

2. Expectation Propagation Algorithm and MU-MIMO Systems

As mentioned in chapter 1, we believe that EP is a low complexity algorithm that is suitable for MIMO data detection. The reason is EP solves the complexity problem caused by optimal data detectors such as belief propagation (BP) or also known as message passing algorithm (MPA). Although there is a performance loss in EP approximation, however as the number of antennas grows, the EP performance can improve significantly. At the same time, due to the exponential increase in complexity in the optimal data detector, it soon becomes prohibitive.

The EP approximates the marginal distribution of the posterior probability by using an exponential family. Thus, the complexity of EP is much lower than the optimal detector algorithm. To support our argument, later in chapter 3 and 4, we provide two applications of EP in massive MU-MIMO communication systems, including complexity analysis, performance evaluation, and theoretical analysis.

In this chapter, we would like to present the details of the single loop EP algorithm [3] and its complexity analysis. Furthermore, we investigate future works for EP, either to improve its performance (double loop EP algorithm) or reduce its complexity (approximation of inverse matrix value). We also briefly discuss about state evolution (SE) which is used to perform a theoretical analysis. Lastly, regarding to the need of implementing massive MU-MIMO for the near future fifth generation wireless system (5G), we briefly introduce the MU-MIMO system model.

2.1 Single Loop Expectation Propagation

Before we discuss about the single loop expectation propagation algorithm, from now on, we refer to the single loop EP algorithm as an EP algorithm. The idea of EP is

to construct a tractable approximation $p(\mathbf{x})$ by a distribution $q(\mathbf{x})$. Given some statistical model with latent variables $\mathbf{x} \in \Omega^\delta$, it can be factorized in the following way

$$p(\mathbf{x}) \sim f(\mathbf{x}) \prod_{i=1}^I t_i(\mathbf{x}) \quad (2.1)$$

where, $f(\mathbf{x})$ belongs to an exponential family \mathcal{F} with sufficient statistics $\Phi(\mathbf{x}) = [\phi_1(\mathbf{x}), \phi_2(\mathbf{x}), \dots, \phi_S(\mathbf{x})]$, and $t_i(\mathbf{x})$, $\forall i \in I$, are non-negative factors. Furthermore, assume \mathcal{F} is the multivariate Gaussian family, $\Phi(\mathbf{x}) = \{x_i, x_i x_j\}_{i,j=1}^\delta$.

The general purpose framework [12] in order to approximate $p(\mathbf{x})$ is achieved when the moments between $p(\mathbf{x})$ and $q(\mathbf{x})$ are equal. This is known as the moment matching (MM) condition, where

$$\mathbb{E}_{q(\mathbf{x})}[\phi_j(\mathbf{x})] = \mathbb{E}_{p(\mathbf{x})}[\phi_j(\mathbf{x})], \forall j \in S. \quad (2.2)$$

When the $p(\mathbf{x})$ and $q(\mathbf{x})$ are defined in the same space and measure, the moment matching condition is equal to finding $q(x)$ that satisfies minimum Kullback-Leibler divergence value.

$$q(\mathbf{x}) = \arg \min_{q'(\mathbf{x}) \in \mathcal{F}} D_{KL}(p(\mathbf{x}) || q'(\mathbf{x})), \quad (2.3)$$

Equation (2.2) can be solved by doing the sequential EP algorithm, iteratively. First, replace each one of $t_i(\mathbf{x})$ factors in (2.1) by a member $\tilde{t}_i(\mathbf{x})$, where $\tilde{t}_i \in \mathcal{F}, \forall i \in I$. We can assume $q(\mathbf{x})$ as an approximation of $p(\mathbf{x})$, hence (2.1) can be rewritten as

$$q(\mathbf{x}) \sim f(\mathbf{x}) \prod_{i=1}^I \tilde{t}_i(\mathbf{x}). \quad (2.4)$$

Define $q^{(0)}(\mathbf{x})$ as initial value of $q(\mathbf{x})$ and $q^{(l)}(\mathbf{x})$ as a $q(\mathbf{x})$ at iteration l , $q(\mathbf{x})$ can be obtained by updating each one of the $\tilde{t}_i(\mathbf{x})$ factors independently. In [3] the iteration of EP algorithm is

- 1) Calculate the cavity distribution

$$q^{(l)\setminus i}(\mathbf{x}) = \frac{q^{(l)}(\mathbf{x})}{\tilde{t}_i(\mathbf{x})} \in \mathcal{F}. \quad (2.5)$$

2) Compute the distribution $\hat{p}_i(\mathbf{x}) \sim t_i(\mathbf{x})q^{(l)\setminus i}(\mathbf{x})$, where $\hat{p}_i(\mathbf{x})$ is approximation value of $p_i(\mathbf{x})$, then compute

$$\mathbb{E}_{\hat{p}_i(\mathbf{x})}[\phi_j(\mathbf{x})], \forall j \in S \quad (2.6)$$

3) The updated factor $\tilde{t}_i^{new}(\mathbf{x})$ can be obtained by

$$\mathbb{E}_{\tilde{t}_i^{new}q^{(l)\setminus i}(\mathbf{x})}[\phi_j(\mathbf{x})] = \mathbb{E}_{\hat{p}_i(\mathbf{x})}[\phi_j(\mathbf{x})], \forall j \in S \quad (2.7)$$

The iteration of EP algorithm will be done if either convergence criterion is met or the maximum number of iteration is reached.

2.1.1 EP Message Passing

According to the principle of expectation propagation [13], the whole EP algorithm can be interpreted as a message passing from estimation module (module A) to demodulation module (module B). The rule of the message passing and the modules of the EP is illustrated in Figure 2.1 and Figure 2.2.

Given that, K is the total user number, N is the total receiver number, $\mathbf{y} = [y_1, y_2, \dots, y_N]$ is the received signal, and $\mathbf{H} = [\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_N]$ denotes the channel vector and its channel response from the k -th transmit antenna to the BS at the n -th receiver antennas. Furthermore, define $a = n$ -th receiver at the BS, the messages from estimation module ($m_{i \rightarrow a}$) and the demodulation module ($m_{a \rightarrow i}$) are given as follow:

$$m_{i \rightarrow a}^{l+1}(x_i) \propto \frac{Proj[P_x(x_i) \prod_b m_{b \rightarrow i}^l(x_i)]}{m_{a \rightarrow i}^l(x_i)} \quad (2.8)$$

$$m_{a \rightarrow i}^l(x_i) \propto \frac{Proj[m_{i \rightarrow a}^l(x_i) \times \int \prod_{j=1, j \neq i}^K m_{j \rightarrow a}^l(x_j) P(y_a | x_j)]}{m_{i \rightarrow a}^l(x_i)}, \quad (2.9)$$

where,

$$P(y_a | x_j) = \frac{1}{\pi \sigma^2} e^{-\frac{|y_a - \sum_{a=1}^K h_a^{\mathbf{H}} x_j|^2}{\sigma^2}} \quad (2.10)$$

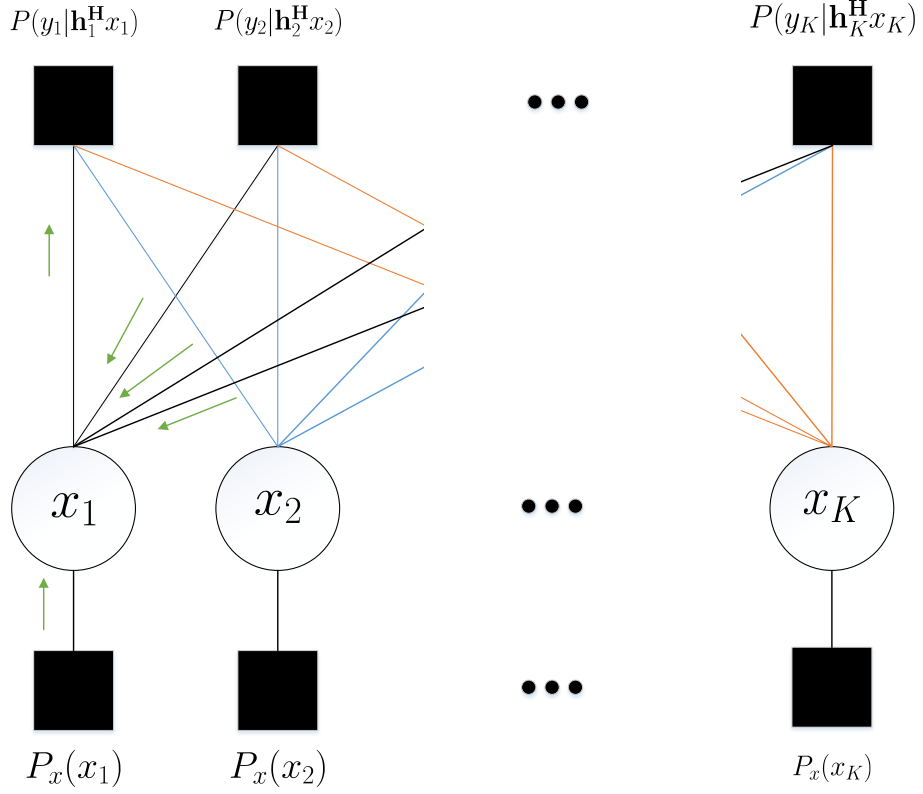


Figure 2.1. EP message passing rule

and $\prod_b m_{b \rightarrow i}^l(x_i) \propto \mathcal{N}_c(x_i^l, \mu_i^l; \Sigma_i^l)$. After projection, $m_{i \rightarrow a}^l(x_k)$ is approximated as

$$\mathcal{N}_c(x_k, \hat{x}_k^l, v_k^l). \quad (2.11)$$

The proven is given as follow. Noting that $\mathbf{y} = \mathbf{h}_i x_i + \sum_{j=1, j \neq i}^I \mathbf{h}_j x_j + \boldsymbol{\eta}$, where $\boldsymbol{\eta} \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$ denotes the Additive White Gaussian Noise (AWGN). Then, we can obtain

$$\int \prod_{j=1, j \neq i}^I m_{j \rightarrow a}^l(x_k) P(y_a | \mathbf{x}) \propto \mathcal{N}_c \left(x_k; \frac{y_a - \sum_{j=1, j \neq k}^K h_{a,j} \hat{x}_j^l}{h_{a,k}}, \frac{\sigma^2 + \sum_{j=1, j \neq k}^K |h_{a,j}|^2 v_j^l}{|h_{a,k}|^2} \right).$$

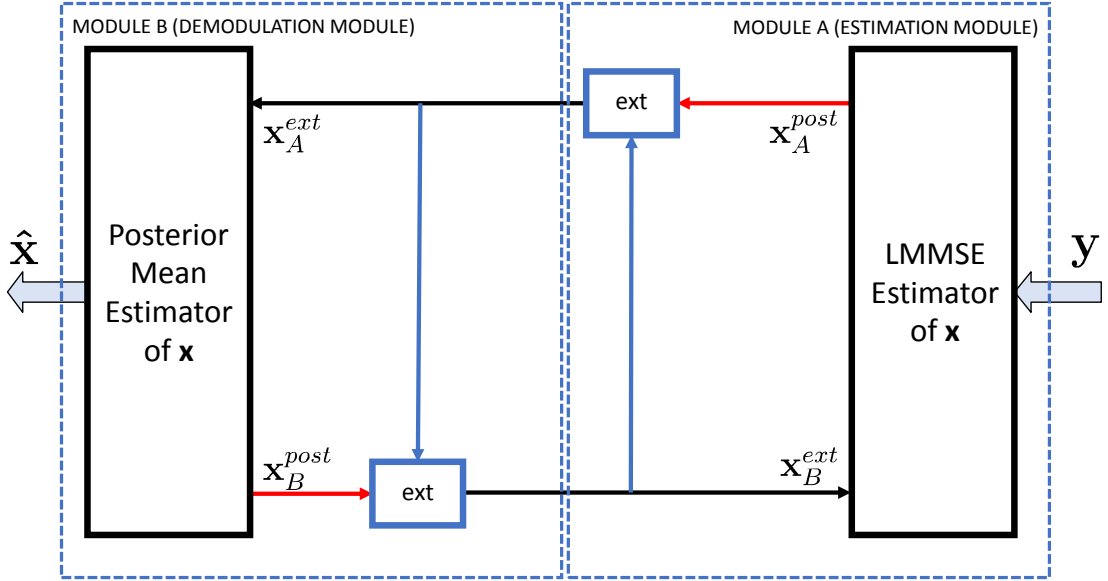


Figure 2.2. Block diagram of EP

Finally, (2.11) can be Rewritten as

$$m_{a \rightarrow i}^t(x_i) \propto \mathcal{N}_c \left(x_k; \frac{y_a - \sum_{j=1, j \neq k}^K h_{a,j} \hat{x}_j^l}{h_{a,k}}, \frac{\sigma^2 + \sum_{j=1, j \neq k}^K |h_{a,j}|^2 v_j^l}{|h_{a,k}|^2} \right). \quad (2.12)$$

2.1.2 Detail of EP Algorithm

EP algorithm is started by reconstructing equation (2.4). The prior input distribution (\tilde{t}_i) is replaced by an independent Gaussian distribution, such that

$$q(\mathbf{x}) \propto \mathcal{N}(\mathbf{y} : \mathbf{H}\mathbf{x}, \sigma^2 \mathbf{I}) \prod_{i=1}^K e^{x_i^{\mathbf{H}} \gamma_i + \gamma_i^{\mathbf{H}} x_i - \lambda_i |x_i|^2} \quad (2.13)$$

where, $\gamma \in \mathbb{R}^K$, and $\lambda \in \mathbb{R}^K, \forall i \in K$. Equation (2.13) fulfills the MMSE approximation to the posterior distribution $p(\mathbf{x}|\mathbf{y})$ as presented in [14], [15]. Performing

Gaussian product Lemma [16] on (2.13), we can define $q(\mathbf{x})$ by its Gaussian mean vector $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$ as given in (2.15a) and (2.15b). Before passing the messages to the next module, the prior information has to be removed from the estimation value as the process is called the calculation of cavity distribution. Calculating the cavity distribution from estimation module is the same as finding the extrinsic values of $\mathbf{x}_A^{\text{ext}}$ and $\mathbf{v}_A^{\text{ext}}$, which are given in (2.16a) and (2.16b), respectively.

In the demodulation module, the expectation and variance of the posterior estimator ($\mathbf{x}_B^{\text{post}}, \mathbf{v}_B^{\text{post}}$) are computed by calculating conditional expectation from the extrinsic information of \mathbf{x}_A . Considering that M is the cardinality of transmitted symbols, for each k -th user, the expectations in (2.17a) and (2.17b) are with respect to $P(x_k|x_{A,k})$, which can be obtained by the Bayes rule

$$P(x_k|x_{A,k}) = \frac{P(x_{A,k}|x_m)P_x(x_m)}{P(x_{A,k})}, \quad (2.14)$$

where

$$P(x_{A,k}|x_m)P_x(x_m) = \frac{1}{M} \frac{1}{\pi v_{A,k}} \exp\left(-\frac{|x_{A,k} - x_m|^2}{v_{A,k}}\right),$$

$$P(x_{A,k}) = \frac{1}{M} \frac{1}{\pi v_{A,k}} \sum_{m=1}^M \exp\left(-\frac{|x_{A,k} - x_m|^2}{v_{A,k}}\right).$$

The results of (2.17a) and (2.17b) are obviously identical to the (2.12), as it can be considered as the posterior message from demodulation module.

The extrinsic values of the demodulation module, i.e., $\mathbf{v}_B^{\text{ext}}$ and $\mathbf{x}_B^{\text{ext}}$, are calculated in (2.18a) and (2.18b), respectively. Noting that, the value of $\mathbf{v}_B^{\text{ext}}$ may return a negative. In this case, we simply use the previous value of $\mathbf{v}_B^{\text{ext}}$ and $\mathbf{x}_A^{\text{ext}}$ as a new pair of updating parameters. After the iteration converges, the conditional expectation given by (2.17a) is expected to be the estimated signals. The complete algorithm is shown in Algorithm 1.

2.1.3 EP Computational Complexity

EP computational complexity becomes the key issue to determine the success of the EP. Specifically, EP is expected to solve the complexity problem of optimal

Initialization: $\gamma_{B \rightarrow A}^0 = 0, \lambda_{B \rightarrow A}^0 = \frac{1}{E_s} \mathbf{I}, d(\mathbf{Q}) = \text{diag}(\mathbf{Q});$

for $l = 1 : L_{max}$ **do**

Estimation Module:

(1) Compute the a posteriori mean/variance of \mathbf{x}_A :

$$\mathbf{v}_{A,l}^{\text{post}} = \Sigma^l = \left(\sigma^{-2} \mathbf{H}^H \mathbf{H} + d(\lambda_{B \rightarrow A}^{l-1}) \right)^{-1} \quad (2.15a)$$

$$\mathbf{x}_{A,l}^{\text{post}} = \boldsymbol{\mu}^l = \Sigma^l \left(\sigma^{-2} \mathbf{H}^H \mathbf{y} + \gamma_{B \rightarrow A}^{l-1} \right) \quad (2.15b)$$

(2) Compute the extrinsic mean/variance of \mathbf{x}_A :

$$\mathbf{v}_{A,l}^{\text{ext}} = \left(\frac{1}{d(\Sigma^l)} - d(\lambda_{B \rightarrow A}^{l-1}) \right)^{-1} \quad (2.16a)$$

$$\mathbf{x}_{A,l}^{\text{ext}} = d(\mathbf{v}_{A \rightarrow B}^l) \left(\frac{\boldsymbol{\mu}^l}{d(\Sigma^l)} - \gamma_{B \rightarrow A}^{l-1} \right) \quad (2.16b)$$

Demodulation Module:

(3) Compute the a posteriori mean/variance of \mathbf{x}_B :

$$\mathbf{x}_{B,l}^{\text{post}} \leftarrow \mathbb{E}\{\mathbf{x} | \mathbf{x}_{A,l}^{\text{ext}}, \mathbf{v}_{A,l}^{\text{ext}}\} \quad (2.17a)$$

$$\mathbf{v}_{B,l}^{\text{post}} \leftarrow \text{Var}\{\mathbf{x} | \mathbf{x}_{A,l}^{\text{ext}}, \mathbf{v}_{A,l}^{\text{ext}}\} \quad (2.17b)$$

(4) Compute the extrinsic mean/variance of \mathbf{x}_B :

$$\mathbf{v}_{B,l}^{\text{ext}} = \lambda_{B \rightarrow A}^l = \left(\frac{1}{\mathbf{v}_{B,l}^{\text{post}}} - \frac{1}{\mathbf{v}_{A,l}^{\text{ext}}} \right)^{-1} \quad (2.18a)$$

$$\mathbf{x}_{B,l}^{\text{ext}} = \gamma_{B \rightarrow A}^l = \left(\frac{\mathbf{x}_{B,l}^{\text{post}}}{\mathbf{v}_{B,l}^{\text{post}}} - \frac{\mathbf{x}_{A,l}^{\text{ext}}}{\mathbf{v}_{A,l}^{\text{ext}}} \right)^{-1} \quad (2.18b)$$

end

Algorithm 1: EP Algorithm

detector algorithm such as belief propagation (BP) [17] as well as maintaining small performance loss. In this section, we provide a comparison regarding to the complexity of BP as a optimal detector algorithm and EP as a low complexity detector. First, we briefly explain the complexity of BP as a optimal detector algorithm. Next, we

compare the EP and BP complexity under the same settings. Finally, we discuss about matrix inversion lemma that can reduce the direct calculation of EP complexity.

EP and BP Complexity

Given that a system model which has transmitter antennas $N_t = 8$, receiver antennas $N_r = 8$, and employs 256 QAM modulation. The complexity of BP is $\mathcal{O}(|\mathcal{A}|^{N_t})$, where \mathcal{A} denotes the cardinality of the QAM modulation. Therefore, the complexity of BP algorithm in the system model given above is $\mathcal{O}(256^8)$. Furthermore, if we increase the transceiver antennas, now become $N_t = 16$, $N_r = 16$, the complexity of BP will be $\mathcal{O}(256^{16})$. Now, we prove that BP complexity is incremental exponentially and soon becomes prohibitive.

EP complexity is dominated by the calculation of inverse matrix in (2.15a). The higher the number of transceiver antennas, the higher the dimensional of the EP inverse matrix. Given the number of iteration ($l = 10$), the complexity of EP is $\mathcal{O}(lN_rN_t^2)$. Under the same system model mentioned above, the complexity of EP is $\mathcal{O}(5120)$, $10^{-14}\%$ of the BP complexity. Furthermore if we continue increase the transceiver antennas, now become $N_t = 16$, $N_r = 16$, the complexity of EP will be $\mathcal{O}(40960)$, still much lower than BP complexity. It is clearly that EP complexity is increase linearly with the dimension of the system instead of exponentially, therefore EP solves the complexity problem of BP.

EP Matrix Inversion Lemma

As mentioned before, EP complexity lies in the inverse matrix Σ . The direct calculation of the inverse matrix will result $\mathcal{O}(lN_t^3)$. If N_r is less than N_t , a better approach to reduce the complexity of EP is by using matrix inversion lemma. Define $\mathbf{P} = d(\boldsymbol{\lambda}_{B \rightarrow A}^{l-1})$, the matrix inversion lemma is given by

$$(\sigma^{-2}\mathbf{H}^H\mathbf{H} + \mathbf{P})^{-1} = \mathbf{P}^{-1} - \sigma^{-2}\mathbf{P}^{-1}\mathbf{H}^H(\mathbf{I} + \sigma^{-2}\mathbf{H}\mathbf{P}^{-1}\mathbf{H}^H)^{-1}\mathbf{H}\mathbf{P}^{-1} \quad (2.19)$$

Now, the complexity of EP is $\mathcal{O}(lN_r^2N_t)$.

2.2 Future Work For Expectation Propagation

In this section, we investigate the future works for EP, either to improve its performance or reduce its complexity. We introduce a double loop algorithm that proposed in [18] as a way to improve the EP performance. The other way to improve the EPA performance is using a combination of Gaussian message and EP message as proposed in [19]. For reducing the complexity of EP algorithm, as proposed in [20] - [21], the approximation of the inverse matrix value in (2.15a) can be used instead of compute the inverse matrix itself.

2.2.1 Double Loop EP

As described in [22], [18], the original EP has a bad fixed point issue. Bad fixed point means after EP approximation value reaches a fixed point, this fixed point can not achieve EP's truly convergence value. This issue causes a performance degradation in EP. To figure out the bad fixed point issue, [18] proposes a double loop EP algorithm.

The idea of double loop algorithm is to perform a jointly solving of the extrinsic computation and expectations calculation. Thus, the converge point can be guaranteed. The double loop algorithm can be divided into two parts which are outer and inner loop. The outer loop is identical to the one for the single loop. The inner loop refers to perform a numerical method iteration which solves the jointly equation of extrinsic and expectations in (2.16a, 2.16b) and (2.17a, 2.17b). According to the Algorithm 1, the jointly equation can be written as

$$\gamma^l + \frac{\mathbf{x}_{B,l}^{post}}{\mathbf{v}_{B,l}^{post}} = \gamma^{l-1} + \frac{\mathbf{x}_{A,l}^{post}}{\mathbf{v}_{A,l}^{post}} \quad (2.20)$$

$$\lambda^l + \frac{1}{\mathbf{v}_{B,l}^{post}} = \lambda^{l-1} + \frac{1}{\mathbf{v}_{A,l}^{post}} \quad (2.21)$$

collecting variable terms and constant terms on the left and right hand side, respectively. The equation 2.20 and 2.21 can be solved iteratively by a numerical method solver such as Newton method.

However, the complexity of the double loop EP algorithm is increased significantly and soon becomes prohibitive as the number of QAM modulation grows large. The reason is the computational burden in the numerical solver. The solver needs an exhaustive iteration in order to find the solution of (2.20). In a large number of constellation, equation (2.20) will be very complex, thus the numerical solver might not efficiently to be used.

Considering the motivation of proposing the EP algorithm i.e. to solve the complexity problem of optimal detector, double loop implementation which causes a high complexity of EP is denying the purpose of the EP itself. We conclude that the double loop EP algorithm is difficult to be implemented. However, in the future works, EP had fixed point still be a promising work that needs to be solved. So, the EP performance can be significantly improved.

2.2.2 Approximation of Inverse matrix

In this section, we briefly describe a way to reduce the EP complexity. Even though EP is known as a low computational complexity detector algorithm, however its complexity still too high especially for a large scale system.

As mentioned above, the complexity of EP lies on the inverse matrix in (2.15a). The way to reduce the EP complexity i.e. approximating the value of the inverse matrix without compute the inverse itself. In [23], conjugate gradient method is proposed to directly find the value of (2.15b) without calculating the inverse matrix in (2.15a). However, (2.16a) also needs the diagonal term of Σ which is now become unknown. As a future work, we note two promising candidates to approximate the value of diagonal term inverse matrix those are fast algorithm to extract diagonal

inverse matrix [20] and probing method for computing the diagonal of a matrix inverse [24]

Furthermore, in case of sparse matrix, the promising candidates are LDU factorization as proposed in [25] and estimation of the diagonal elements of a sparse precision matrix [21].

2.3 EP State Evolution

We then employ the performance analysis framework in [26] to develop the state evolution (SE) of the EPA shown in Algorithm 1. The performance analysis framework is derived from a large scale system. Considering a large scale system, $\mathbf{v}_B^{\text{ext}}$ and $\mathbf{v}_A^{\text{ext}}$ in (2.18a) and (2.16a) can be approximated by their average values v_B and v_A , respectively. Following this assumption, the inputoutput transfer function of estimation module can be derived by substituting (2.15a) and (2.18a) into (2.16a), yielding

$$v_A = \left(K^{-1} \text{tr}(\sigma^{-2} \mathbf{H}^H \mathbf{H} + v_B^{-1} \mathbf{I})^{-1} \right)^{-1} - v_B^{-1}. \quad (2.22)$$

Now, v_A represents the input variance of the demodulation module and can be regarded as the SNR of the equivalent scalar additive white gaussian noise (AWGN) channel

$$y = x + v_A \eta. \quad (2.23)$$

This equivalent scalar AWGN channel can be considered as a k -th channel under K users. Consistent with our assumption, where v_B and v_A are the average values $\mathbf{v}_B^{\text{ext}}$ and $\mathbf{v}_A^{\text{ext}}$; each k -th channel model on (2.17a) will have an identical value. To simplify our explanation on SE, we only adopt k -th channel model (2.23), as an identical parallel channel model in (2.17a). Similarly, we define v as the scalar version of (2.17b). Therefore, v can be calculated by

$$v = \text{Var}\{x|x_A, v_A\} = \text{E}\{|x - \text{E}\{x|x_A, v_A\}|^2\}, \quad (2.24)$$

where the expectation is with respect to $P(x|x_A)$ given by (2.14). Referring to (2.18a), v_B can be defined as

$$v_B = (v^{-1} - v_A)^{-1}. \quad (2.25)$$

The iteration of the EP algorithm is identical to the SE in (2.22) and (2.25). The SE is self-consistent. That is, the iteration of estimation and demodulation module can be traced from (2.22) and (2.25) without iterating the entire algorithm. Let v^* denote the iteration as converging. Regarding the scalar AWGN channel (2.23), the theoretical BER as a function of v_A can be calculated using the Q function. The SE of (2.22) and (2.25) is identical to that in [27] whose fixed points have MSE consistent with the MMSE from [27].

2.4 Massive MU-MIMO Systems

Massive MU-MIMO system is believed to be a key technology for next generation of wireless system. Hundreds or even thousands antennas will be employed in order to fulfill the requirements of 5G technology such as massive connectivity, better quality of service, higher throughput, lower latency, and lower control signaling overhead. From [28], [29], [30], [31], it can be concluded that MU-MIMO systems has some critical advantages i.e. 1) Allows a direct gain in multiple access capacity 2) Line of sight propagation is no longer a problem 3) Increase in spectral efficiency 4) Near optimal simple coherent linear processing techniques

In this thesis, we focus on considering uplink scheme data transceiver. Therefore, the system model of MU-MIMO can be described in 2.3. Suppose base station employs N number of receiver antennas, U denotes the number of users, and each user employ N_t number of transmitter antennas. Define $K = UN_t$, the received signals can be construct as

$$\mathbf{y} = \mathbf{H}\mathbf{x} + \boldsymbol{\eta} \quad (2.26)$$

where, \mathbf{y} is $N \times 1$ received signals, \mathbf{H} is $N \times K$ channel gain, and \mathbf{x} is $K \times 1$ transmitted signals.

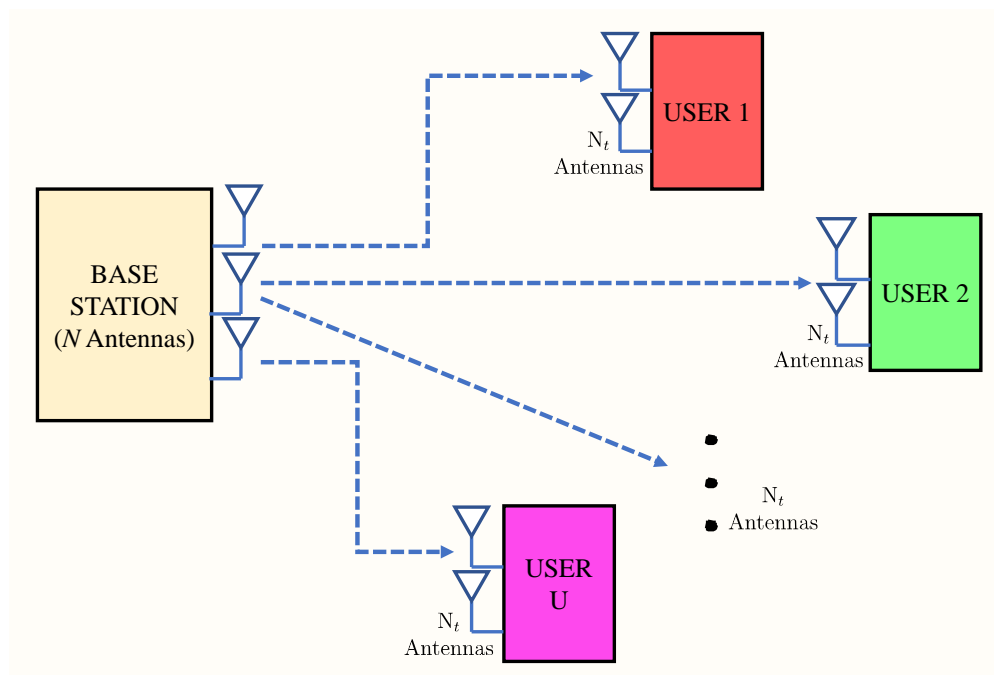


Figure 2.3. Block diagram of EP