} \in \tilde{\mathcal{A}}$ write $\tilde{A} = \times_{t \in \mathbb{Z}} A_t$ with $A_t \in \mathcal{A}$. Then define

$$\tilde{\pi}(\tilde{A}) := \pi\left(\bigcap_{t \in \mathbb{Z}} \tau^{-t} A_t\right).$$

We say that $\tilde{\pi}$ is $\tilde{\tau}$ -ergodic if $\tilde{\pi}(\tilde{\tau}^{-1}\tilde{A}) = \tilde{\pi}(\tilde{A})$ for all $\tilde{A} \in \tilde{\mathcal{A}}$, and $\tilde{\pi}(\tilde{A}) = 0$ or $\tilde{\pi}(\tilde{A}) = 1$ for each $\tilde{A} \in \tilde{\mathcal{A}}$ satisfying $\tilde{\tau}^{-1}\tilde{A} = \tilde{A}$. The proof of the following technical lemma is included in the supplementary material.

Lemma 2. *If π is τ -ergodic, then $\tilde{\pi}$ is $\tilde{\tau}$ -ergodic.*

Note that according to Lemma 2, if the sequence \mathbf{X} is ergodic, then the probability measure $\tilde{P}_{\mathbf{X}}$ on $(\tilde{\mathcal{X}}, \tilde{\mathcal{A}})$ is $\tilde{\tau}$ -ergodic.

Now consider the measurable space $(\mathcal{Z}, \mathcal{C})$ with $\mathcal{Z} = \tilde{\mathcal{X}} \times \mathcal{Y} \times \mathcal{Y}$ and $\mathcal{C} = \tilde{\mathcal{A}} \times \mathcal{B} \times \mathcal{B}$. Let $\lambda \in \Theta$ be fixed. We are going to define a Markov sequence $\mathbf{Z} = (Z_t)_{t \in \mathbb{N}}$ with values in \mathcal{Z} such that \mathbf{Z} has the same distribution as the sequence $((\tau^{k+t} \mathbf{X})_{k \in \mathbb{Z}}, Y_{t-1}, Y_t)_{t \in \mathbb{N}}$ measured with respect to P_{λ} . Using results on invariant measures of Markov processes we are then going to show that the sequence \mathbf{Z} is ergodic and, consequently,

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n g((\tau^{t+k} \mathbf{X})_{k \in \mathbb{Z}}, Y_{t-1}, Y_t) \\ = E_{\lambda}[g((\tau^{t+k} \mathbf{X})_{k \in \mathbb{Z}}, Y_{t-1}, Y_t)] \end{aligned} \quad (4)$$

P_{λ} -almost surely for every measurable $g : \mathcal{Z} \rightarrow \mathbb{R}$ for which $E_{\lambda}|g((\tau^{t+k} \mathbf{X})_{k \in \mathbb{Z}}, Y_{t-1}, Y_t)| < \infty$. As a special case, we obtain the following result which will be of great importance for analyzing the asymptotical properties of Maximum Likelihood estimates.

Theorem 3. *If conditions (A1) and (A2) hold, and $g : \mathcal{X} \times \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ satisfies $E_{\lambda}|g(X_t, Y_{t-1}, Y_t)| < \infty$, then*

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

P_{λ} -almost surely.

In order to establish this theorem, let us first consider how to define \mathbf{Z} . Let μ_{λ} be the distribution on $(\mathcal{Z}, \mathcal{C})$ induced by $((\tau^{k+1} \mathbf{X})_{k \in \mathbb{Z}}, Y_0, Y_1)$ measured with respect to P_{λ} , that is,

$$\begin{aligned} \mu_{\lambda}(\tilde{A} \times \{y_0\} \times \{y_1\}) \\ := P_{\lambda}((\tau^{k+1} \mathbf{X})_{k \in \mathbb{Z}} \in \tilde{A}, Y_0 = y_0, Y_1 = y_1). \end{aligned}$$

This will serve us as the initial distribution of \mathbf{Z} . In order to specify a Markov kernel on $(\mathcal{Z}, \mathcal{C})$, let $Q(\lambda, \mathbf{x})$ with $\mathbf{x} \in \mathcal{X}$ denote the $\ell \times \ell$ -matrix with the (i, j) -th component

$$Q(\lambda, \mathbf{x}, i, j) = m_{\lambda}(x_0, i, j) \lim_{n \rightarrow \infty} \frac{\beta_1^n(\lambda, \mathbf{x}, j)}{\beta_0^n(\lambda, \mathbf{x}, i)}.$$

Note that, for this matrix to be well-defined, it is sufficient to assume that condition (A1) holds. Now we define the Markov kernel Q_{λ} from $(\mathcal{Z}, \mathcal{C})$ into itself,

$$\begin{aligned} Q_{\lambda}((\tilde{\mathbf{x}}, y'_0, y'_1), \tilde{A} \times \{y_0\} \times \{y_1\}) \\ := \begin{cases} Q(\lambda, \mathbf{x}_0, y'_1, y_1) & \text{if } y_0 = y'_1 \text{ and } \tilde{\tau}\tilde{\mathbf{x}} \in \tilde{A}, \\ 0 & \text{otherwise.} \end{cases} \end{aligned}$$

It is not difficult to see that \mathbf{Z} with the initial distribution μ_{λ} and the Markov kernel Q_{λ} has the same distribution as $((\tau^{k+t} \mathbf{X})_{k \in \mathbb{Z}}, Y_{t-1}, Y_t)_{t \in \mathbb{N}}$ measured with respect to P_{λ} . Note that \mathbf{Z} can be regarded as a Markov chain in a random environment (Cogburn, 1984; Orey, 1991). In particular, the pairs (Y_{t-1}, Y_t) with $t \in \mathbb{N}$ form a Markov chain where the transition probabilities from (Y_{t-1}, Y_t) to (Y_t, Y_{t+1}) are governed by the stationary environment $\tau^t \mathbf{X}$.

4.2 Invariant Measures

For the proof that \mathcal{Z} is ergodic and hence (4) applies, we use results on invariant measures of Markov chains. In fact, by the stationarity of \mathbf{X} and Corollary 1 (iii), it is easily verified that the measure μ_λ is *invariant* for Q_λ , that is,

$$\mu_\lambda(C) = \int_{\mathcal{Z}} Q_\lambda(z, C) \mu_\lambda(dz)$$

for every $C \in \mathcal{C}$. It only remains to show that the invariant measure μ_λ for Q_λ is unique (Corollary 2.5.2, Hernández-Lerma and Lasserre, 2003). This fact is stated in the next lemma, the proof of which is included in the supplementary material.

Lemma 3. *Suppose that conditions (A1) and (A2) hold. Then the invariant measure μ_λ for Q_λ is unique.*

5 CONSISTENCY OF MAXIMUM LIKELIHOOD ESTIMATES

In this section, we apply the previous results to study the following problem: Suppose that the distribution of the sequences \mathbf{X} and \mathbf{Y} is governed by P_{λ_0} with $\lambda_0 \in \Theta$ unknown, and we observe realizations of the subsequences $\mathbf{X}_n = (X_1, \dots, X_n)$ and $\mathbf{Y}_n = (Y_1, \dots, Y_n)$. Under which conditions can we identify λ_0 as the sample length n tends to infinity?

According to Corollary 1 (i), we may assume in our analysis that the weighted state feature functions are equal to zero. Therefore, to simplify notation, we write \mathbf{f} instead of $\mathbf{f}^{\text{trans}}$ in the following, so the weighted feature functions in (1) are given by

$$\lambda^T \mathbf{f}(x_n, \mathbf{y}_n) = \sum_{t=1}^n \lambda^T \mathbf{f}(x_t, y_{t-1}, y_t).$$

In order to estimate λ_0 , consider the objective function

$$\mathcal{L}_n(\lambda) = \frac{1}{n} \left(\lambda^T \mathbf{f}(\mathbf{X}_n, \mathbf{Y}_n) - \log Z_\lambda(\mathbf{X}_n) \right) \quad (5)$$

with the partition function $Z_\lambda(\mathbf{X}_n)$ as in (3). Note that $\mathcal{L}_n(\lambda)$ is the average conditional log-likelihood of \mathbf{Y}_n given \mathbf{X}_n based on the finite L-CRF model $P_\lambda^{(0:n)}$. Now consider the estimate $\hat{\lambda}_n$ of λ_0 obtained by maximizing the conditional log-likelihood,

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

If $\mathcal{L}_n(\lambda)$ is strictly concave, then the arg max is unique and can be found using gradient-based search (Sha and Pereira, 2003). Obviously, a necessary and sufficient condition for $\mathcal{L}_n(\lambda)$ to be strictly concave is that the number of labels ℓ is greater than or equal to 2, and

there exists a $\mathbf{y}_n \in \mathcal{Y}^{n+1}$ such that at least one component of $\mathbf{f}(\mathbf{X}_n, \mathbf{y}_n)$ is non-zero.

In the following we investigate conditions under which the estimates $\hat{\lambda}_n$ are strongly consistent, that is,

$$\lim_{n \rightarrow \infty} \hat{\lambda}_n = \lambda_0$$

P_{λ_0} -almost surely (Lehmann, 1999). Sufficient conditions will be given in Theorem 4 below. The key step is to establish conditions under which $\mathcal{L}_n(\lambda)$ converges uniformly to a function $\mathcal{L}(\lambda)$, and $\mathcal{L}(\lambda)$ has a unique maximum in λ_0 . To establish uniform convergence, we need to make the following assumption on the parameter space:

(A3) The parameter space Θ is compact.

In our case where the model parameters are \mathbb{R} -valued vectors, a sufficient condition for (A3) is that Θ is the Cartesian product of finite closed intervals.

5.1 Convergence of the Likelihood Function

First, we show that $\mathcal{L}_n(\lambda)$ converges for every $\lambda \in \Theta$.

Lemma 4. *Suppose that assumptions (A1) and (A2) hold. Then there exists a function $\mathcal{L}(\lambda)$ such that, for every $\lambda \in \Theta$,*

$$\lim_{n \rightarrow \infty} \mathcal{L}_n(\lambda) = \mathcal{L}(\lambda)$$

P_{λ_0} -almost surely.

Proof. Let $\lambda \in \Theta$. We show convergence separately for both terms on the right hand side of (5). For the first term we obtain, according to Theorem 3,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \lambda^T \mathbf{f}(\mathbf{X}_n, \mathbf{Y}_n) = \lambda^T E_{\lambda_0}[\mathbf{f}(X_t, Y_{t-1}, Y_t)]$$

P_{λ_0} -almost surely. For the second term, note that

$$\frac{1}{n} \log Z_\lambda(\mathbf{X}_n) \sim \frac{1}{n} \log \|\mathbf{M}_\lambda(X_1) \dots \mathbf{M}_\lambda(X_n)\|$$

as $n \rightarrow \infty$, where \sim denotes asymptotical equivalence and $\|(m_{ij})\| = \sum_{i,j} |m_{ij}|$. Now, let $L_1(x_1) := \log \|\mathbf{M}_\lambda(x_1)\|$ and, for $t > 1$,

$$L_t(x_t, \dots, x_1) := \log \frac{\|\mathbf{M}_\lambda(x_1) \dots \mathbf{M}_\lambda(x_t)\|}{\|\mathbf{M}_\lambda(x_1) \dots \mathbf{M}_\lambda(x_{t-1})\|}.$$

By the same arguments as in the proof of Theorem 2, the rows of the matrices $\mathbf{M}_\lambda(x_{t-k}) \dots \mathbf{M}_\lambda(x_t)$ and $\mathbf{M}_\lambda(x_{t-k}) \dots \mathbf{M}_\lambda(x_{t-1})$ tend to proportionality as k goes to ∞ . Therefore,

$$L(x_t, x_{t-1}, \dots) := \lim_{k \rightarrow \infty} L_k(x_t, \dots, x_{t-k+1})$$

is well-defined. Putting all together, we obtain

$$\frac{1}{n} \log Z_{\lambda}(\mathbf{X}_n) \sim \frac{1}{n} \sum_{t=1}^n L(X_t, X_{t-1}, \dots),$$

and by the ergodicity of \mathbf{X} the latter expression converges P_{λ_0} -almost surely to $E_{\lambda_0}[L(X_t, X_{t-1}, \dots)]$. Hence the proof is complete. \square

For the proof that the convergence of $\mathcal{L}_n(\lambda)$ to $\mathcal{L}(\lambda)$ is uniform on Θ , we need to consider the gradient of $\mathcal{L}_n(\lambda)$, which is given by

$$\begin{aligned} \nabla \mathcal{L}_n(\lambda) &= \frac{1}{n} \mathbf{f}(\mathbf{X}_n, \mathbf{Y}_n) \\ &\quad - \frac{1}{n} \sum_{\mathbf{y}_n \in \mathcal{Y}^n} \frac{\exp(\lambda^T \mathbf{f}(\mathbf{X}_n, \mathbf{y}_n))}{Z_{\lambda}(\mathbf{X}_n)} \mathbf{f}(\mathbf{X}_n, \mathbf{y}_n). \end{aligned}$$

Lemma 5. *Suppose that assumptions (A1) and (A2) hold. Then, for every $\lambda \in \Theta$,*

$$\begin{aligned} \lim_{n \rightarrow \infty} \nabla \mathcal{L}_n(\lambda) &= E_{\lambda_0}[\mathbf{f}(X_t, Y_{t-1}, Y_t)] \\ &\quad - E_{\lambda}[\mathbf{f}(X_t, Y_{t-1}, Y_t)] \end{aligned}$$

P_{λ_0} -almost surely.

Proof. The convergence of the first term again follows by Theorem 3. For the second term, note that

$$\begin{aligned} \frac{1}{n} \sum_{\mathbf{y}_n \in \mathcal{Y}^n} \frac{\exp(\lambda^T \mathbf{f}(\mathbf{X}_n, \mathbf{y}_n))}{Z_{\lambda}(\mathbf{X}_n)} \mathbf{f}(\mathbf{X}_n, \mathbf{y}_n) \\ = \frac{1}{n} \sum_{t=1}^n E_{\lambda}^{(0:n)}[\mathbf{f}(X_t, Y_{t-1}, Y_t) | \mathbf{X}], \end{aligned}$$

where $E_{\lambda}^{(0:n)}[\mathbf{f}(X_t, Y_{t-1}, Y_t) | \mathbf{X}]$ is the conditional expectation of $\mathbf{f}(X_t, Y_{t-1}, Y_t)$ given \mathbf{X} under the finite L-CRF model $P_{\lambda}^{(0:n)}$. Using arguments similar to the proof of the uniform bound in Theorem 2, one can show that $E_{\lambda}^{(0:n)}[\mathbf{f}(X_t, Y_{t-1}, Y_t) | \mathbf{X}]$ can be replaced by the conditional expectation of $\mathbf{f}(X_t, Y_{t-1}, Y_t)$ given \mathbf{X} under the infinite L-CRF model P_{λ} , that is,

$$\begin{aligned} \frac{1}{n} \sum_{t=1}^n E_{\lambda}^{(0:n)}[\mathbf{f}(X_t, Y_{t-1}, Y_t) | \mathbf{X}] \\ \sim \frac{1}{n} \sum_{t=1}^n E_{\lambda}[\mathbf{f}(X_t, Y_{t-1}, Y_t) | \mathbf{X}]. \end{aligned}$$

See the supplementary material for details. Now, as $E_{\lambda}[\mathbf{f}(X_t, Y_{t-1}, Y_t) | \mathbf{X}] = E_{\lambda}[\mathbf{f}(X_0, Y_{-1}, Y_0) | \tau^t \mathbf{X}]$ for every t , we obtain

$$\frac{1}{n} \sum_{t=1}^n E_{\lambda}[\mathbf{f}(X_t, Y_{t-1}, Y_t) | \mathbf{X}] = E_{\lambda}[\mathbf{f}(X_t, Y_{t-1}, Y_t)]$$

P_{λ} -almost surely by the ergodicity of \mathbf{X} . \square

Now we are ready to prove that $\mathcal{L}_n(\lambda)$ converges to $\mathcal{L}(\lambda)$ uniformly on Θ .

Lemma 6. *Suppose that conditions (A1)-(A3) hold. Then the convergence of $\mathcal{L}_n(\lambda)$ to $\mathcal{L}(\lambda)$ is uniform on Θ , that is,*

$$\lim_{n \rightarrow \infty} \sup_{\lambda \in \Theta} |\mathcal{L}_n(\lambda) - \mathcal{L}(\lambda)| = 0$$

P_{λ_0} -almost surely.

Proof. Since Θ is compact, it is sufficient to show that $\mathcal{L}_n(\lambda)$ is stochastically equicontinuous, i.e., for P_{λ} -almost every $\omega \in \Omega$ and every $\epsilon > 0$, there exists a $\delta > 0$ and an $n_0(\omega)$ such that

$$\sup_{\|\lambda_1 - \lambda_2\| \leq \delta} |\mathcal{L}_n(\lambda_1)(\omega) - \mathcal{L}_n(\lambda_2)(\omega)| \leq \epsilon$$

for all $n \geq n_0(\omega)$. By the Mean value theorem, there exists a (random) $h \in [0, 1]$ such that

$$\begin{aligned} |\mathcal{L}_n(\lambda_1) - \mathcal{L}_n(\lambda_2)| &\leq \|\lambda_1 - \lambda_2\| \\ &\quad \times \|\nabla \mathcal{L}_n((1-h)\lambda_1 + h\lambda_2)\|. \end{aligned}$$

To bound the second factor on the right hand side note that for any λ (not necessarily lying in Θ),

$$\begin{aligned} \|\nabla \mathcal{L}_n(\lambda)\| &\leq \frac{1}{n} \sum_{t=1}^n \|\mathbf{f}(X_t, Y_{t-1}, Y_t)\| \\ &\quad + \frac{1}{n} \sum_{t=1}^n \sum_{i,j \in \mathcal{Y}} \|\mathbf{f}(X_t, i, j)\|. \end{aligned}$$

Let U_n denote this upper bound. By the ergodicity of \mathbf{X} we obtain that U_n converges P_{λ_0} -almost surely to a finite limit, which we denote by U . Now, for P_{λ_0} -almost every $\omega \in \Omega$ there exists an $n_0(\omega)$ such that $\|U_n(\omega) - U\| \leq \epsilon$ for all $n \geq n_0(\omega)$. Substituting this into the above inequality, we obtain

$$|\mathcal{L}_n(\lambda_1)(\omega) - \mathcal{L}_n(\lambda_2)(\omega)| \leq (U + \epsilon) \|\lambda_1 - \lambda_2\|$$

for all $n \geq n_0(\omega)$. Setting $\delta := \epsilon/(U + \epsilon)$, we see that the sequence $\mathcal{L}_n(\lambda)$ is stochastically equicontinuous. \square

5.2 Convergence of the Hessian

Based on the previous results, we are now able to state the following sufficient conditions for strong consistency of $\hat{\lambda}_n$.

Theorem 4. *Suppose that conditions (A1)-(A3) hold, and the limit of $\nabla^2 \mathcal{L}_n(\lambda)$ is finite and strictly negative definite. Then $\mathcal{L}(\lambda)$ is strictly concave on Θ , and*

$$\lim_{n \rightarrow \infty} \hat{\lambda}_n = \lambda_0$$

P_{λ_0} -almost surely.

Proof. According to Lemma 6, we have uniform convergence of $\mathcal{L}_n(\boldsymbol{\lambda})$ to $\mathcal{L}(\boldsymbol{\lambda})$ on Θ . Thus, if the limit of $\nabla^2 \mathcal{L}_n(\boldsymbol{\lambda})$ is strictly negative definite, $\mathcal{L}(\boldsymbol{\lambda})$ is strictly concave and hence has a unique maximum. It only remains to show that this maximum is $\mathcal{L}(\boldsymbol{\lambda}_0)$. Under the assumption that the limit of $\nabla^2 \mathcal{L}_n(\boldsymbol{\lambda})$ is finite, the gradient of $\mathcal{L}(\boldsymbol{\lambda})$ is given by the limit of $\nabla \mathcal{L}_n(\boldsymbol{\lambda})$. According to Lemma 5, this limit is zero if $\boldsymbol{\lambda} = \boldsymbol{\lambda}_0$, hence $\mathcal{L}(\boldsymbol{\lambda}_0)$ is the unique maximum of $\mathcal{L}(\boldsymbol{\lambda})$ on Θ . \square

Let us analyze the asymptotic behavior of $\nabla^2 \mathcal{L}_n(\boldsymbol{\lambda})$. In order to compute the second partial derivatives, let $n \in \mathbb{N}$ and write $\boldsymbol{\lambda}$ as a stacked vector $\boldsymbol{\lambda} = (\boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_n)$ with $\boldsymbol{\lambda}_1 = \dots = \boldsymbol{\lambda}_n$. Correspondingly, consider the stacked feature functions $\mathbf{f}(\mathbf{X}_n, \mathbf{Y}_n) = (\mathbf{f}(X_1, Y_0, Y_1), \dots, \mathbf{f}(X_n, Y_{n-1}, Y_n))$, so that

$$\boldsymbol{\lambda}^T \mathbf{f}(\mathbf{X}_n, \mathbf{Y}_n) = \sum_{t=1}^n \boldsymbol{\lambda}_t^T \mathbf{f}(X_t, Y_{t-1}, Y_t).$$

Note that the first partial derivatives of $\mathcal{L}_n(\boldsymbol{\lambda})$ with respect to $\boldsymbol{\lambda}_t$ are given by

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\lambda}_t} \mathcal{L}_n(\boldsymbol{\lambda}) &= \frac{1}{n} \mathbf{f}(X_t, Y_{t-1}, Y_t) \\ &\quad - \frac{1}{n} E_{\boldsymbol{\lambda}}^{(0:n)} [\mathbf{f}(X_t, Y_{t-1}, Y_t) \mid \mathbf{X}]. \end{aligned}$$

By further differentiating this expression with respect to $\boldsymbol{\lambda}_{t+k}$ (for k such that $1 \leq t+k \leq n$), we obtain

$$\begin{aligned} \frac{\partial^2}{\partial \boldsymbol{\lambda}_t \partial \boldsymbol{\lambda}_{t+k}} \mathcal{L}_n(\boldsymbol{\lambda}) &= -\frac{1}{n} E_{\boldsymbol{\lambda}}^{(0:n)} [\mathbf{f}(X_t, Y_{t-1}, Y_t) \mid \mathbf{X}] \\ &\quad \times E_{\boldsymbol{\lambda}}^{(0:n)} [\mathbf{f}(X_{t+k}, Y_{t+k-1}, Y_{t+k}) \mid \mathbf{X}]^T \\ &\quad + \frac{1}{n} E_{\boldsymbol{\lambda}}^{(0:n)} [\mathbf{f}(X_t, Y_{t-1}, Y_t) \mathbf{f}(X_{t+k}, Y_{t+k-1}, Y_{t+k})^T \mid \mathbf{X}]. \end{aligned}$$

For $k = 0, 1, \dots, n-1$, consider the sum of all second partial derivatives with the lag k ,

$$\hat{\gamma}_{\boldsymbol{\lambda}}^{(n)}(k) = \sum_{t=1}^{n-k} \frac{\partial^2}{\partial \boldsymbol{\lambda}_t \partial \boldsymbol{\lambda}_{t+k}} \mathcal{L}_n(\boldsymbol{\lambda}),$$

and note that the Hessian of $\mathcal{L}_n(\boldsymbol{\lambda})$ can be written as

$$\nabla^2 \mathcal{L}_n(\boldsymbol{\lambda}) = - \left(\hat{\gamma}_{\boldsymbol{\lambda}}^{(n)}(0) + 2 \sum_{k=1}^{n-1} \hat{\gamma}_{\boldsymbol{\lambda}}^{(n)}(k) \right).$$

The following lemma shows that, if conditions (A1) and (A2) are satisfied, the limit of $\nabla^2 \mathcal{L}_n(\boldsymbol{\lambda})$ is finite. The proof is included in the supplementary material. According to the proof of Theorem 4, we obtain that $\boldsymbol{\lambda}_0$ is one solution of $\arg \max_{\boldsymbol{\lambda} \in \Theta} \mathcal{L}(\boldsymbol{\lambda})$, however, this solution is not unique unless the limit of $\nabla^2 \mathcal{L}_n(\boldsymbol{\lambda})$ is non-singular.

Lemma 7. *Suppose that conditions (A1) and (A2) hold. Then*

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

$P_{\boldsymbol{\lambda}_0}$ -almost surely, where

$$\begin{aligned} \gamma_{\boldsymbol{\lambda}}(k) &= \\ &\quad \text{Cov}_{\boldsymbol{\lambda}} [\mathbf{f}(X_t, Y_{t-1}, Y_t), \mathbf{f}(X_{t+k}, Y_{t+k-1}, Y_{t+k})] \end{aligned}$$

is the matrix of covariances between $\mathbf{f}(X_t, Y_{t-1}, Y_t)$ and $\mathbf{f}(X_{t+k}, Y_{t+k-1}, Y_{t+k})$ measured with respect to $P_{\boldsymbol{\lambda}}$. In particular, the limit of $\nabla^2 \mathcal{L}_n(\boldsymbol{\lambda})$ is finite.

The following corollary states a simple necessary condition for the limit of the Hessian $\nabla^2 \mathcal{L}_n(\boldsymbol{\lambda})$ to be non-singular and hence for the solution of $\arg \max_{\boldsymbol{\lambda} \in \Theta} \mathcal{L}(\boldsymbol{\lambda})$ to be unique.

Corollary 2. *Suppose that conditions (A1) and (A2) hold and the vector of feature functions \mathbf{f} has dimensionality d . Then a necessary condition for the limit of $\nabla^2 \mathcal{L}_n(\boldsymbol{\lambda})$ to be non-singular is that for every pair $a \in \mathbb{R}$, $\mathbf{b} \in \mathbb{R}^d$ with $\mathbf{b} \neq \mathbf{0}$ there exist $i, j \in \mathcal{Y}$ and a subset $A \subset \mathcal{X}$ with $P_{\mathbf{X}}(A) > 0$ such that $\mathbf{b}^T \mathbf{f}(x, i, j) \neq a$ for all $x \in A$.*

In particular, the solution of $\arg \max_{\boldsymbol{\lambda} \in \Theta} \mathcal{L}(\boldsymbol{\lambda})$ fails to be unique if any of the components of \mathbf{f} can be expressed as linear combinations of each other. We leave it as an open problem whether the conditions in Corollary 2 are also sufficient for non-singularity. Note that the answer is affirmative when the feature functions $\mathbf{f}(x_t, y_{t-1}, y_t)$ do not depend on y_{t-1} . In this case the matrices $\gamma_{\boldsymbol{\lambda}}(k)$ are zero for all $k > 0$, and hence the limit of $\nabla^2 \mathcal{L}_n(\boldsymbol{\lambda})$ is equal to $\gamma_{\boldsymbol{\lambda}}(0)$.

6 CONCLUSIONS

We have taken a first step to a rigorous study of asymptotic properties of Maximum Likelihood Estimates (MLEs) in Linear-Chain Conditional Random Fields (L-CRFs). Our analysis is based on L-CRFs for infinite sequences, which are defined by the limit distributions of conventional L-CRFs as the length of the observational context tends to infinity. We have derived basic properties of these L-CRFs and shown that ergodicity of the observation sequence implies ergodicity of the joint sequence of observations and labels. Based on these results, we have established uniform convergence of the average conditional log-likelihood and of the gradient, and pointwise convergence of the Hessian. Under the assumption that the limit of the Hessian is non-singular, our results show that the MLEs are strongly consistent. The question under which conditions non-singularity holds is an open problem for future research.

References

- R. Cogburn (1984). The ergodic theory of Markov chains in random environments. *Z. Wahrscheinlichkeitstheorie verw. Gebiete* **66**:109-128.
- I.P. Cornfeld, S.V. Fomin and Y.G. Sinai (1982). *Ergodic theory*. Berlin, Germany: Springer.
- A. Culotta, D. Kulp and A. McCallum (2005). Gene prediction with conditional random fields. *Technical Report UM-CS-2005-028*. University of Massachusetts, Amherst.
- S.R. Foguel (1969). *The ergodic theory of Markov processes*. Princeton, NJ: Van Nostrand.
- O. Hernández-Lerma and J.B. Lasserre (2003). *Markov chains and invariant probabilities*. Basel, Switzerland: Birkhäuser.
- 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)*.
- E.L. Lehmann (1999). *Elements of large-sample theory*. New York, NY: Springer.
- L. Liao, D. Choudhury, D. Fox, H. Kautz and B. Limketkai (2007a). Training conditional random fields using virtual evidence boosting. *Proceedings of the International Joint Conferences on Artificial Intelligence (IJCAI)*.
- L. Liao, D. Fox and H. Kautz (2007b). Extracting places and activities from GPS traces using hierarchical conditional random fields. *International Journal of Robotics Research* **26**(1):119-138
- A. McCallum (2003). Efficiently inducing features of conditional random fields. *Proceedings of the Conference on Uncertainty in Artificial Intelligence (UAI)*.
- F. Omar, M. Sinn, J. Truszkowski, P. Poupart, J. Tung and A. Caine (2010). Comparative analysis of probabilistic models for activity recognition with an instrumented walker. *Proceedings of the Conference on Uncertainty in Artificial Intelligence (UAI)*.
- S. Orey (1991). Markov chains with stochastically stationary transition probabilities. *The Annals of Probability* **19**(3):907-928
- S. Sarawagi and W.W. Cohen (2004). Semi-markov conditional random fields for information extraction. *Proceedings of the Conference on Advances in Neural Information Processing Systems (NIPS)*.
- E. Seneta (2006). *Non-negative matrices and Markov chains. Revised Edition*. New York, NY: Springer.
- F. Sha and F. Pereira (2003). Shallow parsing with conditional random fields. *Proceedings of HLT-NAACL*.
- C. Sutton and A. McCallum (2006). An introduction to Conditional random fields for relational learning. In L. Getoor and B. Taskar (eds.): *Introduction to statistical relational learning*. Cambridge, MA: MIT Press.
- B. Taskar, P. Abbeel and D. Koller (2002). Discriminative probabilistic models for relational data. *Proceedings of the Conference on Uncertainty in Artificial Intelligence (UAI)*.
- H. Wallach (2004). Conditional random fields: an introduction. *Technical Report MS-CIS-04-21*. University of Pennsylvania, Philadelphia.