Gradient Computation In Linear-Chain Conditional Random Fields Using The Entropy Message Passing Algorithm

Velimir M. Ilić^{a,b,*}, Dejan I. Mančev^b, Branimir T. Todorović^b, Miomir S. Stanković^c

^aMathematical Institute of the Serbian Academy of Sciences and Arts, Kneza Mihaila 36, 11000 Beograd, Serbia ^bUniversity of Niš, Faculty of Sciences and Mathematics, Višegradska 33, 18000 Niš, Serbia ^cUniversity of Niš, Faculty of Occupational Safety, Čarnojevića 10a, 18000 Niš, Serbia

Abstract

The paper proposes a numerically stable recursive algorithm for the exact computation of the linear-chain conditional random field gradient. It operates as a forward algorithm over the log-domain expectation semiring and has the purpose of enhancing memory efficiency when applied to long observation sequences. Unlike the traditional algorithm based on the forward-backward recursions, the memory complexity of our algorithm does not depend on the sequence length. The experiments on real data show that it can be useful for the problems which deal with long sequences.

Keywords: conditional random fields, expectation semiring, forward-backward algorithm, gradient computation, graphical models, message passing, sum-product algorithm

Interpaper property
gradient. It operation
gradient. It operation
efficiency when any the memory compuseful for the profit was full for th Conditional random fields (CRFs) [17] are probabilistic discriminative classifiers which can be applied for labeling and segmenting sequential data. When compared with more traditional sequence labeling tools like hidden Markov models (HMMs), the CRFs offer the advantage by relaxing the strong independence assumptions required by HMMs. Additionally, CRFs avoid the label bias problem [17] exhibited by the maximum entropy Markov models and other conditional Markov models based on directed graphical models. However, these improvements are accompanied by a significant cost in time and space needed for the parameter estimation of CRF, especially for real-time problems like labeling very long sequences which appear in computer security [18], [28], bioinformatics [16], [19] and robot navigation systems [15].

The CRF parameter estimation is typically performed by some of the gradient methods, such as iterative scaling, conjugate gradient, or limited memory quasi-Newton methods [12], [17], [21], [23], [25]. All these methods require the computation of the likelihood gradient, which becomes computationally demanding as the sequence length and the number of classes increase. The standard method for gradient computation [17] is based on the internal computation of CRF marginal probabilities by use of the forward-backward (FB) algorithm.

The FB algorithm first appeared in two independent publications [4], [5], but it is better known from the subsequent papers

dejan.mancev@pmf.edu.rs (Dejan I. Mančev),

branimirtodorovic@yahoo.com (Branimir T. Todorović), miromir.stankovic@gmail.com (Miomir S. Stanković)

[2], [3]. It makes use of dynamic programming, running with the asymptotical time complexity $\mathcal{O}(N^2T)$ and with the memory complexity $\mathcal{O}(NT)$, where T denotes the sequence length and N denotes the number of states classified. In spite of the time efficiency, it becomes spatially demanding when the sequence length is exceptionally large [14]. Furthermore, when it is used for the linear-chain gradient computation as in [12], [17], [21], [23], [25] it requires the storage of all CRF transition matrices which increase the total memory complexity for $\mathcal{O}(N^2T).$

The memory complexity can be reduced with modifications of the FB algorithm such as the checkpointing algorithm [11], [24] or with the re-computation of the transition matrices every time they are used (see section 3.3.). However, these techniques increase the computational complexity, while the memory complexity still depends on the sequence length. Another possibility is the use of forward-only algorithm [6], [20], [22], for which the matrices can be computed in runtime. This algorithm runs with constant memory complexity but it is computationally inefficient since it runs with the computational complexity $\mathcal{O}(N^4T).$

In this paper we propose an algorithm for the exact computation of the linear-chain conditional random field gradient. The algorithm is derived as a forward algorithm over the introduced log-domain expectation semiring, which means that its recursive equations can be obtained if real sums and products from an ordinary FB are replaced with products and sums from the log-expectation semiring. Accordingly, it can be seen as a numerically stable version of our previously developed Entropy Message Passing algorithm (EMP) [13], and it will also be called the EMP. Unlike the standard procedure, the EMP does not compute each marginal separately, but computes the gradient in a single forward pass by use of double recursion.

^{*}Research supported by Ministry of Science and Technological Development, Republic of Serbia, Grants No. 174013 and 174026

^{*}Corresponding author. Tel.: +381112630170; fax: +381112186105. Email addresses: velimir.ilic@gmail.com (Velimir M. Ilić),

Since only the forward pass is needed, the EMP can be implemented with the memory complexity being independent of the sequence length, having the advantage over the FB when long sequences are used.

The paper is organized as follows: In section II we explain the *FB* algorithm which operates over a commutative semiring. In section III we introduce the problem of efficient computation of a linear-chain *CRF* gradient and review the standard method based on the *FB* algorithm. The algorithms based on the *EMP* are presented in section IV, where the complexity analysis is given. Finally, the experimental results are presented in section V where two methods are compared and the advantage of the EMP is discussed.

2. The forward-backward algorithm over a commutative semiring

Definition 1. A commutative semiring is a tuple $(\mathbb{K}, \oplus, \otimes, 0, 1)$ where \mathbb{K} is the set with operations \oplus and \otimes such that both \oplus and \otimes are commutative and associative and have identity elements in \mathbb{K} (0 and 1 respectively), and \otimes is distributive over \oplus .

Let $(\mathbb{K}, \oplus, \otimes, 0, 1)$ be a commutative semiring and let $y = \{y_0, \ldots, y_T\}$ be a set of variables taking values from the set \mathcal{Y} of cardinality *N*. We define the *local kernel* functions $u_t : \mathcal{Y}^2 \to \mathbb{K}$ for $t = 1, \ldots, T$, and the *global kernel* function $u : \mathcal{Y}^{T+1} \to \mathbb{K}$, assuming that the following factorization holds

$$u(\mathbf{y}) = \bigotimes_{i=1}^{T} u_i(y_{i-1}, y_i)$$
(1)

for all $\mathbf{y} = (y_0, \dots, y_T) \in \mathcal{Y}^{T+1}$.

The FB algorithm [26], [27] solves two problems

1. The marginalization problem: Computes the sum

$$v_t(y_t, y_{y+1}) = \bigoplus_{y_{\{k-1,k\}^c}} u(y) = \bigoplus_{y_{\{k-1,k\}^c}} \bigotimes_{i=1}^T u_i(y_{i-1}, y_i), \quad (2)$$

2. The normalization problem: Computes the sum

$$Z = \bigoplus_{\mathbf{y}} u(\mathbf{y}) = \bigoplus_{\mathbf{y}} \bigotimes_{i=1}^{T} u_i(y_{i-1}, y_i).$$
(3)

The FB recursively computes the forward vector

$$\alpha_i(y_i) = \bigoplus_{y_{0:i-1}} \bigotimes_{t=1}^{l} u_t(y_{t-1}, y_t), \tag{4}$$

which is initialized to

$$\alpha_0(y_0) = 1, \tag{5}$$

and recursively computed using

$$\alpha_i(y_i) = \bigoplus_{y_{i-1}} u_{i-1}(y_{i-1}, y_i) \otimes \alpha_{i-1}(y_{i-1})$$
(6)

and the backward vector

$$\beta_i(y_i) = \bigoplus_{s_{i+1:T}} \bigotimes_{t=i+1}^{I} u_t(y_{t-1}, y_t),$$
(7)

which is recursively computed using

$$\beta_i(y_i) = \bigoplus_{y_{i+1}} u_{i+1}(y_i, y_{i+1}) \otimes \beta_{i+1}(y_{i+1})$$
(8)

and initialized to

$$\mathbf{y}(y_T) = 1. \tag{9}$$

Once the forward α_{k-1} and backward β_k vectors are computed, we can solve the marginalization problem by use of the formula

 β_T

$$\bigoplus_{y_{\{k-1,k\}^c}} \bigotimes_{i=1}^T u_i(y_{i-1}, y_i) = \alpha_{k-1}(y_{k-1}) \otimes u_k(y_{k-1}, y_k) \otimes \beta_k(y_k)$$
(10)

The normalization problem can be solved with the forward pass only according to

$$\bigoplus_{\mathbf{y}} \bigotimes_{i=1}^{T} u_i(y_{i-1}, y_i) = \bigoplus_{y_T} \alpha_T(y_T).$$
(11)

3. Linear-Chain CRF Training using the Forward Backward Algorithm

Linear-chain CRFs are discriminative probabilistic models over observation sequences $\mathbf{x} = (x_1, \dots, x_T)$ and label sequences $\mathbf{y} = (y_1, \dots, y_T)$, defined with conditional probability

$$p(\mathbf{y}|\mathbf{x};\boldsymbol{\theta}) = \frac{1}{Z(\mathbf{x};\boldsymbol{\theta})} \prod_{i=1}^{T} \mathbf{e}^{\langle \boldsymbol{\theta}, f(y_{i-1}, y_i, \mathbf{x}, i) \rangle}.$$
 (12)

The symbol $\langle \cdot, \cdot \rangle$ denotes the scalar product between an *M*-dimensional parameter vector

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_1, \dots, \theta_M \end{bmatrix} \tag{13}$$

and the *feature vector* on position *i*

$$f(y_{i-1}, y_i, \mathbf{x}, i) = [f_1(y_{i-1}, y_i, \mathbf{x}, i), \dots, f_M(y_{i-1}, y_i, \mathbf{x}, i)].$$
(14)

The normalization factor

$$Z(\boldsymbol{x};\boldsymbol{\theta}) = \sum_{\boldsymbol{y}} \prod_{i=1}^{T} \mathbf{e}^{\langle \boldsymbol{\theta}, f(y_{i-1}, y_i, \boldsymbol{x}, i) \rangle}$$
(15)

is called the *partition function*.

The goal of the *CRF* training is to build up the model (12) from the data set $\{(\mathbf{x}^{(d)}, \mathbf{y}^{(d)})\}_{d=1}^{D}$. The standard method is to maximize the log likelihood of (12):

$$\mathcal{L}(\boldsymbol{\theta}) = \sum_{d=1}^{D} \ln p(\mathbf{y}^{(d)} | \mathbf{x}^{(d)}; \boldsymbol{\theta})$$
(16)

over the parameter vector $\boldsymbol{\theta}$ for the chosen set of feature vectors $f(y_{i-1}, y_i, \boldsymbol{x}, i)$. The maximum can be found with several of the gradient methods [12], [17], [21], [23], [25], which requires the computation of the gradient $\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta})$. The gradient can be expressed, according to (12) and (16), as:

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) = \sum_{d=1}^{D} \sum_{i=1}^{T^{(d)}} \boldsymbol{f}(\boldsymbol{y}_{i-1}^{(d)}, \boldsymbol{y}_{i}^{(d)}, \boldsymbol{x}^{(d)}, i) - \sum_{d=1}^{D} \frac{\nabla_{\boldsymbol{\theta}} \boldsymbol{Z}(\boldsymbol{x}; \boldsymbol{\theta})}{\boldsymbol{Z}(\boldsymbol{x}; \boldsymbol{\theta})}, \quad (17)$$

Figure 1: Forward-backward computation scheme.

where $T^{(d)}$ is the length of the *d*-th observation sequence.

The main problem in the evaluation of the log likelihood gradient (17) is the computation of the quotient between the partition function gradient and the partition function. The partition function gradient can be represented as

$$\nabla_{\boldsymbol{\theta}} Z(\boldsymbol{x}; \boldsymbol{\theta}) = [\nabla_{\theta_1} Z(\boldsymbol{x}; \boldsymbol{\theta}), \dots, \nabla_{\theta_M} Z(\boldsymbol{x}; \boldsymbol{\theta})], \quad (18)$$

where $\nabla_{\theta_m} Z(x; \theta)$ denotes the *m*-th partial derivative, and can be obtained from (15) after the use of the Leibniz's product rule:

$$\nabla_{\theta} Z(\boldsymbol{x}; \boldsymbol{\theta}) = \sum_{y_{0:T}} \prod_{i=1}^{T} \mathbf{e}^{\langle \boldsymbol{\theta}, f(y_{i-1}, y_i, \boldsymbol{x}, i) \rangle} \sum_{k=1}^{T} f(y_{k-1}, y_k, \boldsymbol{x}, k).$$
(19)

The standard method for the computation of the partition function and its gradient [17] is based on the forward-backward algorithm which is reviewed in the following section.

3.1. Sum-product semiring forward-backward algorithm

Definition 2. The sum-product semiring is the tuple $(\mathbb{R}, +, 1, 0)$, where \mathbb{R} is the set of real numbers and the operations defined in a standard way.

The partition function (15) can be obtained as the solution of the normalization problem (11) for factorization:

$$\prod_{i=1}^{T} \mathbf{e}^{\langle \theta, f(y_{i-1}, y_i, \mathbf{x}, i) \rangle}.$$
(20)

as

$$Z(\boldsymbol{x};\boldsymbol{\theta}) = \sum_{y_{0:T}} \prod_{i=1}^{T} \mathbf{e}^{\langle \boldsymbol{\theta}, f(y_{i-1},y_i,\boldsymbol{x},i) \rangle}$$
(21)

The gradient can be computed using the solution for the marginalization problem (10) in the sum-product semiring. First, we change the sum ordering in (19) and split the sum over y to $y_{\{k-1,k\}}$ and $y_{\{k-1,k\}}$ ^c sums, transforming (19) to:

$$\nabla_{\theta} Z(\boldsymbol{x}; \boldsymbol{\theta}) = \sum_{k=1}^{T} \sum_{y_{\{k-1,k\}}} \left(\sum_{y_{\{k-1,k\}^{c}}} \prod_{i=1}^{T} \boldsymbol{e}^{\langle \boldsymbol{\theta}, f(y_{i-1}, y_{i}, \boldsymbol{x}, i) \rangle} \right) \cdot \boldsymbol{f}(y_{k-1}, y_{k}, \boldsymbol{x}, k).$$
(22)

The marginal values,

$$\sum_{\mathcal{Y}_{\{k-1,k\}^c}} \prod_{i=1}^T \mathbf{e}^{\langle \boldsymbol{\theta}, f(y_{i-1}, y_i, \boldsymbol{x}, i) \rangle},$$
(23)

as the marginalization problem over the sum-product semiring, can be found by recursive computation of forward vectors,

$$\alpha_i(y_i) = \sum_{y_{0:i-1}} \prod_{t=1}^i \mathbf{e}^{\langle \boldsymbol{\theta}, f(y_{t-1}, y_t, \mathbf{x}, t) \rangle}$$
(24)

and backward vectors

$$\beta_i(y_i) = \sum_{y_{i+1:T}} \prod_{t=i+1}^{T} \mathbf{e}^{\langle \boldsymbol{\theta}, f(y_{t-1}, y_t, \mathbf{x}, t) \rangle}.$$
 (25)

The FB algorithm over the sum-product semiring suffers from numerical instability since the exponential terms can fall out of the machine precision scope and it is usually replaced with a more stable FB algorithm over the log-domain sumproduct semiring.

3.2. Log-domain sum-product semiring forward-backward algorithm

Definition 3. The log-domain sum-product semiring is the tuple $(\mathbb{R}^*, \oplus, \otimes, -\infty, 0)$ where \mathbb{R}^* is the extended set of real numbers and the operations are defined by

$$a \oplus b = \ln(\boldsymbol{e}^a + \boldsymbol{e}^b) \tag{26}$$

$$a \otimes b = a + b, \tag{27}$$

for all $a, b \in \mathbb{R}$ *.*

The following lemma follows straightforwardly from the definition of the log-domain sum-product semiring.

Lemma 1. Let $a_i \in \mathbb{R}$ for all $1 \le i \le T$. Then, the following equalities hold for log-domain sum-product semiring:

$$\ln\left(\sum_{i=1}^{T} a_i\right) = \bigoplus_{i=1}^{T} \ln a_i, \quad \ln\left(\prod_{i=1}^{T} a_i\right) = \bigotimes_{i=1}^{T} \ln a_i.$$
(28)

In log-domain the local kernels have the form:

$$u_i(y_{i-1}, y_i) = \langle \boldsymbol{\theta}, \, \boldsymbol{f}(y_{i-1}, y_i, \boldsymbol{x}, i) \rangle, \tag{29}$$

for i = 1, ..., T. According to Lemma 1 and expression (4), the forward vector in the log-domain sum-product semiring is the logarithm of the forward vector in the sum-product semiring:

$$\alpha_{i}(y_{i}) = \bigoplus_{y_{0:i-1}} \bigotimes_{t=1}^{i} u_{t}(y_{t-1}, y_{t}) =$$
$$= \ln \left(\sum_{y_{0:i-1}} \prod_{t=1}^{i} \mathbf{e}^{\langle \boldsymbol{\theta}, f(y_{t-1}, y_{t}, \mathbf{x}, t) \rangle} \right).$$
(30)

The forward vector α_0 is initialized to 0 which is the identity for \otimes ,

$$\alpha_0(y_0) = 0, \tag{31}$$

and it is recursively computed using

$$\alpha_i(y_i) = \bigoplus_{y_{i-1}} \left(u_i(y_{i-1}, y_i) + \alpha_{i-1}(y_{i-1}) \right).$$
(32)

Similarly to Lemma 1 and to expression (7), the backward vector in the log-domain sum-product semiring is the logarithm of the backward vector in sum-product semiring:

$$\beta_{i}(y_{i}) = \bigoplus_{y_{i+1:T}} \bigotimes_{t=i+1}^{T} \mathbf{e}^{\langle \theta, f(y_{t-1}, y_{t}, \mathbf{x}, t) \rangle} =$$
$$= \ln \sum_{y_{i+1:T}} \prod_{t=i+1}^{T} \mathbf{e}^{\langle \theta, f(y_{t-1}, y_{t}, \mathbf{x}, t) \rangle},$$
(33)

being initialized to

$$\beta_T(y_T) = 0, \tag{34}$$

and recursively computed using

$$\beta_i(y_i) = \bigoplus_{y_{i+1}} \left(u_{i+1}(y_i, y_{i+1}) + \beta_{i+1}(y_{i+1}) \right).$$
(35)

If the log-domain addition is performed using the definition, $a \oplus b = \ln(\mathbf{e}^a + \mathbf{e}^b)$, the numerical precision is being lost when computing \mathbf{e}^a and \mathbf{e}^b . But, as noted in ([23]), \oplus can be computed as

$$a \oplus b = a + \ln\left(1 + \mathbf{e}^{(b-a)}\right) = b + \ln\left(1 + \mathbf{e}^{(a-b)}\right),\tag{36}$$

which can be much more numerically stable, particularly if we pick the version of the identity with the smaller exponent.

The logarithm of the normalization constant (21) is according to Lemma 1

$$\ln Z(\boldsymbol{x};\boldsymbol{\theta}) = \ln \sum_{\boldsymbol{y}} \prod_{i=1}^{T} \mathbf{e}^{\langle \boldsymbol{\theta}, f(y_{i-1}, y_i, \boldsymbol{x}, i) \rangle} = \bigoplus_{\boldsymbol{y}} \bigotimes_{i=1}^{T} u_i(y_{i-1}, y_i),$$
(37)

and it can be computed using the solution of normalization problem in the log-domain sum-product semiring with forward algorithm according to (11)

$$\ln Z(\boldsymbol{x};\boldsymbol{\theta}) = \bigoplus_{y_T} \alpha_T(y_T). \tag{38}$$

According to Lemma 1, the marginal values (23) in the logdomain sum-product semiring have the form

$$v_{k}(y_{k-1}, y_{k}) = \ln \sum_{y_{\{k-1,k\}^{c}}} \prod_{i=1}^{T} \mathbf{e}^{\langle \theta, f(y_{i-1}, y_{i}, \mathbf{x}, i) \rangle} =$$
$$= \bigoplus_{y_{\{k-1,k\}^{c}}} \bigotimes_{i=1}^{T} u_{i}(y_{i-1}, y_{i}).$$
(39)

The marginal values can efficiently be computed according to the solution of the marginalization problems (10):

$$\psi_k(y_{k-1}, y_k) = \alpha_{k-1}(y_{k-1}) \otimes u_k(y_{k-1}, y_k) \otimes \beta_k(y_k), \quad (40)$$

where $\alpha_{k-1}(y_{k-1})$ and $\beta_k(y_k)$ are computed with the *FB* algorithm over the log-domain sum-product semiring, using equations (31)-(35). Then, by taking the logarithm of the *m*-th component in gradient expression (22), we get

$$\ln \nabla_{\theta_m} Z(\boldsymbol{x}; \boldsymbol{\theta}) =$$

$$= \bigoplus_{k=1}^T \bigoplus_{y_{\{k-1,k\}}} v_k(y_{k-1}, y_k) \otimes \ln f_m(y_{k-1}, y_k, \boldsymbol{x}, k), \quad (41)$$

for m = 1, ..., M. Finally, the quotient between the partition function and its gradient can be computed according to

$$\frac{\nabla_{\theta} Z(\boldsymbol{x}; \boldsymbol{\theta})}{Z(\boldsymbol{x}; \boldsymbol{\theta})} = \mathbf{e}^{\ln \nabla_{\theta} Z(\boldsymbol{x}; \boldsymbol{\theta}) - \ln Z(\boldsymbol{x}; \boldsymbol{\theta})}.$$
(42)

Algorithm 1: Log-domain FB algorithm

 $\begin{array}{c|c} \operatorname{input} : x, \ \theta, \ f(y_{k-1}, y_k, x, k); \ y_{k-1}, y_k \in \mathcal{Y}, \ k = 1, \dots, T; \\ \operatorname{output} : \nabla_{\theta} Z(x; \theta) / Z(x; \theta); \\ /* \ \text{Matrices initialization } */ \\ 1 \ \operatorname{for} k \leftarrow 1 \ \operatorname{to} T \ \operatorname{do} \\ 2 \ & foreach \ y_{k-1} \ in \ \mathcal{Y} \ \operatorname{do} \\ 3 \ & foreach \ y_k \ in \ \mathcal{Y} \ \operatorname{do} \\ 4 \ & u_k(y_{k-1}, y_k) = \sum_{m \in \mathcal{A}_k(y_{k-1}, y_k)} \theta_m \cdot f_m(y_{k-1}, y_k, x, k) \end{array}$

/* Forward phase */

foreach
$$v_0$$
 in \mathcal{Y} do

6
$$\alpha_0(y_0) \leftarrow 0;$$

7 for
$$k \leftarrow 1$$
 to T do

8 foreach
$$y_k$$
 in \mathcal{Y} do

9
$$\left\lfloor \alpha_k(y_k) \leftarrow \bigoplus_{y_{k-1}} \left(u_k(y_{k-1}, y_k) + \alpha_{k-1}(y_{k-1}) \right) \right.$$

/* Backward phase */

10 foreach
$$y_T$$
 in \mathcal{Y} do

11
$$\beta_T(y_T) \leftarrow 0;$$

12 for
$$k \leftarrow T - 1$$
 to 0 do
13 foreach y_k in \mathcal{Y} do
14 $\beta_k(y_k) \leftarrow \bigoplus_{y_{k+1}} (u_{k+1}(y_k, y_{k+1}) + \beta_{k+1}(y_{k+1}));$

 $\ln Z(\mathbf{w}, \boldsymbol{\theta}) = \Phi$ $\omega_{-}(\mathbf{w}_{-})$

$$\lim_{z \to \infty} Z(\mathbf{x}, \boldsymbol{\theta}) = \bigoplus_{y_T} u_T(y_T)$$

$$\inf_{z \to \infty} Z(\mathbf{x}, \boldsymbol{\theta}) = \bigoplus_{y_T} u_T(y_T)$$

$$\inf_{z \to \infty} Z(\mathbf{x}; \boldsymbol{\theta}) = \bigcup_{y_T} u_T(y_T)$$

3.3. Time and Memory Complexity

The time and memory complexity of the algorithm for the computation of the partition function and its derivatives by the FB algorithm is given in Table 1. The time complexity is defined as the number of operations required for the execution of the algorithm for a given pseudo code. In our analysis we consider real operations (addition and multiplication), log-domain

	\oplus	+	×	ln	Mem
		N ² T 4	N ² m 4) r) m
и	-	N^2TA	$N^2 T A$	-	N^2T
α	N^2T	N^2T	_	_	NT
β	N^2T	N^2T	_	_	NT
v	-	$2N^2T$	_	_	1
lnf	-	-	_	N^2TA	1
$\ln Z(x; \theta)$	Ν	_	-	-	1
$\ln \nabla_{\theta_m} Z(\boldsymbol{x}; \boldsymbol{\theta})$	N^2TA	N^2TA	_	_	М
Asymptotical	N^2TA	$2N^2TA$	N^2TA	N^2TA	$N^2T + M$

Table 1: Time and memory complexity of the log-domain FB algorithm.

operations (recall that log-domain multiplication is defined as real addition) and the number of computed logarithms. The memory complexity is defined as the number of 32-bit registers needed to store variables during algorithm execution. The complexity expressions are simplified by taking the quantities in expressions to tend to infinity, and keeping only the leading terms. In discussion, we will use big \mathcal{O} notation [7].

In applications, the feature functions $f(y_{k-1}, y_k, x, k)$ map the input space for a fixed sequence x into sparse vectors, which has nonzero values only at positions

$$\mathcal{A}_k(y_{k-1}, y_k) = \left\{ m \; ; \; f_m(y_{k-1}, y_k, \boldsymbol{x}, k) \text{ is nonzero } \right\}, \quad (43)$$

which allows the complexity reduction by performing the computation only for nonzero elements. In our analysis we will use the average number of nonzero elements defined as

$$A = \frac{\sum_{k=1}^{T} \mathcal{A}_k(y_{k-1}, y_k)}{T}.$$
 (44)

As Table 1 shows, the computationally most demanding part of the algorithm is the termination phase, which requires $\mathcal{O}(N^2TA)$ log-additions (recall that one log-addition requires the computation of the exponent and logarithm). The memory complexity of the algorithm is $\mathcal{O}(N^2T + M)$, governed by the space needed for storing the matrices u_i . The dependence of the memory complexity on the sequence length can significantly decrease computational performances of the algorithm if a long sequence is used, since it can cause overflows from the internal system memory to the disk storage, as shown in section 5.

The memory complexity can be reduced by the recomputation of the matrices u_i in the backward pass (line 14) and in the termination step (line 19), but this leads to the increased total number of additions and multiplications for $2N^2TA$, while the memory complexity still depends on the sequence length since all forward and backward vectors need to be stored. The further improvement can be achieved if one notes that the backward vectors are computed during the termination step since they are used only once in line 19. In this case, each backward vectors can be deleted after use in line 19 and all backward vectors can be stored at the memory location not depending on the sequence length. Then, the matrices u_i can be recomputed only once in the termination step, where they are used for the computation of the backward vector but, again, all forward vectors need to be stored and the memory complexity is $\mathcal{O}(NT + M)$, still depending on the sequence length.

The problem of memory complexity of the forwardbackward algorithm for the HMM has already been studied by Khreich et al. in [14]. In this paper they have proposed the algorithm for the computation of marginal probabilities called forward filtering backward smoothing (EFFBS), which runs with the memory complexity independent of the sequence length, $\mathcal{O}(N)$, with the same asymptotical computational complexity as the standard forward-backward algorithm. However, the algorithm is based on the HMM assumption that the transition matrix is constant, and as such cannot be applied to CRFs. Khreich et al. also gave a good review of the previously developed techniques for memory reduction such as checkpointing and forward-only algorithm, which try to reduce the memory complexity of the FB algorithm at the cost of computational overhead, and these techniques can be modified to deal with CRFs.

The checkpointing algorithm [11], [24] divides the input sequence into \sqrt{T} and during the forward pass only stores the first forward vector in each sub-sequence (checkpoint vectors). In the backward pass, the forward values for each sub-sequence are sequentially recomputed, beginning with checkpoint vectors. In this way, the computational complexity required for the computation of the forward and backward vectors is increased to $\mathcal{O}(2T - N^2 \sqrt{T})$, while the matrices u_i should also be recomputed, which leads to greater total computational cost. On the other hand, the memory complexity, although reduced to $\mathcal{O}(N\sqrt{T})$, still depends on the sequence length.

In the forward-only algorithm [6], [14], [20], [22], the expression of the form is obtained from three-dimensional matrices which are recursively computed. For the *HMM*, the computation can be realized in the constant memory space independent of the sequence length $\mathcal{O}(N^2 + N)$ and with time complexity $\mathcal{O}(N^4T)$. However, if it is applied to the *CRF*, its time complexity increases to $\mathcal{O}(N^4MT)$, which is significantly slower than the *FB* algorithm.

In the following section we derive a forward-only algorithm which operates with the time complexity of order $O(N^2(M + A)T)$, while keeping the memory complexity independent of

the sequence length.

4. Log-domain expectation semiring forward algorithm

In this section we consider a memory-efficient algorithm for CRF gradient computation which operates as a FB algorithm over an expectation semiring and develop its numerically stable log-domain version. In our previous work [13], we have developed the Entropy Message Passing EMP, which operates as a forward algorithm over the entropy semiring, which is the special case the expectation semiring. Although the algorithms presented in this paper are more general, in the following text they will be called the EMP, since they operate in the same manner as the algorithm from [13].

4.1. Expectation semiring forward algorithm

Definition 4. The expectation semiring of an order M is a tuple $\langle \mathbb{R} \times \mathbb{R}^M, \oplus, \odot, (0, \mathbf{0}), (1, \mathbf{0}) \rangle$, where the operations \oplus and \odot are defined with:

$$(z_1, \boldsymbol{h}_1) \oplus (z_2, \boldsymbol{h}_2) = (z_1 + z_2, \boldsymbol{h}_1 + \boldsymbol{h}_2), \qquad (45)$$

$$(z_1, \boldsymbol{h}_1) \odot (z_2, \boldsymbol{h}_2) = (z_1 z_2, z_1 \boldsymbol{h}_2 + z_2 \boldsymbol{h}_1),$$
 (46)

for all (z_1, h_1) , (z_2, h_2) from $\mathbb{R} \times \mathbb{R}^M$, and **0** denotes zero vector. The first component of an ordered pair is called a z-part, while the second one is an h-part.

For M = 1, the expectation semiring reduces to the entropy semiring considered in [13]. According to the addition rule, the z and h components of sum of two pairs are the sums of z and h components respectively, which gives us the following lemma.

Lemma 2. Let $(z_i, h_i) \in \mathbb{R} \times \mathbb{R}^M$ for all $1 \le i \le T$. Then, the following equality holds in the expectation semiring:

$$\bigoplus_{i=1}^{T} (z_i, z_i \boldsymbol{h}_i) = \left(\sum_{i=1}^{T} z_i, \sum_{i=1}^{T} \boldsymbol{h}_i\right).$$
(47)

Note that if the pairs have the form (z, zh), the multiplication acts as

$$(z_1, \boldsymbol{h}_1) \odot (z_2, \boldsymbol{h}_2) = (z_1 z_2, z_1 z_2 (\boldsymbol{h}_1 + \boldsymbol{h}_2).$$
(48)

This can be generalized with the following lemma.

Lemma 3. Let $(z_i, z_i h_i) \in \mathbb{R} \times \mathbb{R}^M$ for all $1 \le i \le T$. Then, the following equality holds in the expectation semiring:

$$\bigotimes_{i=1}^{T} (z_i, z_i \boldsymbol{h}_i) = \left(\prod_{i=1}^{T} z_i, \prod_{i=1}^{T} z_i \cdot \sum_{j=1}^{T} \boldsymbol{h}_j \right)$$
(49)

According to lemma 3, if the local kernels have the form:

$$u_{i}(y_{i-1}, y_{i}) = \left(\mathbf{e}^{\langle \boldsymbol{\theta}, f(y_{i-1}, y_{i}, \boldsymbol{x}, i) \rangle}, \mathbf{e}^{\langle \boldsymbol{\theta}, f(y_{i-1}, y_{i}, \boldsymbol{x}, i) \rangle} \cdot \boldsymbol{f}(y_{i-1}, y_{i}, \boldsymbol{x}, i)\right),$$
(50)

for $i = 1, \dots, T$, the global kernel is, according to the Lemma 3

$$\bigotimes_{i=1}^{T} u_i(y_{i-1}, y_i) = \Big(\prod_{i=1}^{T} \mathbf{e}^{\langle \boldsymbol{\theta}, f(y_{i-1}, y_i, \mathbf{x}, i) \rangle}, \qquad (51)$$
$$\prod_{i=1}^{T} \mathbf{e}^{\langle \boldsymbol{\theta}, f(y_{i-1}, y_i, \mathbf{x}, i) \rangle} \cdot \sum_{j=1}^{T} f(y_{j-1}, y_j, \mathbf{x}, j)\Big).$$

By applying lemma 2 to the expression (51), we can obtain the partition function (15),

$$Z(\boldsymbol{x};\boldsymbol{\theta}) = \sum_{\boldsymbol{y}} \prod_{i=1}^{T} \mathbf{e}^{\langle \boldsymbol{\theta}, f(y_{i-1}, y_i, \boldsymbol{x}, i) \rangle},$$
(52)

and its gradient (19):

$$\nabla_{\theta_m} Z(\boldsymbol{x}; \boldsymbol{\theta}) = \sum_{\boldsymbol{y}} \prod_{i=1}^{T} \mathbf{e}^{\langle \boldsymbol{\theta}, f(y_{i-1}, y_i, \boldsymbol{x}, i) \rangle} \cdot \sum_{j=1}^{T} f(y_{j-1}, y_j, \boldsymbol{x}, j), \quad (53)$$

as z and h parts of the sum

$$\bigoplus_{\mathbf{y}} \bigotimes_{i=1}^{T} u_i(y_{i-1}, y_i) = (Z(\mathbf{x}; \boldsymbol{\theta}), \nabla_{\boldsymbol{\theta}} Z(\mathbf{x}; \boldsymbol{\theta})).$$
(54)

The expression (54) can be computed as the normalization problem (11) by use of the forward algorithm over the expectation semiring. Note that z-parts of addition and multiplication acts as addition and multiplication in the sum-product semiring. Accordingly, the z-parts of forward vectors will be the same as the forward vectors in the sum-product semiring, and their computation is numerically unstable. In the following subsection we a develop numerically stable forward algorithm which operates over a log-domain expectation semiring.

4.2. Log-domain expectation semiring forward algorithm

The log-domain expectation semiring is a combination of the log-domain sum-product semiring and the expectation semiring. It can be obtained if real addition and multiplication in the definition of expectation semiring operations are replaced with their log-domain counterparts.

Before we define the log-domain expectation semiring, we introduce some usefull notation. Firstly, recall that log-domain addition and multiplication are defined with

$$a \oplus b = \ln(\mathbf{e}^a + \mathbf{e}^b) \tag{55}$$

$$a \otimes b = a + b. \tag{56}$$

The log-product between the scalar $z \in \mathbb{R}$ and the vector $\boldsymbol{h} = (h[1], \dots, h[M]) \in \mathbb{R}^M$ is defined as the vector $z \otimes h$:

$$z \otimes \boldsymbol{h} = z \otimes (h[1], \dots, h[M]) = (z \otimes h[1], \dots, z \otimes h[M]),$$
(57)

the logarithm of the vector $[h_1, \ldots, h_M] \in \mathbb{R}^M$ is defined as

$$\ln[h_1,\ldots,h_M] = [\ln h_1,\ldots,\ln h_M].$$
(58)

The vector $-\infty$ is defined as a vector all of whose coordinates are $-\infty$.

Figure 2: *EMP* computation scheme.

Definition 5. The log-domain expectation semiring of an order M is a tuple $\langle \mathbb{R} \times \mathbb{R}^M, \mathbb{O}, \mathbb{O}, (-\infty, -\infty), (0, -\infty) \rangle$, where the operations \mathbb{O} and \mathbb{O} are defined with:

$$(z_1, \boldsymbol{h}_1) \oplus (z_2, \boldsymbol{h}_2) = (z_1 \oplus z_2, \boldsymbol{h}_1 \oplus \boldsymbol{h}_2), \qquad (59)$$

$$(z_1, \boldsymbol{h}_1) \odot (z_2, \boldsymbol{h}_2) = (z_1 \otimes z_2, (z_1 \otimes \boldsymbol{h}_2) \oplus (z_2 \otimes \boldsymbol{h}_1)),$$
 (60)

for all (z_1, h_1) , (z_2, h_2) from $\mathbb{R} \times \mathbb{R}^M$. Similar to the expectation semiring, the first component of an ordered pair is called a *z*-part, while the second one is an h-part.

The following lemma is the log-domain version of Lemma 2.

Lemma 4. Let $(z_i, z_i h_i) \in \mathbb{R} \times \mathbb{R}^M$ for all $1 \le i \le T$. Then, the following equality holds in the log-domain expectation semiring:

$$\bigoplus_{i=1}^{T} (z_i, \boldsymbol{h}_i) = \Big(\bigoplus_{i=1}^{T} z_i, \bigoplus_{i=1}^{T} \boldsymbol{h}_i \Big),$$
(61)

where

$$\bigoplus_{i=1}^{T} a_i = \ln\left(\sum_{i=1}^{T} e^{a_i}\right).$$
(62)

Similar to the expectation semiring, if the pairs have the form $(z, z \otimes h)$ the multiplication acts as

$$(z_1, z_1 \otimes \boldsymbol{h}_1) \otimes (z_2, z_2 \otimes \boldsymbol{h}_2) = (z_1 \otimes z_2, z_1 \otimes z_2 \otimes (\boldsymbol{h}_1 \oplus \boldsymbol{h}_2)).$$
(63)

The following lemma is the log-domain version of Lemma 3.

Lemma 5. Let $(z_i, z_i \otimes h_i) \in \mathbb{R} \times \mathbb{R}^M$ for all $1 \le i \le T$. Then, the following equality holds in the log-domain expectation semiring:

$$\bigotimes_{i=1}^{T} (z_i, \ z_i \otimes \boldsymbol{h}_i) = \Big(\bigotimes_{i=1}^{T} z_i, \ \bigotimes_{i=1}^{T} z_i \otimes \bigoplus_{j=1}^{T} \boldsymbol{h}_j \Big), \tag{64}$$

where

$$\bigotimes_{i=1}^{T} a_i = \sum_{i=1}^{T} a_i.$$
 (65)

Let for $i = 1, \ldots, T$

$$\psi_i(y_i, y_{i-1}) = \langle \boldsymbol{\theta}, \boldsymbol{f}(y_{i-1}, y_i, \boldsymbol{x}, i) \rangle.$$
(66)

Then, the logarithm partition function (15) can be written as

$$\ln Z(\boldsymbol{x};\boldsymbol{\theta}) = \bigoplus_{\boldsymbol{y}} \bigotimes_{i=1}^{T} \psi_i(y_i, y_{i-1}).$$
(67)

The logarithm of the *m*-th partial derivative can be written as

$$\ln \nabla_{\theta_m} Z(\boldsymbol{x}; \boldsymbol{\theta}) = \ln \Big(\sum_{y_{0:T}} \prod_{i=1}^T \mathbf{e}^{\langle \boldsymbol{\theta}, f(y_{i-1}, y_i, \boldsymbol{x}, i) \rangle} \sum_{k=1}^T f_m(y_{k-1}, y_k, \boldsymbol{x}, k) \Big),$$
(68)

or, using the operations from the log-domain sum-product semiring,

$$\ln \nabla_{\theta} Z(\boldsymbol{x}; \boldsymbol{\theta}) = \bigoplus_{y_{0:T}} \bigotimes_{i=1}^{T} \psi_i(y_i, y_{i-1}) \otimes \bigoplus_{k=1}^{T} \ln \boldsymbol{f}(y_{k-1}, y_k, \boldsymbol{x}, k).$$
(69)

If the local kernels have the form:

$$u_{i}(y_{i-1}, y_{i}) = (\psi_{i}(y_{i}, y_{i-1}), \psi_{i}(y_{i}, y_{i-1}) \otimes \ln f(y_{i-1}, y_{i}, \mathbf{x}, i)),$$
(70)

for $i = 1, \dots, T$, the global kernel is, according to Lemma 5

$$\bigoplus_{i=1}^{T} u_i(y_{i-1}, y_i) = \left(\bigotimes_{i=1}^{T} \psi_i(y_i, y_{i-1}), \\ \bigotimes_{i=1}^{T} \psi_i(y_i, y_{i-1}) \otimes \bigoplus_{j=1}^{T} \ln f(y_{j-1}, y_j, \boldsymbol{x}, j)\right). \quad (71)$$

Furthermore, Lemma (4) for addition in the expectation semiring implies that the sum of ordered pairs is the ordered pair of the sums so the partition function and its gradient can be found as the z and h part of the sum:

$$\bigoplus_{\mathbf{y}} \bigoplus_{i=1}^{T} u_i(y_{i-1}, y_i) = (\ln Z(\mathbf{x}; \boldsymbol{\theta}), \ln \nabla_{\boldsymbol{\theta}} Z(\mathbf{x}; \boldsymbol{\theta})).$$
(72)

Expression (72) can be computed as the normalization problem (11) by use of the forward algorithm over the log-domain expectation semiring (*log-domain EMP algorithm*). The forward algorithm is initialized to the log-domain expectation semiring identity for the multiplication:

$$\alpha_0(y_0) = (0, -\infty),$$
 (73)

for all $y_0 \in \mathcal{Y}$. After that, we compute other forward vectors using the recurrent formula

$$\alpha_i(y_i) = \bigoplus_{y_{i-1}} u_i(y_{i-1}, y_i) \odot \alpha_{i-1}(y_{i-1}),$$
(74)

where the local factors are given with (70).

According to the rules for the addition and multiplication in the expectation semiring, the z and h parts of recursive equation (74) are:

$$\alpha_i^{(z)}(y_i) = \bigoplus_{y_{i-1}} \psi_i(y_i, y_{i-1}) \otimes \alpha_{i-1}^{(z)}(y_{i-1})$$
(75)

$$\alpha_{i}^{(h)}(y_{i}) = \bigoplus_{y_{i-1}} \psi_{i}(y_{i}, y_{i-1}) \otimes \alpha_{i-1}^{(h)}(y_{i-1}) \oplus \bigoplus_{y_{i-1}} \psi_{i}(y_{i}, y_{i-1}) \otimes \alpha_{i-1}^{(z)}(y_{i-1}) \otimes \ln f(y_{i-1}, y_{i}, \mathbf{x}, i)$$
(76)

for each $y_i \in \mathcal{Y}$, i = 1, ..., T. Finally, the normalization problem can be solved by summation

$$\bigoplus_{y_{0:T}} \bigoplus_{i=1}^{T} u_i(y_{i-1}, y_i) = \bigoplus_{y_T} \alpha_T(y_T),$$
(77)

whose z part is a partition function and the h part is its gradient:

$$\ln Z(\boldsymbol{x};\boldsymbol{\theta}) = \bigoplus_{y_T} \alpha_T^{(z)}(y_T), \quad \ln \nabla_{\boldsymbol{\theta}} Z(\boldsymbol{x};\boldsymbol{\theta}) = \bigoplus_{y_T} \alpha_T^{(h)}(y_T).$$
(78)

Hence, the algorithm consists of two parts: i) *forward pass*, at which the forward vectors are initialized according to (73) and recursively computed by (75)-(76), and during each step the corresponding matrix u_i is computed and ii) *termination*, at which the final summation of the forward algorithm is performed according to (78) and the partition function and its derivatives are obtained.

Figure 2 describes the EMP computation scheme. Recall that the forward-backward based computation requires that all forward and backward vectors be computed and stored until the partition function and the derivatives are obtained in the termination step. When the EMP is used, the computation terminates when the last forward vector is computed by use of the formulas (75) and (76). This can be realized in the fixed memory space with the size independent of the sequence length since the vectors $\alpha_{i-1}^{(z)}$, $\alpha_{i-1}^{(h)}$ and the matrices ψ_i should be computed only once in i-1-th iteration and, after having been used for the computation of $\alpha_i^{(z)}$ and $\alpha_i^{(h)}$, they can be deleted. The pseudo code is given in the table Algorithm 2. Here, the computation is performed using only two pairs of vectors $(\hat{\alpha}^{(z)}, \hat{\alpha}^{(h)})$ and $(\alpha^{(z)}, \alpha^{(h)})$. Note that the coordinates of the *h*-parts, $\hat{\alpha}^{(h)}(y_i)$ and $\alpha^{(h)}(y_i)$, are vectors which carry the information about the gradient and the *m*-th components of these vectors are denoted with $\hat{\alpha}_{[m]}^{(h)}(y_i)$ and $\alpha_{[m]}^{(h)}(y_i)$.

In comparison to the FB algorithm which needs the memory size of $\mathcal{O}(N^2T + M)$, the *EMP* has a memory complexity $\mathcal{O}(N^2 + NM)$, no longer depending on the sequence length T as in the FB algorithm. The additional cost is paid in time complexity which is increased for the term N^2TM . This is the consequence of the non-sparse computation of the EMP h component in line 13. Recall that the FB can be completely sparse implemented and, since $A \ll M$ in most of the application, the FB time complexity is lower. However, the sparsity can be reduced using the conditionally trained hidden Markov model assumption considered in [29], which we used in our implementation. With the reduced sparsity, the time complexity of the *EMP* is decreased and it becomes closer to the *FB* algorithm. When long sequences are used, the *EMP* becomes dominating since the FB needs to use the external memory. This assertion is justified in the following section, where we compare the two algorithms on a real data example.

5. Experiments

The intrusion detector learning task is to build a predictive model capable of distinguishing between "bad" connections, Algorithm 2: Log-domain EMP algorithm

input :
$$x, \theta, f(y_{j-1}, y_j, x, j); j = 1, ..., T, y_{j-1}, y_j \in \mathcal{Y};$$

output: $\nabla_{\theta} Z(x; \theta) / Z(x; \theta)$;

/* Forward algorithm */

1 foreach y_0 in \mathcal{Y} do $\alpha^{(z)}(y_0) \leftarrow 1$ 2 for $m \leftarrow 1$ to M do 3 $\alpha_{[m]}^{(h)}(y_0) \leftarrow -\infty;$ 4 5 for $i \leftarrow 1$ to T do foreach y_i in \mathcal{Y} do 6 **foreach** y_{i-1} in \mathcal{Y} **do** 7 $\psi(y_{i-1}, y_i) = \langle \boldsymbol{\theta}, \boldsymbol{f}(y_{i-1}, y_i, \boldsymbol{x}, i) \rangle;$ 8 **foreach** y_i in \mathcal{Y} **do** 9 $\hat{\alpha}^{(z)}(y_i) = \bigoplus_{y_{i-1}} (\psi(y_{i-1}, y_i) + \alpha^{(z)}(y_{i-1}));$ 10 for $m \leftarrow 1$ to M do 11 12 for each y_i in \mathcal{Y} do $\hat{\alpha}_{[m]}^{(h)}(y_i) \leftarrow \bigoplus_{y_{i-1}} (\psi(y_{i-1}, y_i) + \alpha_{[m]}^{(h)}(y_{i-1}));$ 13 foreach y_{i-1} in \mathcal{Y} do 14 foreach y_i in \mathcal{Y} do 15 $\gamma(y_{i-1}) \leftarrow \psi(y_{i-1},y_i) + \alpha^{(z)}(y_{i-1});$ 16 **foreach** *m* in $\mathcal{A}(y_{i-1}, y_i)$ **do** 17 $\begin{array}{l} lnf \leftarrow \ln f_m(y_{i-1}, y_i, \boldsymbol{x}, i) \\ \hat{\alpha}_{[m]}^{(h)}(y_i) \leftarrow \hat{\alpha}_{[m]}^{(h)}(y_i) \oplus (\gamma(y_{i-1}) + lnf); \end{array}$ 18 foreach y_i in \mathcal{Y} do 19 $\alpha^{(z)}(y_i) \leftarrow \hat{\alpha}^{(z)}(y_i);$ 20 21 22 /* Termination */ 23 $\ln Z \leftarrow \bigoplus_{y_T} \alpha^{(z)}(y_T)$ 24 $\ln \nabla_m Z \leftarrow \bigoplus_{y_T} \alpha^{(h)}_{[m]}(y_T)$ **25** for $m \leftarrow 1$ to M do $\nabla_{\theta_m} Z(\boldsymbol{x}; \boldsymbol{\theta}) / Z(\boldsymbol{x}; \boldsymbol{\theta}) \leftarrow \mathbf{e}^{\ln \nabla_m Z - \ln Z}$

called intrusions or attacks, and "good" normal connections. Conditional random fields have proven to be very effective in detecting intrusion [30].

As we have already mentioned, in the standard CRF training based on the FB algorithm the storage requirements are high when long train sequences are used. This may cause overflows from the internal system memory to disk storage which decreases computational performances, since accessing paged memory data on a typical disk drive is significantly slower than accessing data in RAM [14],[28]. On the other hand, the *EMP* runs with a small fixed memory and it becomes preferable for long sequences.

In Figure 3 we show the time and memory usage of both al-

	\oplus	+	×	ln	Mem
ψ	_	N^2TA	N^2TA	_	N^2
$\hat{lpha}^{(z)}$	N^2T	N^2T	_	_	Ν
γ	-	N^2T	-	-	1
$\hat{\alpha}_m^{(h)}$	$N^2T(M+A)$	$N^2T(M+A)$	_	N^2TA	NM
$\alpha^{(z)}$	_	-	_	_	Ν
lnf	-	_	-	N^2TA	1
$\alpha_r^{(h)}$	_	_	_	_	NM
$\ln Z(x; \theta)$	Ν	-	-	_	1
$\ln \nabla_{\theta_m} Z(\boldsymbol{x}; \boldsymbol{\theta})$	NM	-	_	_	M
Asymptotical	$N^2T(M+A)$	$N^2T(M+2A)$	N^2TA	N^2TA	$N^2 + NM$

Table 2: Time and memory complexity of the log-domain *EMP* algorithm.



Figure 3: CPU and RAM usage of FB and EMP algorithms

gorithms as functions of the sequence length. The experiments are performed using a computer with 3GB RAM and IntelCore 2 Duo CPU 2.33GHz. In our experiments, we used a KDE corpus [31] for sequences up to 5 million, while the sequences longer than 5 million are created by the concatenation of the KDE corpus on itself. We consider four different implementation cases depending on the sequence length:

Case I: This corresponds to short sequences, with the length shorter than 4 million. This case corresponds to the basic version of the FB algorithm (Algorithm 1). In this case, the se-

quence is stored in RAM and the RAM usage of both algorithms linearly grows with the sequence length (Figure 3a). However, the *FB* algorithm uses $O(N^2T + M)$ memory for storing intermediate results and its RAM usage grows faster in comparison to the *EMP*, which needs fixed-size additional space $O(N^2 + NM)$. RAM usage growth reflects on the computational performances of *FB* algorithm, which runs faster then the *EMP* for the sequences with the length up to 4.5 million. As Figure 3a shows, at the sequence length of 3 millions *FB* RAM usage becomes considerable and the *FB* growth becomes nonlinear due to the memory paging. Finally, at the sequence length of about 4.5 million the *EMP* becomes faster than *FB*. One possibility for *FB* memory reduction is recomputation of transition matrices which is done in the Case II.

Case II: This corresponds to middle length sequences, between 4 and 5 million. At this case the sequence and all intermediate results are stored in RAM, but the transition matrices are recomputed every time they are used. Similar to the Case I, as the sequence becomes longer, the memory required for storing the forward vectors increases and, for sequences longer then 5 million, *FB* algorithm becomes slower than the *EMP* (see Figure 3b).

Case III: This corresponds to long sequences between 5 and 25 million. As in Case II, transition matrices are recomputed and all another intermediate results are stored in RAM, but the sequence cannot fit in RAM and needs to be stored on the secondary memory. In the FB the sequences have to be read twice from secondary memory, once in the forward and once in the backward phase. On the other hand, the *EMP* uses a single forward pass and reads the sequence only once, which makes it faster than the FB (Figure 3c).

Case IV: This corresponds to very long sequences loger than 25 million. In this case, similar to the Case *III*, the sequence is stored on the secondary memory. To avoid the *FB* performance decreasing due to a large number of intermediate variables stored in *RAM*, the portion of variables is stored on the secondary memory, which keeps the RAM usage constant, no longer dependent on the sequence length (Figure 3c). This increases the number of accessions to the secondary memory, which further decreases *FB* performances in comparison to Case *III*. On the other hand, the *EMP* does not need to store additional data on the secondary memory and has the same time growth as in Case *III*, while using a small constant memory.

The previous results can vary with different operating systems and used hardware. Nevertheless, the access to secondary memory is very expensive operation and the algorithm with a low memory complexity has the advantage, when all data cannot fit in RAM, since the secondary memory accesses can be avoided.

6. Conclusion

In this paper, we have developed a numerically stable algorithm for the computation of the linear-chain CRF gradient. As opposed to the standard way of finding a CRF gradient by use of the forward-backward algorithm, the calculation by the proposed algorithm requires only the forward pass and can be realized with the memory independent of the observation sequence length. This makes the algorithm useful in the long sequence labeling tasks found in computer security [18], [28], bioinformatics [16],[19], and robot navigation systems [15].

The proposed algorithm operates as a forward algorithm over the log-domain expectation semiring, which can be seen as a modification of the expectation semiring used in the automata theory and probabilistic context free grammars [32], [33]. As mentioned in the paper, the use of the expectation semiring leads to numerically unstable algorithms and its log-domain counterpart can also be applied to numerically stable solutions of problems considered in [32], [33].

7. References

References

- S.M. Aji and R.J. McEliece. The generalized distributive law. *Information Theory, IEEE Transactions on*, vol. 46(2), pages 325 –343, 2000.
- [2] L. Bahl, J. Cocke, F. Jelinek, and J. Raviv. Optimal decoding of linear codes for minimizing symbol error rate (corresp.). *Information Theory, IEEE Transactions on*, vol. 20(2), pages 284 – 287, 1974.
- [3] L. Baum. An inequality and associated maximization technique in statistical estimation for probabilistic functions of Markov processes. *Inequalities*, vol. 3, pages 1–8, 1972.
- [4] Leonard E. Baum and Ted Petrie. Statistical inference for probabilistic functions of finite state markov chains. *The Annals of Mathematical Statistics*, vol. 37(6), pages 1554–1563, 1966.
- [5] R. Chang and J. Hancock. On receiver structures for channels having memory. *Information Theory, IEEE Transactions on*, vol. 12(4), pages 463 – 468, 1966.
- [6] Alexander Churbanov and Stephen Winters-Hilt. Implementing em and viterbi algorithms for hidden markov model in linear memory. BMC Bioinformatics, vol. 9(224), 2008.
- [7] Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, and Clifford Stein. *Introduction to Algorithms*. McGraw-Hill Science / Engineering / Math, 2nd edition, 2003.
- [8] Cortes C., Mohri M., Rastogi A., and Riley M., On the computation of the relative entropy of probabilistic automata, *International Journal of Foundations of Computer Science*, vol. 19, pages 219-242, 2007.
- [9] Robert G. Cowell, Philip A. Dawid, Steffen L. Lauritzen, and David J. Spiegelhalter. *Probabilistic Networks and Expert Systems (Information Science and Statistics)*. Springer, New York, 2003.
- [10] Eisner J., Parameter estimation for probabilistic finite-state transducers. in Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, pages 1–8, 2002.
- [11] J. Alicia Grice, Richard Hughey, and Don Speck. Reduced space sequence alignment. *Computer Applications in the Biosciences*, vol. 13(1), pages 45–53, 1997.
- [12] Rahul Gupta. Conditional random fields, 2006. Technique Report, IIT Bombay.
- [13] Velimir M. Ilić, Miomir S. Stanković, and Branimir T. Todorović. Entropy message passing algorithm. *IEEE Transactions on Information Theory*, vol 57(1), pages 375–380, 2011.
- [14] Wael Khreich, Eric Granger, Ali Miri, and Robert Sabourin. On the memory complexity of the forward-backward algorithm. *Pattern Recognition Letters*, vol 31(2), pages 91–99, 2010.
- [15] Sven Koenig and Reid G. Simmons. Unsupervised learning of probabilistic models for robot navigation. In *in Proceedings of the IEEE International Conference on Robotics and Automation*, pages 2301–2308, 1996.
- [16] Anders Krogh, Michael Brown, I. Saira Mian, Kimmen Sjlander, and David Haussler. Hidden markov models in computational biology: applications to protein modeling. *Journal of Molecular Biology*, vol. 235, pages 1501–1531, 1994.
- [17] John Lafferty, Andrew Mccallum, and Fernando Pereira. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In Proc. 18th International Conf. on Machine Learning, pages 282– 289, 2001.
- [18] Terran D. Lane. Machine learning techniques for the computer security domain of anomaly detection, PhD thesis, Purdue University, USA, 2000.
- [19] Irmtraud M. Meyer and Richard Durbin. Gene structure conservation aids similarity based gene prediction. *Nucleic acids research*, vol. 32(2), pages 776–783, 2004.
- [20] I. Miklós and I. M. Meyer. A linear memory algorithm for baum-welch training. BMC bioinformatics, vol. 6, 2005.
- [21] Fei Sha and Fernando Pereira. Shallow parsing with conditional random fields. In NAACL '03, vol.1, pages 213–220, 2003.

- [22] S. Sivaprakasam and K. Sam Shanmugan. A forward-only recursion based hmm for modeling burst errors in digital channels. *In GLOBECOM* '95, vol.2, pages 1054–1058, 1995.
- [23] Charles A. Sutton. *Efficient training methods for conditional random fields*. PhD thesis, 2008.
- [24] Christopher Tarnas and Richard Hughey. Reduced space hidden markov model training. *BIOINFORMATICS*, vol. 14(5), pages 401–406, 1998.
- [25] S. V. N. Vishwanathan, Nicol N. Schraudolph, Mark W. Schmidt, and Kevin P. Murphy. Accelerated training of conditional random fields with stochastic gradient methods. *In ICML*, pages 969–976, 2006.
- [26] F. R. Kschischang, B. J. Frey, and H.-A. Loeliger. Factor graphs and the sum-product algorithm. *IEEE Trans. Inform. Theory*, vol. 47, pp. 498-519, 2001.
- [27] S. M. Aji and R. J. McEliece. The generalized distributive law. *IEEE Trans. Inform. Theory*, vol. 46, pp. 325-343, Mar. 2000.
- [28] C. Warrender, S. Forrest, and B. Pearlmutter. Detecting intrusions using system calls: alternative data models. In IEEE Symposium on Security and Privacy, pages 133–145, 1999.
- [29] Shaolei Feng, R. Manmatha, and Andrew McCallum and Kevin P. Exploring the Use of Conditional Random Field Models and HMMs for Historical Handwritten Document Recognition. *DIAL '06*, pages 30-37, 2006.
- [30] Kapil Kumar Gupta, Baikunth Nath, and Kotagiri Ramamohanarao. Conditional Random Fields for Intrusion Detection. In AINAW '07, vol.1, pages 203-208, 2007.
- [31] KDE-CUP-99 Task description:
- http://kdd.ics.uci.edu/databases/kddcup99/task.html
- [32] J. Eisner, "Parameter estimation for probabilistic finite-state transducers," In Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, Philadelphia, July 2002, pp. 18.
- [33] Z. Li, J. Eisner. First- and second-order expectation semirings with applications to minimum-risk training on translation forests, *In Proceedings* of the 2009 Conference on Empirical Methods in Natural Language Processing, Volume 1 - Volume 1, EMNLP '09, pp. 40–51, Morristown, NJ, USA, 2009. Association for Computational Linguistics.