

DENSITY EVOLUTION FOR EXPECTATION PROPAGATION

John MacLaren Walsh, Ph. D.

Drexel University
Department of Electrical & Computer Engineering
Philadelphia, PA 19104

ABSTRACT

Expectation propagation (EP) [1, 2, 3, 4] is a theoretical extension of the belief propagation family of message passing algorithms [5, 6] for statistical inference which allows for efficient handling of models with continuous random variables as well as second or higher order correlation via the use of standard exponential families of probability measures [7, 8, 9]. Here we provide theoretically rigorous justifications for the use of density evolution [10, 11] to analyze the convergence and performance behavior of the family of algorithms in the large system regime by extending and expanding on the corresponding results for belief propagation decoding and turbo decoding.

Index Terms— expectation propagation, Bayes procedures, distributed iterative decoding and estimation, belief propagation

1. INTRODUCTION AND NOTATION

In this paper we will mathematically study the performance and convergence behavior of a family of iterative algorithms for Bayesian statistical inference known as expectation propagation (henceforth EP) [1, 2]. These algorithms extend belief propagation [5, 6] to efficiently handle continuous random variables and to allow for second and higher order statistical dependence among the random variables. Due to the rich family of algorithms which can be considered as special cases of belief propagation and thus EP, e.g. the turbo decoder/equalizer [13], the LDPC decoder [14], the Kalman filter, and the forward backward algorithm, a wealth of literature exists studying the performance and convergence properties of instances of EP.

One such body of convergence and performance analysis work which has had incredible theoretical impact is called density evolution [15, 11, 10]. Density evolution studies the performance and convergence of the turbo and LDPC decoders in the limit that the block length grows arbitrarily large. Of particular interest in this paper is the extension of the ideas of density evolution [15, 11, 10] from the decoding of turbo and LDPC codes to EP. With this in mind, we provide a review of EP in section 2, followed by a list of conditions under which density evolution can be applied to EP in section 3. The paper concludes with a list of possible practical applications of this theoretical work.

2. EXPECTATION PROPAGATION

EP is a distributed iterative method for statistical inference, which approximates the a posteriori distribution for some high-dimensional vector of parameters θ given observations \mathbf{r} . The parameters θ are known to be members of the set \mathcal{Q} which is the Cartesian product of the sets $\{\mathcal{Q}_i\}$

$$\mathcal{Q} = \mathcal{Q}_1 \times \mathcal{Q}_2 \times \cdots \times \mathcal{Q}_K$$

so that θ is a vector $\theta = [\theta_1, \theta_2, \dots, \theta_K]$, with $\theta_i \in \mathcal{Q}_i$ for all $i \in \{1, \dots, K\}$. We will assume that \mathcal{Q} has been endowed with a measure $d\theta$ that is a product of measures $d\theta_i$ on the \mathcal{Q}_i s. Typically, $d\theta_i$ will be a counting measure if \mathcal{Q}_i is countable, and will be Lebesgue measure otherwise.

In statistical inference problems to which EP may be applied, there is a multiplicative factoring of the Radon Nikodym derivative of the joint probability distribution for \mathbf{r} and θ with respect to $d\theta$,

$$p_{\mathbf{r},\theta}(\mathbf{r}, \theta) := \prod_{a=1}^A f_{a,\mathbf{r}}(\theta_a), \quad \theta_a \subset \theta$$

where $\theta_a \subset \theta$ indicates that θ_a is the (smaller dimensional) vector created by removing some of the elements of θ . EP exploits this multiplicative factoring in order to approximate $p_{\theta|\mathbf{r}}(\theta|\mathbf{r})$ with the product of minimal standard exponential family densities $g_{a,\lambda_a^{(k)}(\mathbf{r})}(\theta_a)$

$$p_{\theta|\mathbf{r}}(\theta|\mathbf{r}) \approx c(\lambda^{(k)}) \prod_{a=1}^A g_{a,\lambda_a^{(k)}(\mathbf{r})}(\theta_a) \quad (1)$$

Here $\lambda^{(k)} := [\lambda_1^{(k)}, \dots, \lambda_a^{(k)}, \dots, \lambda_A^{(k)}]$ is a vector of adjustable parameters which EP iteratively refines over time k to improve the approximation in (1) for a particular observed value $\mathbf{r} = \mathbf{r}$. Also, c is a normalization constant defined by

$$c(\lambda^{(k)}) := \left(\int_{\mathcal{Q}} \prod_{a=1}^A \exp(\lambda_{a,\mathbf{r}}^{(k)} \cdot \mathbf{t}_a(\theta_a)) d\theta \right)^{-1}$$

and the g s are standard exponential family probability densities [7, 8]

$$g_{a,\lambda_a^{(k)}(\mathbf{r})}(\theta_a) := \exp(\mathbf{t}_a(\theta_a) \cdot \lambda_a^{(k)}(\mathbf{r}) - \psi_{\mathbf{t}_a}(\lambda_a^{(k)}(\mathbf{r})))$$

with the log partition function defined by

$$\psi_{\mathbf{t}_a}(\lambda_a^{(k)}(\mathbf{r})) := \log \left(\int \exp(\mathbf{t}_a(\theta_a) \cdot \lambda_a^{(k)}(\mathbf{r})) d\theta_a \right)$$

EP attempts to choose the best approximation in (1) by iteratively refining the $\lambda_{a,\mathbf{r}}$ s by solving the equations

$$\begin{aligned} & \frac{\int \mathbf{t}_a(\theta_a) \exp(\lambda_{a,\mathbf{r}} \cdot \mathbf{t}_a(\theta_a)) \prod_{c \neq a} \exp(\mathbf{t}_c(\theta_c) \cdot \lambda_{c,\mathbf{r}}) d\theta}{\int \exp(\lambda_{a,\mathbf{r}} \cdot \mathbf{t}_a(\theta_a)) \prod_{c \neq a} \exp(\mathbf{t}_c(\theta_c) \cdot \lambda_{c,\mathbf{r}}) d\theta} \\ &= \frac{\int \mathbf{t}_a(\theta_a) f_a(\theta_a) \prod_{c \neq a} \exp(\mathbf{t}_c(\theta_c) \cdot \lambda_{c,\mathbf{r}}) d\theta}{\int f_a(\theta_a) \prod_{c \neq a} \exp(\mathbf{t}_c(\theta_c) \cdot \lambda_{c,\mathbf{r}}) d\theta} \end{aligned} \quad (2)$$

for $\lambda_{a,\mathbf{r}}$ in terms of $\{\lambda_{c,\mathbf{r}} | c \neq a\}$. The order in which these equations are solved is referred to as scheduling, and there are several

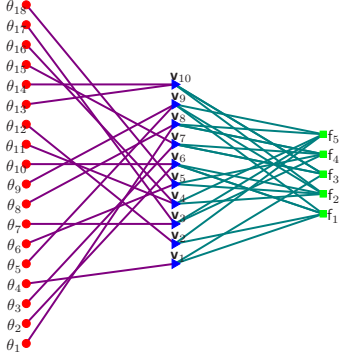


Fig. 1. A parameter basis factor graph.

possibilities. In **parallel scheduling**, (2) is solved to get $\lambda_{a,r}^{(k+1)}$ from $\lambda_{c,r}^{(k)}$, $c \neq a$ for each a . In **serial scheduling**, (2) is solved to get $\lambda_{a,r}^{(k+1)}$ from $\lambda_{c,r}^{(k)}$, $\lambda_{c',r}^{(k+1)}$, $c < a$, $c' > a$ for each a . Finally, in **random scheduling**, (2) is solved to update $\lambda_{a,r}$ from the most recently updated values of $\lambda_{c,r}$, $c \neq a$.

In this article, we will focus on parallel scheduling since we wish to emphasize the distributed, parallelized, nature of EP. In fact, under some additional assumptions, solving (2) iteratively can be equated to a message passing algorithm on a statistics factor graph. Situations in which these assumptions are satisfied will constitute the cases of interest for this article, so we wish to highlight them now.

As. 1 (Sufficiency): The factors $f_{a,r}(\theta_a)$ depend on the parameters only through $\mathbf{t}_a(\theta_a)$, so that $f_{a,r}(\theta_a) = \hat{f}_{a,r}(\mathbf{t}_a(\theta_a))$ for all θ_a .

Assumption 1 refocuses our interest on the functions $\{\mathbf{t}_a(\theta_a)\}$ which the factors depend on the parameters through. Typically, the same function $\mathbf{t}_i(\theta_i)$ will appear in several vector functions \mathbf{t}_a (i.e. with different a s). Take each element t_i of the vectors $\{\mathbf{t}_a | a \in \{1, \dots, A\}\}$ and collect them without repetition into the vector \mathbf{t} , so that $\mathbf{t}_a \subset \mathbf{t}$ for all a . Next, break \mathbf{t} up into the disjoint concatenation of the vectors \mathbf{v}_j with $j \in \{1, \dots, S\}$, in such a way so that $\mathbf{t}_i, \mathbf{t}_1 \in \mathbf{v}_j$ and $\mathbf{t}_1 \in \mathbf{t}_a$ implies that $\mathbf{t}_i \in \mathbf{t}_a$. We call the \mathbf{v}_j s elementary basis functions. Finally, denote the vector of arguments of \mathbf{v}_j by ϑ_j , so that $\vartheta_j \subset \theta$, and $\vartheta_j \subset \theta_a$ for all a such that $\mathbf{v}_j \subset \mathbf{t}_a$.

We may now depict the interdependence between the parameters $\{\theta_i\}$, the elementary basis functions $\{\mathbf{v}_j\}$, and the factors f_a with a tri-partite graph called a parameter basis factor graph. In this graph, the left nodes are the parameters $\{\theta_i\}$, the middle nodes are the elementary basis functions $\{\mathbf{v}_j\}$, and the right nodes are the factors f_a . An edge connects θ_i with \mathbf{v}_j if and only if $\theta_i \in \vartheta_j$. Similarly, an edge connects \mathbf{v}_j with f_a if and only if $\mathbf{v}_j \subset \mathbf{t}_a$. In the following development, we will denote the indices $a \in \{1, \dots, A\}$ of factor nodes which share an edge with the basis node \mathbf{v}_j by $\mathcal{F}(j)$. Furthermore, we will denote the indices $j \in \{1, \dots, S\}$ of the basis nodes which share an edge with the factor node f_a by $\mathcal{S}(a)$. We shall call the bipartite subgraph formed by only the basis and factor nodes and edges between them the basis factor graph. An example of a parameter basis factor graph is shown in figure 1.

As. 2 (Reciprocity): There is a product measure $d\mathbf{t} = \prod_i d\mathbf{v}_i$ on $\mathbf{t}(\mathcal{Q})$ such that the Radon Nikodym derivative μ of the in-

duced measure via the inverse image $\mathbf{t}^{-1}(\mathcal{A})$ from $d\theta$, with respect to the $d\mathbf{t}$ measure of \mathcal{A} factors into the product of functions of \mathbf{v}_j , i.e. $\mu(\mathbf{t}) = \prod_j \mu_j(\mathbf{v}_j)$.

Perhaps the easiest way to satisfy assumption 2 is to have each parameter θ_i appear in only one ϑ_j , so that the parameter nodes in the parameter basis factor graph all have degree one. This is the situation depicted in Figure 1.

Under as. 1 and 2, EP can be shown to be equivalent to a message passing algorithm on the basis factor graph. In particular, these assumptions imply that (2) simplifies to

$$\frac{\int \mathbf{v}_j(\vartheta_j) f_a(\theta_a) \exp \left(\sum_{1 \in \mathcal{S}(a)} \mathbf{v}_1(\vartheta_1) \cdot \mathbf{m}_{1 \rightarrow a} \right) d\theta_a}{\int f_a(\theta_a) \exp \left(\sum_{1 \in \mathcal{S}(a)} \mathbf{v}_1(\vartheta_1) \cdot \mathbf{m}_{1 \rightarrow a} \right) d\theta_a} = \frac{\int \mathbf{v}_j(\vartheta_j) \exp(\mathbf{v}_j(\vartheta_j) \cdot (\mathbf{n}_{a \rightarrow j} + \mathbf{m}_{j \rightarrow a})) d\vartheta_j}{\int \exp(\mathbf{v}_j(\vartheta_j) \cdot (\mathbf{n}_{a \rightarrow j} + \mathbf{m}_{j \rightarrow a})) d\vartheta_j} \quad (3)$$

where we have introduced the left going messages $\mathbf{n}_{a \rightarrow j} = [\lambda_{a,r}]_j$, that are the components of $\lambda_{a,r}$ multiplying the elements \mathbf{v}_j in \mathbf{t}_a . Equation (3) also introduces the right going messages

$$\mathbf{m}_{j \rightarrow a} = \sum_{c \in \mathcal{F}(j) \setminus \{a\}} \mathbf{n}_{a \rightarrow j} \quad (4)$$

Emphasizing again the sequential aspects of the algorithm, under the parallel scheduling, for each $a \in \{1, \dots, A\}$ equation (3) is solved to update $\mathbf{n}_{a \rightarrow j}$ for all $j \in \mathcal{S}(a)$. Then, for each $j \in \{1, \dots, S\}$ equation (4) is solved to update $\mathbf{m}_{j \rightarrow a}$ for all $a \in \mathcal{F}(j)$. This two step sequence is then repeated either for a fixed number of iterations or until convergence is reached.

3. DENSITY EVOLUTION FOR HOMOGENOUS EP

Here we provide an extension of density evolution analysis from belief propagation decoding [10] and turbo decoding [12, 11] to general EP. To simplify the results, we consider some special cases of EP which satisfy some additional assumptions, which collectively we will refer to as *the homogeneity assumptions*.

As. 3 (Regular Parameter Nodes): Each parameter in the parameter basis factor graph has degree one. In other word, each parameter θ_i appears as an argument of only one elementary basis vector \mathbf{v}_j .

As. 4 (Regular Basis Nodes): Each elementary basis vector \mathbf{v}_j appears in 1 degree 1 factor node (f_j) and d_p degree d_f factor nodes, and is made of the same function, so that $\mathbf{v}_j(\vartheta_j) = \mathbf{v}(\vartheta_j)$ for all ϑ_j for some vector function \mathbf{v} .

As. 5 (Symmetry): For the purposes of characterizing the performance and convergence of the particular instance of EP, it suffices to condition on $\theta = \theta_0 \in \mathcal{R}$ where

$$\mathcal{R} := \{\theta^0 | \vartheta_j = \theta_1 \forall j, 1 \in \{1, \dots, S\}\}$$

Assumption 5 is equivalent to the assumption that the all zero code word was transmitted in error control/correction decoding problems. In estimation problems, it relies on symmetry of the performance of the algorithm by requiring, say, that the distribution of the messages passed to \mathbf{v}_j after k iterations be the same when $\boldsymbol{\theta} = \boldsymbol{\theta}^0 \in \mathcal{R}$ as when $\boldsymbol{\theta} = \boldsymbol{\theta}^1$ not a member of \mathcal{R} but $\boldsymbol{\theta}^0 = \boldsymbol{\theta}^1$.

As. 6 (Marginal Messages): The marginal a posteriori density

$$\frac{\int f_{\mathbf{a}}(\mathbf{r}_{\mathbf{a}}, \boldsymbol{\theta}_{\mathbf{a}}) \prod_{i \in S(\mathbf{a}) \setminus \{j\}} \exp(\mathbf{v}(\boldsymbol{\theta}_i) \cdot \boldsymbol{\gamma}_i) d\boldsymbol{\theta}_c \setminus \boldsymbol{\theta}_j}{\int f_{\mathbf{a}}(\mathbf{r}_{\mathbf{a}}, \boldsymbol{\theta}_{\mathbf{a}}) \prod_{i \in S(\mathbf{a}) \setminus \{j\}} \exp(\mathbf{v}(\boldsymbol{\theta}_i) \cdot \boldsymbol{\gamma}_i) d\boldsymbol{\theta}_c}$$

is a minimal standard exponential family with sufficient statistics $\mathbf{v}(\cdot)$.

As. 7 (Independent Factors): The factors of degree d_f depend on different observations which are conditionally independent given $\boldsymbol{\theta} = \boldsymbol{\theta}_0$, so that the observation vector \mathbf{r} can be broken down into the disjoint¹ concatenation and reordering of smaller vectors $\{\mathbf{r}_{\mathbf{a}} | \mathbf{a} \in \{S+1, \dots, A\}\}$, which are conditionally i.i.d. given $\boldsymbol{\theta} = \boldsymbol{\theta}_0$

$$\mathbf{p}_{\mathbf{r}|\boldsymbol{\theta}}(\mathbf{r}|\boldsymbol{\theta}) = \prod_{\mathbf{a}=1}^A \mathbf{p}_{\mathbf{r}_{\mathbf{a}}|\boldsymbol{\theta}_{\mathbf{a}}}(\mathbf{r}_{\mathbf{a}}|\boldsymbol{\theta}_{\mathbf{a}})$$

with $f_{\mathbf{a},\mathbf{r}}$ a function of only $\mathbf{r}_{\mathbf{a}}$ and $\boldsymbol{\theta}_{\mathbf{a}}$, i.e. $f_{\mathbf{a},\mathbf{r}}(\boldsymbol{\theta}_{\mathbf{a}}) = f_{\mathbf{a},\mathbf{r}_{\mathbf{a}}}(\boldsymbol{\theta}_{\mathbf{a}})$, and $\mathbf{p}_{\mathbf{r}_{\mathbf{a}}|\boldsymbol{\theta}_{\mathbf{a}}}(\mathbf{r}_{\mathbf{a}}|\boldsymbol{\theta}_{\mathbf{a}}) = h(\mathbf{r}_{\mathbf{a}}|\boldsymbol{\theta}_{\mathbf{a}}) \forall \mathbf{r}_{\mathbf{a}}, \boldsymbol{\theta}_{\mathbf{a}}$ for all $\mathbf{a} \in \{S+1, \dots, A\}$.

As. 8 (Regular Factor Nodes): Each degree 1 factor node is the same function, so that

$$f_{\mathbf{a},\mathbf{r}_{\mathbf{a}}}(\boldsymbol{\theta}_{\mathbf{a}}) = w_{\mathbf{r}_{\mathbf{a}}}(\boldsymbol{\theta}_{\mathbf{a}}), \quad \hat{f}_{\mathbf{a},\mathbf{r}_{\mathbf{a}}}(\mathbf{t}_{\mathbf{a}}) = \hat{w}_{\mathbf{r}_{\mathbf{a}}}(\mathbf{t}_{\mathbf{a}})$$

for all $\mathbf{a} \in \{1, \dots, S\}$. Furthermore, each degree d_f factor $f_{\mathbf{a}}$ is the same function, so that

$$f_{\mathbf{a},\mathbf{r}_{\mathbf{a}}}(\boldsymbol{\theta}_{\mathbf{a}}) = f_{\mathbf{r}_{\mathbf{a}}}(\boldsymbol{\theta}_{\mathbf{a}}), \quad \hat{f}_{\mathbf{a},\mathbf{r}_{\mathbf{a}}}(\mathbf{t}_{\mathbf{a}}) = \hat{f}_{\mathbf{r}_{\mathbf{a}}}(\mathbf{t}_{\mathbf{a}})$$

for some functions f and \hat{f} for all \mathbf{a} . Also, the functions f and \hat{f} are insensitive to the ordering of the elementary basis vectors, so that

$$f_{\mathbf{r}_{\mathbf{a}}}(\boldsymbol{\pi}_{\boldsymbol{\theta}}(\boldsymbol{\theta}_{\mathbf{a}})) = f_{\mathbf{r}_{\mathbf{a}}}(\boldsymbol{\theta}_{\mathbf{a}}), \quad \hat{f}_{\mathbf{r}_{\mathbf{a}}}(\boldsymbol{\pi}_{\mathbf{v}}(\mathbf{t}_{\mathbf{a}})) = \hat{f}_{\mathbf{r}_{\mathbf{a}}}(\mathbf{t}_{\mathbf{a}})$$

where $\boldsymbol{\pi}_{\boldsymbol{\theta}}$ is a permutation on the elementary basis vector ordering of $\boldsymbol{\theta}_{\mathbf{a}}$, i.e. $\boldsymbol{\pi}_{\boldsymbol{\theta}}([\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_{S_{\mathbf{a}}}]) = [\boldsymbol{\theta}_{j_1}, \boldsymbol{\theta}_{j_2}, \dots, \boldsymbol{\theta}_{j_{S_{\mathbf{a}}}}]$ and $\boldsymbol{\pi}_{\mathbf{v}}([\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_{S_{\mathbf{a}}}]) = [\mathbf{v}_{j_1}, \mathbf{v}_{j_2}, \dots, \mathbf{v}_{j_{S_{\mathbf{a}}}}]$ for $(j_1, \dots, j_{S_{\mathbf{a}}})$ is any reordering of the integers from 1 to $S_{\mathbf{a}}$.

As. 9 (Randomly Chosen Edges): The subset of edges in the basis factor graph connected to degree d_f factor nodes are chosen by picking randomly from $\mathcal{G}(N, K, d_p, d_f)$, the set of bipartite graphs with N left nodes each of degree d_p and K right nodes each of degree d_f .

To highlight the purpose of these assumptions we must first become acquainted with some terminology. Define the *computation neighborhood* of a basis node $\{\mathbf{v}_j\}$ ² $\mathcal{C}_k(\mathbf{v}_j)$ of depth k to be sub-graph of the basis factor graph of nodes and edges no more than k edges away from \mathbf{v}_j . An example of a computation neighborhood is shown in Fig. 2.

¹Note that the degree 1 factor nodes have indices $\mathbf{a} \in \{1, \dots, S\}$ while the degree d_f factor nodes have indices $\mathbf{a} \in \{S+1, \dots, A\}$.

²This is adapted from the term *decoding neighborhood* from [10] because we wish to consider estimation problems as well as decoding problems.

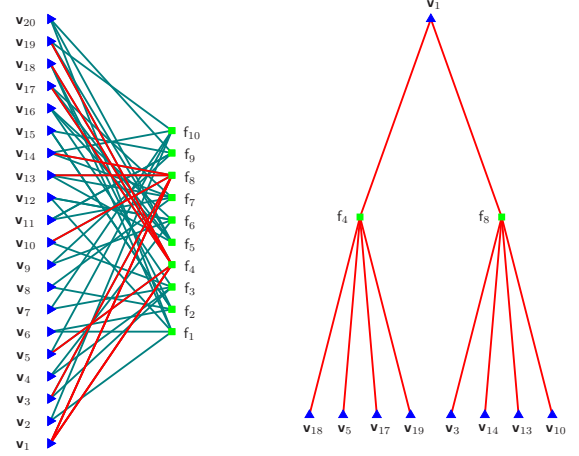


Fig. 2. The computation neighborhood of \mathbf{v}_1 of length 2.

The computation neighborhood is an important concept in EP because $\mathcal{C}_{2k}(\mathbf{v}_j)$ contains all of the factor nodes which influence the estimate of the a posteriori density provided by the EP algorithm after k iterations.

As we shall prove in the upcoming theorem, As. 3-9 see to it that, with probability $\rightarrow 1$ as the size of the parameter basis factor graph $K \rightarrow \infty$, for each k , the messages passed in the basis factor graph during the k th iteration are identically distributed according to a single probability distribution. Furthermore, again with probability $\rightarrow 1$ as $K \rightarrow \infty$, at any given iteration k all of the incoming messages to any particular basis node \mathbf{v}_j (i.e. all incoming messages at the root of the computation neighborhood $\mathcal{C}_k(\mathbf{v}_j)$) are independent. This allows one to determine the performance of EP at any given parameter node by keeping track of the way a single probability distribution evolves over iterations.

To state this result mathematically, it is necessary to introduce the operators \otimes^{d_p-1} , $\mathbf{y} \otimes$, $\mathfrak{F}_{\boldsymbol{\theta}_0}$ which we presently define. Given a probability density function \mathbf{g} , we will denote by $\otimes^{d_p-1} \mathbf{g}$ the $d_p - 1$ fold convolution of \mathbf{g} with itself. Furthermore, denote by $\mathbf{y} \otimes$ convolution with the probability density, given $\boldsymbol{\theta} = \boldsymbol{\theta}_0 \in \mathcal{R}$, for the vector

$$\boldsymbol{\Lambda}_{\mathbf{v}}^{-1}(\boldsymbol{\lambda}) \left(\frac{\int \mathbf{v}(\boldsymbol{\theta}_j) w(\mathbf{r}_j, \boldsymbol{\theta}_j) d\boldsymbol{\theta}_j}{\int w(\mathbf{r}_j, \boldsymbol{\theta}_j) d\boldsymbol{\theta}_j} \right)$$

where w was the function for the degree one factor nodes. Here, we have introduced the canonical transformation from the theory of statistical exponential families [7] which is bijective for minimal sufficient exponential families

$$\boldsymbol{\Lambda}_{\mathbf{v}_j}(\boldsymbol{\lambda}) := \frac{\int \mathbf{v}_j(\boldsymbol{\theta}_j) \exp(\mathbf{v}_j(\boldsymbol{\theta}_j) \cdot \boldsymbol{\lambda}) d\boldsymbol{\theta}_j}{\int \exp(\mathbf{v}_j(\boldsymbol{\theta}_j) \cdot \boldsymbol{\lambda}) d\boldsymbol{\theta}_j}$$

Finally, denote by $\mathfrak{F}_{\mathbf{a} \rightarrow j, \boldsymbol{\theta}_0}$ the map which takes the probability density $\mathbf{g}_{\boldsymbol{\lambda}}$ to the probability density, given $\boldsymbol{\theta} = \boldsymbol{\theta}_0 \in \mathcal{R}$, for the vector

$$\boldsymbol{\Lambda}_{\mathbf{v}_j}^{-1} \left(\frac{\int_{\mathcal{Q}_{\mathbf{a}}} \mathbf{v}_j(\boldsymbol{\theta}_j) f_{\mathbf{a},\mathbf{r}}(\boldsymbol{\theta}_{\mathbf{a}}) \prod_{i \in S(\mathbf{a}) \setminus \{j\}} \exp(\mathbf{v}(\boldsymbol{\theta}_i) \cdot \boldsymbol{\lambda}_i) d\boldsymbol{\theta}_{\mathbf{a}}}{\int_{\mathcal{Q}_{\mathbf{a}}} f_{\mathbf{a},\mathbf{r}}(\boldsymbol{\theta}_{\mathbf{a}}) \prod_{i \in S(\mathbf{a}) \setminus \{j\}} \exp(\mathbf{v}(\boldsymbol{\theta}_i) \cdot \boldsymbol{\lambda}_i) d\boldsymbol{\theta}_{\mathbf{a}}} \right)$$

where λ_i , $i \in \mathcal{S}(a) \setminus \{j\}$ are i.i.d. samples from the probability distribution g_λ .

With these definitions we are now ready to prove that density evolution is a viable way of assessing the performance and convergence of EP under the As. 3-9 in an asymptotic sense.

Thm. 1 (Density Evolution for EP): Let the number of parameters and observations grow in such a way that As. 3-9 are all satisfied for each K as $K \rightarrow \infty$. Then, with probability $\rightarrow 1$ as $K \rightarrow \infty$, at any given iteration k , given $\theta = \theta_0 \in \mathcal{R}$, the incoming messages $\{n_{a \rightarrow j} | a \in \mathcal{F}(j)\}$ to any basis node v_j are independent and identically distributed samples from a probability distribution q_k obeying the recursion

$$q_{k+1} = \left(\mathfrak{F}_{\theta_0} \circ \gamma \otimes \circ \bigotimes_{d_p-1} \right) q_k \quad (5)$$

where \mathfrak{F}_{θ_0} can be chosen as any $\mathfrak{F}_{a \rightarrow j, \theta_0}$.

Proof: It was shown in [10] that for any fixed k , the assumption 9 guarantees that the computation neighborhood $C_k(v_j)$ is a tree with probability $\rightarrow 1$ as $K \rightarrow \infty$. Only those factor nodes which are in the computation tree $C_{2k}(v_j)$ affect those messages passed to v_j during the k th iteration. Furthermore, As. 5 and 7 require the observations r_a included in this computation tree to be independently and identically distributed given $\theta = \theta_0 \in \mathcal{R}$. This gives that the message passed at the top edges of the computation tree (to v_j) are independent. As. 3, 4, and 8 which guarantee that all of the factors and basis functions, and thus the message creating functions at the factor nodes, are the same, then imply that not only are these incoming messages to any basis node v_j independent, but they are also identically distributed. Thus, the probability density for the message at the output of the factor node can be calculated by taking the probability distribution corresponding to $d_f - 1$ i.i.d. samples from the input message distribution, and putting it through the factor message update (3). Solving (3) is then equivalent to calculating $H_{a \rightarrow j}(r, \lambda)$ where λ are the samples from the incoming message probability distribution.

At the basis node, the message passing rules (4) require that all but one of the incoming messages be summed to create each output message. Since we have just proved that these messages are independent and identically distributed, this operation creates an outgoing (to the factor nodes) message with probability density that is the convolution of the incoming distribution with itself $d_p - 1$ times and then with the degree one factor node message density. \square

Cor. 1 (Parameter Distribution Estimate): Given $\theta = \theta_0 \in \mathcal{R}$, the estimate $\hat{p}_{\theta_j | r}^k$ provided by EP for the a posteriori distribution $p_{\theta_j | r}$ after k iterations is given by

$$\hat{p}_{\theta_j | r}^j = \exp(v(\theta_j) \cdot \lambda - \psi_v(\lambda))$$

where λ is distributed according to the probability density function h_j defined by $h_k := \bigotimes_{d_p} q_k$ and q_k is defined by the recursion (5).

Theorem 1 and Corollary 1 are immensely useful, because they show that determining the performance and convergence of EP under the homogeneity As. 3-9 is equivalent to studying the convergence properties of the iterated map (5). This is, in essence, the theory of density evolution.

4. CONCLUSIONS

This paper showed that density evolution style analysis can be applied to general expectation propagation under some homogeneity and symmetry assumptions. Future work will then apply density evolution to determine the performance of particular instances of expectation propagation, for example ones used for distributed estimation in sensor networks.

5. REFERENCES

- [1] T. P. Minka, *A Family of Algorithms for Approximate Bayesian Inference*, Ph.D. thesis, Massachusetts Institute of Technology, 2001.
- [2] T. P. Minka, "Expectation propagation for approximate Bayesian inference," in *Uncertainty in AI'01*, 2001.
- [3] J. M. Walsh and P. A. Regalia, "Iterative constrained maximum likelihood estimation via expectation propagation," in *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, Toulouse, France, May 2006.
- [4] J. M. Walsh, "Dual Optimality Frameworks for Expectation Propagation," in *IEEE Conference on Signal Processing Advances in Wireless Communications (SPAWC)*, Cannes, France, June 2006.
- [5] J. Pearl, *Probabilistic reasoning in intelligent systems : networks of plausible inference*, Morgan Kaufmann Publishers, 1988.
- [6] 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, Feb. 2001.
- [7] Lawrence D. Brown, *Fundamentals of Statistical Exponential Families with Applications in Statistical Decision Theory*, Institute of Mathematical Statistics, 1986.
- [8] S. Amari, *Methods of Information Geometry*, vol. 191, AMS Translations of Mathematical Monographs, 2004.
- [9] M. J. Wainwright and M. I. Jordan, "Graphical models, exponential families, and variational inference," Tech. Rep., Department of Statistics, University of California, Berkeley.
- [10] T. J. Richardson and R. L. Urbanke, "The capacity of low-density parity-check codes under message-passing decoding," *IEEE Trans. Inform. Theory*, vol. 47, no. 2, pp. 599–618, Feb. 2001.
- [11] S. ten Brink, "Convergence behavior of iteratively decoded parallel concatenated codes.," *IEEE Trans. Commun.*, vol. 49, pp. 1727–1737, Oct. 2001.
- [12] H. El Gamal and A. R. Hammons, Jr., "Analyzing the turbo decoder using the gaussian approximation," *IEEE Trans. Inform. Theory*, vol. 47, pp. 671–686, Feb. 2001.
- [13] C. Berrou, A. Glavieux, and P. Thitimajshima, "Near shannon limit error-correcting coding and decoding: Turbo-codes.," in *ICC 93*, Geneva, May 1993, vol. 2, pp. 1064–1070.
- [14] R. G. Gallager, "Low-density parity-check codes.," *IRE Trans. Information Theory*, vol. 2, pp. 21–28, 1962.
- [15] D. Divsalar, S. Dolinar, and F. Pollara, "Iterative turbo decoder analysis based on density evolution.," *IEEE J. Select. Areas Commun.*, vol. 19, pp. 891–907, May 2001.