
Mixing Properties of Conditional Markov Chains with Unbounded Feature Functions

Mathieu Sinn

IBM Research - Ireland
Mulhuddart, Dublin 15
mathsinn@ie.ibm.com

Bei Chen

McMaster University
Hamilton, Ontario, Canada
bei.chen@math.mcmaster.ca

Abstract

Conditional Markov Chains (also known as Linear-Chain Conditional Random Fields in the literature) are a versatile class of discriminative models for the distribution of a sequence of hidden states conditional on a sequence of observable variables. Large-sample properties of Conditional Markov Chains have been first studied in [1]. The paper extends this work in two directions: first, mixing properties of models with unbounded feature functions are being established; second, necessary conditions for model identifiability and the uniqueness of maximum likelihood estimates are being given.

1 Introduction

Conditional Random Fields (CRF) are a widely popular class of discriminative models for the distribution of a set of hidden states conditional on a set of observable variables. The fundamental assumption is that the hidden states, conditional on the observations, form a Markov random field [2,3]. Of special importance, particularly for the modeling of sequential data, is the case where the underlying undirected graphical model forms a simple linear chain. In the literature, this subclass of models is often referred to as Linear-Chain Conditional Random Fields. This paper adopts the terminology of [4] and refers to such models as Conditional Markov Chains (CMC).

Large-sample properties of CRFs and CMCs have been first studied in [1] and [5]. [1] defines CMCs of infinite length and studies ergodic properties of the joint sequences of observations and hidden states. The analysis relies on fundamental results from the theory of weak ergodicity [6]. The exposition is restricted to CMCs with bounded feature functions which precludes the application, e.g., to models with linear features and Gaussian observations. [5] considers weak consistency and central limit theorems for models with a more general structure. Ergodicity and mixing of the models is assumed, but no explicit conditions on the feature functions or on the distribution of the observations are given. An analysis of model identifiability in the case of finite sequences can be found in [7].

The present paper studies mixing properties of Conditional Markov Chains with unbounded feature functions. The results are fundamental for analyzing the consistency of Maximum Likelihood estimates and establishing Central Limit Theorems (which are very useful for constructing statistical hypothesis tests, e.g., for model misspecifications and the significance of features). The paper is organized as follows: Sec. 2 reviews the definition of infinite CMCs and some of their basic properties. In Sec. 3 the ergodicity results from [1] are extended to models with unbounded feature functions. Sec. 4 establishes various mixing properties. A key result is that, in order to allow for unbounded feature functions, the observations need to follow a distribution such that Hoeffding-type concentration inequalities can be established. Furthermore, the mixing rates depend on the tail behaviour of the distribution. In Sec. 5 the mixture properties are used to analyze model identifiability and consistency of the Maximum Likelihood estimates. Sec. 6 concludes with an outlook on open problems for future research.

2 Conditional Markov Chains

Preliminaries. We use \mathbb{N} , \mathbb{Z} and \mathbb{R} to 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} , and \mathcal{Y} be a finite set. Furthermore, consider a probability space $(\Omega, \mathcal{F}, \mathbb{P})$ and let $\mathbf{X} = (X_t)_{t \in \mathbb{Z}}$, $\mathbf{Y} = (Y_t)_{t \in \mathbb{Z}}$ be sequences of measurable mappings from Ω into \mathcal{X} and \mathcal{Y} , respectively. Here,

- \mathbf{X} is an infinite sequence of *observations* ranging in the domain \mathcal{X} ,
- \mathbf{Y} is an aligned sequence of *hidden states* taking values in the finite set \mathcal{Y} .

For now, the distribution of \mathbf{X} is arbitrary. Next we define Conditional Markov Chains, which parameterize the conditional distribution of \mathbf{Y} given \mathbf{X} .

Definition. Consider a vector \mathbf{f} of real-valued functions $f : \mathcal{X} \times \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$, called the *feature functions*. Throughout this paper, we assume that the following condition is satisfied:

(A1) All feature functions are finite: $|f(x, i, j)| < \infty$ for all $x \in \mathcal{X}$ and $i, j \in \mathcal{Y}$.

Associated with the feature functions is a vector $\boldsymbol{\lambda}$ of real-valued *model-weights*. The key in the definition of Conditional Markov Chains is the matrix $\mathbf{M}(x)$ with the (i, j) -th component

$$m(x, i, j) = \exp(\boldsymbol{\lambda}^T \mathbf{f}(x, i, j)).$$

In terms of statistical physics, $m(x, i, j)$ measures the *potential* of the transition between the hidden states i and j from time $t-1$ to t , given the observation x at time t . Next, for a sequence $\mathbf{x} = (x_t)_{t \in \mathbb{Z}}$ in \mathcal{X} and time points $s, t \in \mathbb{Z}$ with $s \leq t$, introduce the vectors

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

and write $\alpha_s^t(\mathbf{x}, i)$ and $\beta_s^t(\mathbf{x}, j)$ to denote the i th respectively j th components. Intuitively, $\alpha_s^t(\mathbf{x}, i)$ measures the potential of the hidden state i at time t given the observations x_s, \dots, x_t and assuming that at time $s-1$ all hidden states have potential equal to 1. Similarly, $\beta_s^t(\mathbf{x}, j)$ is the potential of j at time s assuming equal potential of all hidden states at time t . Now let $t \in \mathbb{Z}$ and $k \in \mathbb{N}$, and define the distribution of the labels Y_t, \dots, Y_{t+k} conditional on \mathbf{X} ,

$$\begin{aligned} \mathbb{P}(Y_t = y_t, \dots, Y_{t+k} = y_{t+k} \mid \mathbf{X}) &:= \prod_{i=1}^k m(X_{t+i}, y_{t+i-1}, y_{t+i}) \\ &\times \lim_{n \rightarrow \infty} \frac{\alpha_{t-n}^t(\mathbf{X}, y_t) \beta_{t+k}^{t+k+n}(\mathbf{X}, y_{t+k})}{\boldsymbol{\alpha}_{t-n}^t(\mathbf{X})^T \boldsymbol{\beta}_t^{t+k+n}(\mathbf{X})}. \end{aligned}$$

Note that, under assumption (A1), the limit on the right hand side is well-defined (see Theorem 2 in [1]). Furthermore, the family of all marginal distributions obtained this way satisfies the consistency conditions of Kolmogorov's Extension Theorem. Hence we obtain a unique distribution for \mathbf{Y} conditional on \mathbf{X} parameterized by the feature functions \mathbf{f} and the model weights $\boldsymbol{\lambda}$. Intuitively, the distribution is obtained by conditioning the marginal distributions of \mathbf{Y} on the finite observational context $(X_{t-n}, \dots, X_{t+k+n})$, and then letting the size of the context go to infinity.

Basic properties. We introduce the following notation: For any matrix $\mathbf{P} = (p_{ij})$ with strictly positive entries let $\phi(\mathbf{P})$ denote the mixing coefficient

$$\phi(\mathbf{P}) = \min_{i,j,k,l} \frac{p_{ik} p_{jl}}{p_{jk} p_{il}}.$$

Note that $0 \leq \phi(\mathbf{P}) \leq 1$. This coefficient will play a key role in the analysis of mixing properties. The following proposition summarizes fundamental properties of the distribution of \mathbf{Y} conditional on \mathbf{X} , which directly follow from the above definition (also see Corollary 1 in [1]).

Proposition 1. *Suppose that condition (A1) holds true. Then \mathbf{Y} conditional on \mathbf{X} forms a time-inhomogeneous Markov chain. Moreover, if \mathbf{X} is strictly stationary, then the joint distribution of the aligned sequences (\mathbf{X}, \mathbf{Y}) is strictly stationary. The conditional transition probabilities $P_t(\mathbf{x}, i, j) := \mathbb{P}(Y_t = j | Y_{t-1} = i, \mathbf{X} = \mathbf{x})$ of \mathbf{Y} given $\mathbf{X} = \mathbf{x}$ have the following form:*

$$P_t(\mathbf{x}, i, j) = m(x_t, i, j) \lim_{n \rightarrow \infty} \frac{\beta_t^n(\mathbf{x}, j)}{\beta_{t-1}^n(\mathbf{x}, i)}.$$

In particular, a lower bound for $P_t(\mathbf{x}, i, j)$ is given by

$$P_t(\mathbf{x}, i, j) \geq \frac{m(x_t, i, j) (\min_{k \in \mathcal{Y}} m(x_{t+1}, i, k))}{|\mathcal{Y}| (\max_{k \in \mathcal{Y}} m(x_t, j, k)) (\max_{k, l \in \mathcal{Y}} m(x_{t+1}, k, l))},$$

and the matrix of transition probabilities $\mathbf{P}_t(\mathbf{x})$, with the (i, j) -th component given by $P_t(\mathbf{x}, i, j)$, satisfies $\phi(\mathbf{P}_t(\mathbf{x})) = \phi(\mathbf{M}(x_t))$.

3 Ergodicity

In this section we establish conditions under which the aligned sequences (\mathbf{X}, \mathbf{Y}) are jointly ergodic. Let us first recall the definition of ergodicity of \mathbf{X} (see [8]): By \mathcal{X} we denote the space of sequences $\mathbf{x} = (x_t)_{t \in \mathbb{Z}}$ in \mathcal{X} , and by \mathcal{A} the corresponding product σ -field. Consider the probability measure $P_{\mathbf{X}}$ on $(\mathcal{X}, \mathcal{A})$ defined by $P_{\mathbf{X}}(\mathbf{A}) := \mathbb{P}(\mathbf{X} \in \mathbf{A})$ for $\mathbf{A} \in \mathcal{A}$. Finally, let τ denote the operator on \mathcal{X} which shifts sequences one position to the left: $\tau \mathbf{x} = (x_{t+1})_{t \in \mathbb{Z}}$. Then ergodicity of \mathbf{X} is formally defined as follows:

(A2) \mathbf{X} is ergodic, that is, $P_{\mathbf{X}}(\mathbf{A}) = P_{\mathbf{X}}(\tau^{-1} \mathbf{A})$ for every $\mathbf{A} \in \mathcal{A}$, and $P_{\mathbf{X}}(\mathbf{A}) \in \{0, 1\}$ for every set $\mathbf{A} \in \mathcal{A}$ satisfying $\mathbf{A} = \tau^{-1} \mathbf{A}$.

As a particular consequence of the invariance $P_{\mathbf{X}}(\mathbf{A}) = P_{\mathbf{X}}(\tau^{-1} \mathbf{A})$, we obtain that \mathbf{X} is strictly stationary. Now we are able to formulate the key result of this section, which will be of central importance in the later analysis. For simplicity, we state it for functions depending on the values of \mathbf{X} and \mathbf{Y} only at time t . The generalization of the statement is straight-forward. In our later analysis, we will use the theorem to show that the time average of feature functions $f(X_t, Y_{t-1}, Y_t)$ converges to the expected value $\mathbb{E}[f(X_t, Y_{t-1}, Y_t)]$.

Theorem 1. *Suppose that conditions (A1) and (A2) hold, and $g : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ is a function which satisfies $\mathbb{E}[|g(X_t, Y_t)|] < \infty$. Then*

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n g(X_t, Y_t) = \mathbb{E}[g(X_t, Y_t)] \quad \mathbb{P}\text{-almost surely.}$$

Proof. Consider the sequence $\mathbf{Z} = (Z_t)_{t \in \mathbb{N}}$ given by $Z_t := (\tau^{t-1} \mathbf{X}, Y_t)$, where we write τ^{t-1} to denote the $(t-1)$ th iterate of τ . Note that Z_t represents the hidden state at time t together with the entire aligned sequence of observations $\tau^{t-1} \mathbf{X}$. In the literature, such models are known as Markov sequences in random environments (see [9]). The key step in the proof is to show that \mathbf{Z} is ergodic. Then, for any function $h : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ with $\mathbb{E}[|h(Z_t)|] < \infty$, the time average $\frac{1}{n} \sum_{t=1}^n h(Z_t)$ converges to the expected value $\mathbb{E}[h(Z_t)]$ \mathbb{P} -almost surely. Applying this result to the composition of the function g and the projection of $(\tau^{t-1} \mathbf{X}, Y_t)$ onto (X_t, Y_t) completes the proof. The details of the proof that \mathbf{Z} is ergodic can be found in an extended version of this paper, which is included in the supplementary material. \square

4 Mixing properties

In this section we are going to study mixing properties of the aligned sequences (\mathbf{X}, \mathbf{Y}) . To establish the results, we will assume that the distribution of the observations \mathbf{X} satisfies conditions under which certain concentration inequalities hold true:

(A3) Let $A \subset \mathcal{A}$ be a measurable set, with $p := \mathbb{P}(X_t \in A)$ and $S_n(\mathbf{x}) := \frac{1}{n} \sum_{t=1}^n \mathbf{1}(x_t \in A)$ for $\mathbf{x} \in \mathcal{X}$. Then there exists a constant γ such that, for all $n \in \mathbb{N}$ and $\epsilon > 0$,

$$\mathbb{P}(|S_n(\mathbf{X}) - p| \geq \epsilon) \leq \exp(-\gamma \epsilon^2 n).$$

If \mathbf{X} is a sequence of independent random variables, then (A3) follows by Hoeffding's inequality. In the dependent case, concentration inequalities of this type can be obtained by imposing Martingale or mixing conditions on \mathbf{X} (see [12,13]). Furthermore, we will make the following assumption, which relates the feature functions to the tail behaviour of the distribution of \mathbf{X} :

(A4) Let $h : [0, \infty) \rightarrow [0, \infty)$ be a differentiable decreasing function with $h(z) = O(z^{-(1+\kappa)})$ for some $\kappa > 0$. Furthermore, let

$$F(x) := \sum_{j,k \in \mathcal{Y}} |\boldsymbol{\lambda}^T \mathbf{f}(x, j, k)|$$

for $x \in \mathcal{X}$. Then $\mathbb{E}[h(F(X_t))^{-1}]$ and $\mathbb{E}[h'(F(X_t))^{-1}]$ both exist and are finite.

The following theorem establishes conditions under which the expected conditional covariances of square-integrable functions are summable. The result is obtained by studying ergodic properties of the transition probability matrices.

Theorem 2. *Suppose that conditions (A1) - (A3) hold true, and $g : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ is a function with finite second moment, $\mathbb{E}[|g(X_t, Y_t)|^2] < \infty$. Let $\gamma_{t,k}(\mathbf{X}) = \text{Cov}(g(X_t, Y_t), g(X_{t+k}, Y_{t+k}) | \mathbf{X})$ denote the covariance of $g(X_t, Y_t)$ and $g(X_{t+k}, Y_{t+k})$ conditional on \mathbf{X} . Then, for every $t \in \mathbb{Z}$:*

$$\lim_{n \rightarrow \infty} \sum_{k=1}^n \mathbb{E}[|\gamma_{t,k}(\mathbf{X})|] < \infty.$$

Proof. Without loss of generality we may assume that g can be written as $g(x, y) = g(x)\mathbf{1}(y = i)$. Hence, using Hölder's inequality, we obtain

$$\mathbb{E}[|\gamma_{t,k}(\mathbf{X})|] \leq \mathbb{E}[|g(X_t)|] \mathbb{E}[|g(X_{t+k})|] \mathbb{E}[|\text{Cov}(\mathbf{1}(Y_t = i), \mathbf{1}(Y_{t+k} = i) | \mathbf{X})|].$$

According to the assumptions, we have $\mathbb{E}[|g(X_t)|] = \mathbb{E}[|g(X_{t+k})|] < \infty$, so we only need to bound the expectation of the conditional covariance. Note that

$$\text{Cov}(\mathbf{1}(Y_t = i), \mathbf{1}(Y_{t+k} = i) | \mathbf{X}) = \mathbb{P}(Y_t = i, Y_{t+k} = i | \mathbf{X}) - \mathbb{P}(Y_t = i | \mathbf{X}) \mathbb{P}(Y_{t+k} = i | \mathbf{X}).$$

Recall the definition of $\phi(\mathbf{P})$ before Corollary 1. Using probabilistic arguments, it is not difficult to show that the ratio of $\mathbb{P}(Y_t = i, Y_{t+k} = i | \mathbf{X})$ to $\mathbb{P}(Y_t = i | \mathbf{X}) \mathbb{P}(Y_{t+k} = i | \mathbf{X})$ is greater than or equal to $\phi(\mathbf{P}_{t+1}(\mathbf{X}) \dots \mathbf{P}_{t+k}(\mathbf{X}))$, where $\mathbf{P}_{t+1}(\mathbf{X}), \dots, \mathbf{P}_{t+k}(\mathbf{X})$ denote the transition matrices introduced in Proposition 1. Hence,

$$|\text{Cov}(\mathbf{1}(Y_t = i), \mathbf{1}(Y_{t+k} = i) | \mathbf{X})| \leq \mathbb{P}(Y_t = i, Y_{t+k} = i | \mathbf{X}) [1 - \phi(\mathbf{P}_{t+1}(\mathbf{X}) \dots \mathbf{P}_{t+k}(\mathbf{X}))].$$

Now, using results from the theory of weak ergodicity (see Chapter 3 in [6]), we can establish

$$\frac{1 - \sqrt{\phi(\mathbf{P}_{t+1}(\mathbf{x}) \dots \mathbf{P}_{t+k}(\mathbf{x}))}}{1 + \sqrt{\phi(\mathbf{P}_{t+1}(\mathbf{x}) \dots \mathbf{P}_{t+k}(\mathbf{x}))}} \leq \prod_{j=1}^k \frac{1 - \sqrt{\phi(\mathbf{P}_{t+j}(\mathbf{x}))}}{1 + \sqrt{\phi(\mathbf{P}_{t+j}(\mathbf{x}))}}$$

for all $\mathbf{x} \in \mathcal{X}$. Using Bernoulli's inequality and the fact $\phi(\mathbf{P}_{t+j}(\mathbf{x})) = \mathbf{M}(x_{t+j})$ established in Proposition 1, we obtain $\phi(\mathbf{P}_{t+1}(\mathbf{x}) \dots \mathbf{P}_{t+k}(\mathbf{x})) \geq 1 - 4 \prod_{j=1}^k [1 - \phi(\mathbf{M}(x_{t+j}))]$. Consequently,

$$|\text{Cov}(\mathbf{1}(Y_t = i), \mathbf{1}(Y_{t+k} = i) | \mathbf{X})| \leq 4 \prod_{j=1}^k [1 - \phi(\mathbf{M}(X_{t+j}))].$$

With the notation introduced in assumption (A3), let $\delta > 0$ and $A \subset \mathcal{X}$ with $p > 0$ be such that $x \in A$ implies $\phi(\mathbf{M}(x)) \geq \delta$. Furthermore, let ϵ be a constant with $0 < \epsilon < p$. In order to bound $|\text{Cov}(\mathbf{1}(Y_t = i), \mathbf{1}(Y_{t+k} = i) | \mathbf{X})|$ for a given $k \in \mathbb{N}$, we distinguish two different cases: If $|S_k(\mathbf{X}) - p| < \epsilon$, then we obtain

$$4 \prod_{j=1}^k (1 - \phi(\mathbf{M}(X_{t+j}))) \leq 4(1 - \delta)^{k(p-\epsilon)}.$$

If $|S_k(\mathbf{X}) - p| \geq \epsilon$, then we use the trivial upper bound 1. According to assumption (A3), the probability of the latter event is bounded by an exponential, and hence

$$\mathbb{E}[|\text{Cov}(\mathbf{1}(Y_t = i), \mathbf{1}(Y_{t+k} = i) | \mathbf{X})|] \leq 4(1 - \delta)^{k(p-\epsilon)} + \exp(-\gamma \epsilon^2 k).$$

Obviously, the sum of all these expectations is finite, which completes the proof. \square

The next theorem bounds the difference between the distribution of \mathbf{Y} conditional on \mathbf{X} and finite approximations of it. Introduce the following notation: For $t, k \geq 0$ with $t + k \leq n$ let

$$\begin{aligned} \mathbb{P}^{(0:n)}(Y_t = y_t, \dots, Y_{t+k} = y_{t+k} \mid \mathbf{X} = \mathbf{x}) \\ := \prod_{i=1}^k m(x_{t+i}, y_{t+i-1}, y_{t+i}) \lim_{n \rightarrow \infty} \frac{\alpha_0^t(\mathbf{x}, y_t) \beta_{t+k}^n(\mathbf{x}, y_{t+k})}{\alpha_0^t(\mathbf{x})^T \beta_t^n(\mathbf{x})}. \end{aligned}$$

Accordingly, write $\mathbb{E}^{(0:n)}$ for expectations taken with respect to $\mathbb{P}^{(0:n)}$. As emphasized by the superscripts, these quantities can be regarded as marginal distributions of \mathbf{Y} conditional on the observations at times $t = 0, 1, \dots, n$. To simplify notation, the following theorem is stated for 1-dimensional conditional marginal distributions, however, the extension to the general case is straight-forward.

Theorem 3. *Suppose that conditions (A1) - (A4) hold true. Then the limit*

$$\mathbb{P}^{(0:\infty)}(Y_t = i \mid \mathbf{X}) := \lim_{n \rightarrow \infty} \mathbb{P}^{(0:n)}(Y_t = i \mid \mathbf{X})$$

is well-defined \mathbb{P} -almost surely. Moreover, there exists a measurable function $C(\mathbf{x})$ of $\mathbf{x} \in \mathcal{X}$ with finite expectation, $E[|C(\mathbf{X})|] < \infty$, and a function $h(z)$ satisfying the conditions in (A4), such that

$$|\mathbb{P}^{(0:\infty)}(Y_t = i \mid \mathbf{X}) - \mathbb{P}^{(0:n)}(Y_t = i \mid \mathbf{X})| \leq C(\tau^t \mathbf{X}) h(n - t).$$

Proof. Define $\mathbf{G}_n(\mathbf{x}) := \mathbf{M}(x_{t+1}) \dots \mathbf{M}(x_n)$ and write $g_n(\mathbf{x}, i, j)$ for the (i, j) -th component of $\mathbf{G}_n(\mathbf{x})$. Note that $\beta_t^n(\mathbf{x}) = \mathbf{G}_n(\mathbf{x})(1, 1, \dots, 1)^T$. According to Lemma 3.4 in [6], there exist numbers $r_{ij}(\mathbf{x})$ such that

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

for all $k \in \mathcal{Y}$. Consequently, the ratio of $\beta_t^n(\mathbf{x}, i)$ to $\beta_t^n(\mathbf{x}, j)$ converges to $r_{ij}(\mathbf{x})$, and hence

$$\lim_{n \rightarrow \infty} \frac{\alpha_0^t(\mathbf{x}, i) \beta_t^n(\mathbf{x}, i)}{\alpha_0^t(\mathbf{x})^T \beta_t^n(\mathbf{x})} = \frac{1}{\mathbf{q}_i(\mathbf{x})^T \mathbf{r}_i(\mathbf{x})}$$

where we use the notation $\mathbf{q}_i(\mathbf{x}) = \alpha_0^t(\mathbf{x}) / \alpha_0^t(\mathbf{x}, i)$ and $\mathbf{r}_i(\mathbf{x})$ denotes the vector with the k th component $r_{ki}(\mathbf{x})$. This proves the first part of the theorem. In order to prove the second part, note that $|x - y| \leq |x^{-1} - y^{-1}|$ for any $x, y \in (0, 1]$, and hence

$$|\mathbb{P}^{(0:\infty)}(Y_t = i \mid \mathbf{X}) - \mathbb{P}^{(0:n)}(Y_t = i \mid \mathbf{X})| \leq \left| \mathbf{q}_i(\mathbf{X})^T \mathbf{r}_i(\mathbf{X}) - \frac{\alpha_0^t(\mathbf{X})^T \beta_t^n(\mathbf{X})}{\alpha_0^t(\mathbf{X}, i) \beta_t^n(\mathbf{X}, i)} \right|.$$

To bound the latter expression, introduce the vectors $\underline{\mathbf{r}}_i^n(\mathbf{x})$ and $\bar{\mathbf{r}}_i^n(\mathbf{x})$ with the k th components

$$\underline{r}_{ki}^n(\mathbf{x}) = \min_{l \in \mathcal{Y}} \left(\frac{g_n(\mathbf{x}, k, l)}{g_n(\mathbf{x}, i, l)} \right) \quad \text{and} \quad \bar{r}_{ki}^n(\mathbf{x}) = \max_{l \in \mathcal{Y}} \left(\frac{g_n(\mathbf{x}, k, l)}{g_n(\mathbf{x}, i, l)} \right).$$

It is easy to see that $\mathbf{q}_i(\mathbf{x})^T \underline{\mathbf{r}}_i^n(\mathbf{x}) \leq \mathbf{q}_i(\mathbf{x})^T \mathbf{r}_i(\mathbf{x}) \leq \mathbf{q}_i(\mathbf{x})^T \bar{\mathbf{r}}_i^n(\mathbf{x})$, and

$$\mathbf{q}_i(\mathbf{x})^T \underline{\mathbf{r}}_i^n(\mathbf{x}) \leq \frac{\alpha_0^t(\mathbf{x})^T \beta_t^n(\mathbf{x})}{\alpha_0^t(\mathbf{x}, i) \beta_t^n(\mathbf{x}, i)} \leq \mathbf{q}_i(\mathbf{x})^T \bar{\mathbf{r}}_i^n(\mathbf{x}).$$

Hence,

$$\left| \mathbf{q}_i(\mathbf{X})^T \mathbf{r}_i(\mathbf{X}) - \frac{\alpha_0^t(\mathbf{X})^T \beta_t^n(\mathbf{X})}{\alpha_0^t(\mathbf{X}, i) \beta_t^n(\mathbf{X}, i)} \right| \leq \left| \mathbf{q}_i(\mathbf{X})^T (\bar{\mathbf{r}}_i^n(\mathbf{X}) - \underline{\mathbf{r}}_i^n(\mathbf{X})) \right|.$$

Due to space limitations, we only give a sketch of the proof how the latter quantity can be bounded. For details, see the extended version of this paper in the supplementary material. The first step is to show the existence of a function $C_1(\mathbf{x})$ with $\mathbb{E}[|C_1(\mathbf{X})|] < \infty$ such that $|\underline{r}_{ki}^n(\mathbf{X}) - \bar{r}_{ki}^n(\mathbf{X})| \leq C_1(\tau^t \mathbf{X})(1 - \zeta)^{n-t}$ for some $\zeta > 0$. With the function $F(x)$ introduced in assumption (A4), we define $C_2(x) := \exp(F(x))$ for $x \in \mathcal{X}$ and arrive at

$$|\mathbb{P}^{(0:\infty)}(Y_t = i \mid \mathbf{X}) - \mathbb{P}^{(0:n)}(Y_t = i \mid \mathbf{X})| \leq |\mathcal{Y}|^2 C_1(\tau^t \mathbf{X}) C_2(\mathbf{X}_t) (1 - \zeta)^{n-t}.$$

The next step is to construct a function $C_3(x)$ satisfying the following two conditions: (i) If $C_2(x)(1 - \zeta)^k \geq 1$, then $C_3(x)h(k) \geq 1$. (ii) If $C_2(x)(1 - \zeta)^k < 1$, then $C_3(x)h(k) \geq C_2(x)(1 - \zeta)^k$. Since the difference between two probabilities cannot exceed 1, we obtain

$$|\mathbb{P}^{(0:\infty)}(Y_t = i | \mathbf{X}) - \mathbb{P}^{(0:n)}(Y_t = i | \mathbf{X})| \leq |\mathcal{Y}|^2 C_1(\tau^t \mathbf{X}) C_3(X_t) h(n - t).$$

The last step is to show that $\mathbb{E}[|C_3(X_t)|] < \infty$. \square

The following result will play a key role in the later analysis of empirical likelihood functions.

Theorem 4. *Suppose that conditions (A1) - (A4) hold, and the function $g : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ satisfies $\mathbb{E}[|g(X_t, Y_t)|] < \infty$. Then*

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n \mathbb{E}^{(0:n)}[g(X_t, Y_t) | \mathbf{X}] = \mathbb{E}[g(X_t, Y_t)] \quad \mathbb{P}\text{-almost surely.}$$

Proof. Without loss of generality we may assume that g can be written as $g(x, y) = g(x)\mathbf{1}(y = i)$. Using the result from Theorem 3, we obtain

$$\left| \sum_{t=1}^n \mathbb{E}^{(0:n)}[g(X_t, Y_t) | \mathbf{X}] - \sum_{t=1}^n \mathbb{E}^{(0:\infty)}[g(X_t, Y_t) | \mathbf{X}] \right| \leq \sum_{t=1}^n |g(X_t)| |C(\tau^t \mathbf{X})| h(n - t),$$

where $h(z)$ is a function satisfying the conditions in assumption (A4). See the extended version of this paper in the supplementary material for more details. Using the facts that \mathbf{X} is ergodic and the expectations of $|g(X_t)|$ and $|C(\tau^t \mathbf{X})|$ are finite, we obtain

$$\lim_{n \rightarrow \infty} \frac{1}{n} \left| \sum_{t=1}^n \mathbb{E}^{(0:n)}[g(X_t, Y_t) | \mathbf{X}] - \sum_{t=1}^n \mathbb{E}^{(0:\infty)}[g(X_t, Y_t) | \mathbf{X}] \right| = 0.$$

By similar arguments to the proof of the first part of Theorem 3 one can show that the difference $|\mathbb{E}^{(0:\infty)}[g(X_t, Y_t) | \mathbf{X}] - \mathbb{E}[g(X_t, Y_t) | \mathbf{X}]|$ tends to 0 as $t \rightarrow \infty$. Thus,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \left| \sum_{t=1}^n \mathbb{E}^{(0:\infty)}[g(X_t, Y_t) | \mathbf{X}] - \sum_{t=1}^n \mathbb{E}[g(X_t, Y_t) | \mathbf{X}] \right| = 0.$$

Now, noting that $E[g(X_t, Y_t) | \mathbf{X}] = E[g(X_0, Y_0) | \tau^t \mathbf{X}]$ for every t , the theorem follows by the ergodicity of \mathbf{X} . \square

5 Maximum Likelihood learning and model identifiability

In this section we apply the previous results to analyze asymptotic properties of the empirical likelihood function. The setting is the following: Suppose that we observe 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{Y} conditional on \mathbf{X} follows a conditional Markov chain with fixed feature functions \mathbf{f} and unknown model weights λ_* . We assume that λ_* lies in some parameter space Θ , the structure of which will become important later. To emphasize the role of the model weights, we will use subscripts, e.g., \mathbb{P}_λ and \mathbb{E}_λ , to denote the conditional distributions. Our goal is to identify the unknown model weights from the finite samples, \mathbf{X}_n and \mathbf{Y}_n . In order to do so, introduce the shorthand notation $\mathbf{f}(\mathbf{x}_n, \mathbf{y}_n) = \sum_{t=1}^n \mathbf{f}(x_t, y_{t-1}, y_t)$ for $\mathbf{x}_n = (x_0, \dots, x_n)$ and $\mathbf{y}_n = (y_0, \dots, y_n)$. Consider the normalized conditional likelihood,

$$\mathcal{L}_n(\lambda) = \frac{1}{n} \left(\lambda^T \mathbf{f}(\mathbf{X}_n, \mathbf{Y}_n) - \log \sum_{\mathbf{y}_n \in \mathcal{Y}^{n+1}} \exp(\lambda^T \mathbf{f}(\mathbf{X}_n, \mathbf{y}_n)) \right).$$

Note that, in the context of finite Conditional Markov Chains, $\mathcal{L}_n(\lambda)$ is the exact likelihood of \mathbf{Y}_n conditional on \mathbf{X}_n . The Maximum Likelihood estimate of λ_* is given by

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

If $\mathcal{L}_n(\lambda)$ is strictly concave, then the arg max is unique and can be found using gradient-based search (see [14]). It is easy to see that $\mathcal{L}_n(\lambda)$ is strictly concave if and only if $|\mathcal{Y}| > 1$, and there exists a sequence \mathbf{y}_n such that at least one component of $\mathbf{f}(\mathbf{X}_n, \mathbf{y}_n)$ is non-zero. In the following, we study strong consistency of the Maximum Likelihood estimates, a property which is of central importance in large sample theory (see [15]). As we will see, a key problem is the identifiability and uniqueness of the model weights.

5.1 Asymptotic properties of the likelihood function

In addition to the conditions (A1)-(A4) stated earlier, we will make the following assumptions:

- (A5) The feature functions have finite second moments: $\mathbb{E}_{\lambda_*} [|f(X_t, Y_{t-1}, Y_t)|^2] < \infty$.
- (A6) The parameter space Θ is compact.

The next theorem establishes asymptotic properties of the likelihood function $\mathcal{L}_n(\lambda)$.

Theorem 5. *Suppose that conditions (A1)-(A6) are satisfied. Then the following holds true:*

- (i) *There exists a function $\mathcal{L}(\lambda)$ such that $\lim_{n \rightarrow \infty} \mathcal{L}_n(\lambda) = \mathcal{L}(\lambda)$ \mathbb{P}_{λ_*} -almost surely for every $\lambda \in \Theta$. Moreover, the convergence of $\mathcal{L}_n(\lambda)$ to $\mathcal{L}(\lambda)$ is uniform on Θ , that is, $\lim_{n \rightarrow \infty} \sup_{\lambda \in \Theta} |\mathcal{L}_n(\lambda) - \mathcal{L}(\lambda)| = 0$ \mathbb{P}_{λ_*} -almost surely.*
- (ii) *The gradients satisfy $\lim_{n \rightarrow \infty} \nabla \mathcal{L}_n(\lambda) = \mathbb{E}_{\lambda_*} [\mathbf{f}(X_t, Y_{t-1}, Y_t)] - \mathbb{E}_{\lambda} [\mathbf{f}(X_t, Y_{t-1}, Y_t)]$ \mathbb{P}_{λ_*} -almost surely for every $\lambda \in \Theta$.*
- (iii) *If the limit of the Hessian $\nabla^2 \mathcal{L}_n(\lambda)$ is finite and non-singular, then the function $\mathcal{L}(\lambda)$ is strictly concave on Θ . As a consequence, the Maximum Likelihood estimates are strongly consistent:*

$$\lim_{n \rightarrow \infty} \hat{\lambda}_n = \lambda_* \quad \mathbb{P}_{\lambda_*}\text{-almost surely.}$$

Proof. The statements are obtained analogously to Lemma 4-6 and Theorem 4 in [1], using the auxiliary results for Conditional Markov Chains with unbounded feature functions established in Theorem 1, Theorem 2, and Theorem 4. \square

Next, we study the asymptotic behaviour of the Hessian $\nabla^2 \mathcal{L}_n(\lambda)$. In order to compute the derivatives, introduce the vectors $\lambda_1, \dots, \lambda_n$ with $\lambda_t = \lambda$ for $t = 1, \dots, n$. This allows us to write $\lambda^T \mathbf{f}(\mathbf{X}_n, \mathbf{Y}_n) = \sum_{t=1}^n \lambda_t^T \mathbf{f}(X_t, Y_{t-1}, Y_t)$. Now, regard the argument λ of the likelihood function as a stacked vector $(\lambda_1, \dots, \lambda_n)$. Same as in [1], this gives us the expressions

$$\frac{\partial^2}{\partial \lambda_t \partial \lambda_{t+k}^T} \mathcal{L}_n(\lambda) = \frac{1}{n} \text{Cov}_{\lambda}^{(0:n)} [\mathbf{f}(X_t, Y_{t-1}, Y_t), \mathbf{f}(X_{t+k}, Y_{t+k-1}, Y_{t+k})^T | \mathbf{X}]$$

where $\text{Cov}_{\lambda}^{(0:n)}$ is the covariance with respect to the measure $\mathbb{P}_{\lambda}^{(0:n)}$ introduced before Theorem 3. Using these expressions, the Hessian of $\mathcal{L}_n(\lambda)$ can be written as

$$\nabla^2 \mathcal{L}_n(\lambda) = - \left(\sum_{t=1}^n \frac{\partial^2}{\partial \lambda_t \partial \lambda_t^T} \mathcal{L}_n(\lambda) + 2 \sum_{k=1}^{n-1} \sum_{t=1}^{n-k} \frac{\partial^2}{\partial \lambda_t \partial \lambda_{t+k}^T} \mathcal{L}_n(\lambda) \right).$$

The following theorem establishes an expression for the limit of $\nabla^2 \mathcal{L}_n(\lambda)$. It differs from the expression given in Lemma 7 of [1], which is incorrect.

Theorem 6. *Suppose that conditions (A1) - (A5) hold. Then*

$$\lim_{n \rightarrow \infty} \nabla^2 \mathcal{L}_n(\lambda) = - \left(\gamma_{\lambda}(0) + 2 \sum_{k=1}^{\infty} \gamma_{\lambda}(k) \right) \quad \mathbb{P}_{\lambda_*}\text{-almost surely}$$

where $\gamma_{\lambda}(k) = \mathbb{E}[\text{Cov}_{\lambda}(\mathbf{f}(X_t, Y_{t-1}, Y_t), \mathbf{f}(X_{t+k}, Y_{t+k-1}, Y_{t+k}) | \mathbf{X})]$ is the expectation of the conditional covariance (with respect to \mathbb{P}_{λ}) between $\mathbf{f}(X_t, Y_{t-1}, Y_t)$ and $\mathbf{f}(X_{t+k}, Y_{t+k-1}, Y_{t+k})$ given \mathbf{X} . In particular, the limit of $\nabla^2 \mathcal{L}_n(\lambda)$ is finite.

Proof. The key step is to show the existence of a positive measurable function $U_{\lambda}(k, \mathbf{x})$ such that

$$\lim_{n \rightarrow \infty} \sum_{k=1}^{n-1} \sum_{t=1}^{n-k} \left| \frac{\partial^2}{\partial \lambda_t \partial \lambda_{t+k}^T} \mathcal{L}_n(\lambda) \right| \leq \lim_{n \rightarrow \infty} \sum_{k=1}^{n-1} \mathbb{E}[U_{\lambda}(k, \mathbf{X})]$$

with the limit on the right hand side being finite. Then the rest of the proof is straight-forward: Theorem 4 shows that, for fixed $k \geq 0$,

$$\lim_{n \rightarrow \infty} \sum_{t=1}^{n-k} \frac{\partial^2}{\partial \lambda_t \partial \lambda_{t+k}^T} \mathcal{L}_n(\boldsymbol{\lambda}) = \gamma_{\boldsymbol{\lambda}}(k) \quad \mathbb{P}_{\boldsymbol{\lambda}_*}\text{-almost surely.}$$

Hence, in order to establish the theorem, it suffices to show that

$$\lim_{n \rightarrow \infty} \sum_{k=1}^{n-1} \left| \gamma_{\boldsymbol{\lambda}}(k) - \sum_{t=1}^{n-k} \frac{\partial^2}{\partial \lambda_t \partial \lambda_{t+k}^T} \mathcal{L}_n(\boldsymbol{\lambda}) \right| \leq \epsilon$$

for all $\epsilon > 0$. Now let $\epsilon > 0$ be fixed. According to Theorem 2 we have $\sum_{k=1}^{\infty} |\gamma_{\boldsymbol{\lambda}}(k)| < \infty$. Hence we can find a finite N such that

$$\lim_{n \rightarrow \infty} \sum_{k=N}^{n-1} |\gamma_{\boldsymbol{\lambda}}(k)| + \lim_{n \rightarrow \infty} \sum_{k=N}^{n-1} \mathbb{E}[U_{\boldsymbol{\lambda}}(k, \mathbf{X})] \leq \epsilon.$$

On the other hand, Theorem 4 shows that

$$\lim_{n \rightarrow \infty} \sum_{k=1}^{N-1} \left| \gamma_{\boldsymbol{\lambda}}(k) - \sum_{t=1}^{n-k} \frac{\partial^2}{\partial \lambda_t \partial \lambda_{t+k}^T} \mathcal{L}_n(\boldsymbol{\lambda}) \right| = 0.$$

For details on how to construct $U_{\boldsymbol{\lambda}}(k, \boldsymbol{x})$, see the extended version of this paper. \square

5.2 Model identifiability

Let us conclude the analysis by investigating conditions under which the limit of the Hessian $\nabla^2 \mathcal{L}_n(\boldsymbol{\lambda})$ is non-singular. Note that $\nabla^2 \mathcal{L}_n(\boldsymbol{\lambda})$ is negative definite for every n , so also the limit is negative definite, but not necessarily strictly negative definite. Using the result in Theorem 6, we can establish the following statement:

Corollary 1. *Suppose that assumptions (A1)-(A5) hold true. Then the following conditions are necessary for the limit of $\nabla^2 \mathcal{L}_n(\boldsymbol{\lambda})$ to be non-singular:*

- (i) *For each feature function $f(x, i, j)$, there exists a set $A \subset \mathcal{X}$ with $\mathbb{P}(X_t \in A) > 0$ such that, for every $x \in A$, we can find $i, j, k, l \in \mathcal{Y}$ with $f(x, i, j) \neq f(x, k, l)$.*
- (ii) *There does not exist a non-zero vector $\boldsymbol{\lambda}$ and a subset $A \subset \mathcal{X}$ with $\mathbb{P}(X_t \in A) = 1$ such that $\boldsymbol{\lambda}^T \mathbf{f}(x, i, j)$ is constant for all $x \in \mathcal{X}$ and $i, j \in \mathcal{Y}$.*

Condition (i) essentially says: features $f(x, i, j)$ must not be constant in i and j . Condition (ii) says that features must not be expressible as linear combinations of each other. Clearly, features violating condition (i) can be assigned arbitrary model weights without any effect on the conditional distributions. If condition (ii) is violated, then there are infinitely many ways for parameterizing the same model. In practice, some authors have found positive effects of including redundant features (see, e.g., [16]), however, usually in the context of a learning objective with an additional penalizer.

6 Conclusions

We have established ergodicity and various mixing properties of Conditional Markov Chains with unbounded feature functions. The main insight is that similar results to the setting with bounded feature functions can be obtained, however, under additional assumptions on the distribution of the observations. In particular, the proof of Theorem 2 has shown that the sequence of observations needs to satisfy conditions under which Hoeffding-type concentration inequalities can be obtained. The proof of Theorem 3 has revealed an interesting interplay between mixing rates, feature functions, and the tail behaviour of the distribution of observations. By applying the mixing properties to the empirical likelihood functions we have obtained necessary conditions for the Maximum Likelihood estimates to be strongly consistent. We see a couple of interesting problems for future research: establishing Central Limit Theorems for Conditional Markov Chains; deriving bounds for the asymptotic variance of Maximum Likelihood estimates; constructing tests for the significance of features; generalizing the estimation theory to an infinite number of features; finally, finding sufficient conditions for the model identifiability.

References

- [1] Sinn, M. & Poupart, P. (2011) Asymptotic theory for linear-chain conditional random fields. In *Proc. of the 14th International Conference on Artificial Intelligence and Statistics (AISTATS)*.
- [2] Lafferty, J., McCallum, A. & Pereira, F. (2001) Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proc. of the 18th IEEE International Conference on Machine Learning (ICML)*.
- [3] Sutton, C. & McCallum, A. (2006) An introduction to conditional random fields for relational learning. In: Getoor, L. & Taskar, B. (editors), *Introduction to Statistical Relational Learning*. Cambridge, MA: MIT Press.
- [4] Hofmann, T., Schölkopf, B. & Smola, A.J. (2008) Kernel methods in machine learning. *The Annals of Statistics*, Vol. 36, No. 3, 1171-1220.
- [5] Xiang, R. & Neville, J. (2011) Relational learning with one network: an asymptotic analysis. In *Proc. of the 14th International Conference on Artificial Intelligence and Statistics (AISTATS)*.
- [6] Seneta, E. (2006) *Non-Negative Matrices and Markov Chains. Revised Edition*. New York, NY: Springer.
- [7] Wainwright, M.J. & Jordan, M.I. (2008) Graphical models: exponential families, and variational inference. *Foundations and Trends[®] in Machine Learning*, Vol. 1, Nos. 1-2, 1-305.
- [8] Cornfeld, I.P., Fomin, S.V. & Sinai, Y.G. (1982) *Ergodic Theory*. Berlin, Germany: Springer.
- [9] Orey, S. (1991) Markov chains with stochastically stationary transition probabilities. *The Annals of Probability*, Vol. 19, No. 3, 907-928.
- [10] Hernández-Lerma, O. & Lasserre, J.B. (2003) *Markov Chains and Invariant Probabilities*. Basel, Switzerland: Birkhäuser.
- [11] Foguel, S.R. (1969) *The Ergodic Theory of Markov Processes*. Princeton, NJ: Van Nostrand.
- [12] Samson, P.-M. (2000) Concentration of measure inequalities for Markov chains and Φ -mixing processes. *The Annals of Probability*, Vol. 28, No. 1, 416-461.
- [13] Kontorovich, L. & Ramanan, K. (2008) Concentration inequalities for dependent random variables via the martingale method. *The Annals of Probability*, Vol. 36, No. 6, 2126-2158.
- [14] Sha, F. & Pereira, F. (2003) Shallow parsing with conditional random fields. In *Proc. of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics (HLT-NAACL)*.
- [15] Lehmann, E.L. (1999) *Elements of Large-Sample Theory*. New York, NY: Springer.
- [16] Hoefel, G. & Elkan, C. (2008) Learning a two-stage SVM/CRF sequence classifier. In *Proc. of the 17th ACM International Conference on Information and Knowledge Management (CIKM)*.