

On the Matrix Inversion Approximation Based on Neumann Series in Massive MIMO Systems

Dengkui Zhu*, Boyu Li*, and Ping Liang*[†]

*RF DSP Inc., 30 Corporate Park, Suite 210, Irvine, CA 92606, USA,

e-mail: dkzhu@rfdsp.com, byli@rfdsp.com, pliang@rfdsp.com [†]Department of Electrical

Engineering, University of California - Riverside, Riverside, CA 92521, USA,

e-mail: liang@ee.ucr.edu

Abstract

Zero-Forcing (ZF) has been considered as one of the potential practical precoding and detection method for massive MIMO systems. One of the most important advantages of massive MIMO is the capability of supporting a large number of users in the same time-frequency resource, which requires much larger dimensions of matrix inversion for ZF than conventional multi-user MIMO systems. In this case, Neumann Series (NS) has been considered for the Matrix Inversion Approximation (MIA), because of its suitability for massive MIMO systems and its advantages in hardware implementation. The performance-complexity trade-off and the hardware implementation of NS-based MIA in massive MIMO systems have been discussed. In this paper, we analyze the effects of the ratio of the number of massive MIMO antennas to the number of users on the performance of NS-based MIA. In addition, we derive the approximation error estimation formulas for different practical numbers of terms of NS-based MIA. These results could offer useful guidelines for practical massive MIMO systems.

I. INTRODUCTION

Massive Multiple-Input Multiple-Output (MIMO) systems were firstly introduced in [1], and have drawn great interest from both academia and industry. In such systems, each Base Station (BS) is equipped with dozens to hundreds of antennas to serve tens of users in the same time-frequency resource. Therefore, they can achieve significantly higher spatial multiplexing gains than conventional multi-user MIMO systems, which offers one of the most important advantages of massive MIMO systems, the potential capability to offer linear capacity growth without increasing power or bandwidth [1]–[4].

It has been shown that, for massive MIMO systems where the number of antennas M , e.g., $M = 128$, is much larger than the number of served users K , e.g., $K = 16$, [2], [4], Zero-Forcing (ZF) precoding and detection can achieve performance very close to the channel capacity for the downlink and uplink respectively [2]. As a result, ZF has been considered as one of the potential practical precoding and detection method for massive MIMO systems [2], [4]–[6].

For the hardware implementation of ZF, despite of the very large number of M , the main complexity is the inverse of a $K \times K$ matrix [2], [7], [8]. Unfortunately, for massive MIMO systems, although K is much smaller than M , it is still much larger than conventional multi-user MIMO systems. As a result, in this case, the computation of the exact inversion of the $K \times K$ matrix could result in very high complexity [8], which may cause large processing delay so that the demands of the channel coherence time is not met. Due to this reason, Neumann Series (NS) has been considered to carry out the Matrix Inversion Approximation (MIA), because it is well suited for massive MIMO systems and it is advantageous for hardware implementation [2], [7], [8].

Despite of the advantages, some potential application issues of the NS-based MIA have also been identified. Firstly, for a finite M/K ratio, the NS may not converge, resulting the failure of the algorithm [2], [7]. What M/K ratio could offer high convergence probability is still not clear. Secondly, for the NS-based MIA to achieve good performance with quick convergence, the $K \times K$ matrix needs to be diagonally dominant [8]. In order to satisfy this condition, $M \gg K$ is required [2], [8]. Similarly, what M/K ratio could provide high probability of diagonally dominant is also not clear. Moreover, with a larger number of terms, the NS-based MIA offers closer performance to the exact inversion [7], [8]. However, the larger number of terms results in more processing cycles. Hence, for practical hardware implementation, the number of terms cannot be very large. Although the approximation error analysis was carried out and a residual error upper bound of the NS-based MIA was derived [8], the approximation error analysis with high accuracy has not been derived.

In this paper, we address the three problems listed above. Specifically, we firstly derived a M/K ratio condition that offers high convergence probability. Then, we derived another M/K ratio condition that provides high probability for the $K \times K$ matrix to be diagonally dominant. Finally, we carry out the approximation error analysis with high accuracy for practical numbers of terms for the NS-based MIA in hardware implementation.

The remainder of this paper is organized as follows. In Section II, the basis of the NS-based MIA in

massive MIMO systems is briefly reviewed. The M/K ratio condition that provides high convergence probability is derived in Section III. Then, another M/K ratio condition that offers high diagonally dominant probability for the $K \times K$ matrix is derived in Section IV. In Section V, the approximation error analysis with high accuracy for practical numbers of terms for the NS-based MIA is carried out. Finally, after a discussion in Section VI, conclusions are drawn in Section VII.

II. BASIS OF NS-BASED MIA IN MASSIVE MIMO SYSTEMS

Consider a massive MIMO wireless system where the BS is equipped with M antennas to serve K single-antenna users in the same time-frequency resource. Then, for the uplink, the $M \times K$ channel matrix is represented by $\mathbf{H} = [h_{mk}]$, where h_{mk} denotes the channel coefficient between the m th antenna and the k th user, with $m = 1, \dots, M$, and $k = 1, \dots, K$. Similarly to [2], [7], [8], the analysis in this paper assumes that the h_{mk} elements are in uncorrelated Rayleigh flat fading, i.e., independent and identically distributed (i.i.d.) zero-mean unit-variance complex Gaussian variables. Note that, for the Time-Division Duplexing (TDD) mode, due to the channel reciprocity, the downlink has the same channel matrix \mathbf{H} as the uplink, as long as the transmission duration is within the channel coherence time [1]–[6].

In order to carry out ZF precoding for the downlink or the ZF detection for the uplink, the pseudo-inverse of \mathbf{H} needs to be calculated [2], [4]–[6], which is written as

$$\mathbf{H}^\dagger = (\mathbf{H}^H \mathbf{H})^{-1} \mathbf{H}^H. \quad (1)$$

Let $\mathbf{G} = \mathbf{H}^H \mathbf{H}$. Then, in (1), despite of the very large number of M , e.g., 256, in massive MIMO systems, the main complexity of the hardware implementation lies in the inversion of the $K \times K$ matrix \mathbf{G} [2], [7], [8]. To exploit the large spatial multiplexing gains of massive MIMO systems, although much smaller than M , the number of K is much larger than conventional multi-user MIMO systems, e.g., $K = 16$. As a result, the complexity of calculating \mathbf{G}^{-1} may be too high for hardware implementation. To address this issue, NS has been considered to carry out the MIA, because it is advantageous in hardware implementation and it is suitable for massive MIMO systems [2], [7], [8]. Specifically, it can be written as

$$\mathbf{G}_N^{-1} \approx \sum_{n=0}^{N-1} (\mathbf{I}_K - \Theta \mathbf{G})^n \Theta, \quad (2)$$

where N denotes the number of terms used in the NS, and Θ is a $K \times K$ diagonal matrix. Note that for

(2) to work, the requirement below has to be satisfied

$$\lim_{n \rightarrow \infty} (\mathbf{I}_K - \Theta \mathbf{G})^n \rightarrow \mathbf{0}_K. \quad (3)$$

Note that \mathbf{G} is a complex central Wishart matrix because the elements of \mathbf{H} are i.i.d. complex Gaussian random variables [9]. Let $\alpha = M/K$. As K and M grow, as derived in [10], the largest and the smallest eigenvalues of \mathbf{G} converge respectively to

$$\begin{aligned} \lambda_{\max}(\mathbf{G}) &\rightarrow M \left(1 + \frac{1}{\sqrt{\alpha}}\right)^2, \\ \lambda_{\min}(\mathbf{G}) &\rightarrow M \left(1 - \frac{1}{\sqrt{\alpha}}\right)^2. \end{aligned} \quad (4)$$

As a result, if Θ is chosen as [2], which is

$$\Theta = \frac{\alpha}{M(1+\alpha)} \mathbf{I}_K = \frac{1}{M+K} \mathbf{I}_K, \quad (5)$$

then,

$$\begin{aligned} \lambda_{\max}(\Theta \mathbf{G}) &\rightarrow 1 + \frac{2\sqrt{\alpha}}{1 + \sqrt{\alpha}}, \\ \lambda_{\min}(\Theta \mathbf{G}) &\rightarrow 1 - \frac{2\sqrt{\alpha}}{1 + \sqrt{\alpha}}. \end{aligned} \quad (6)$$

Therefore, the eigenvalues of $(\mathbf{I}_K - \Theta \mathbf{G})$ lie approximately in the range of $[-2\sqrt{\alpha}/(1+\alpha), 2\sqrt{\alpha}/(1+\alpha)]$ [2], [7]. Since $2\sqrt{\alpha}/(1+\alpha) \leq 1$ when $\alpha \geq 1$, the convergence of (3) is satisfied with the choice (5). Moreover, when α is very large, $2\sqrt{\alpha}/(1+\alpha) \rightarrow 0$, which means that (3) converges very quickly. Hence, a small number of N in (2) can offer close performance to the exact inverse.

Unfortunately, for finite M and K values, the eigenvalues of the product $\Theta \mathbf{G}$ for a particular channel realization can lie outside the range of $[-2\sqrt{\alpha}/(1+\alpha), 2\sqrt{\alpha}/(1+\alpha)]$ [2], [7], which results in the failure of (3). To address this issue, an attenuation factor δ where $0 < \delta < 1$ was introduced in [2], so (5) changes to

$$\Theta = \frac{\delta}{M+K} \mathbf{I}_K. \quad (7)$$

However, the proper choice of δ is hard to be determined. On the one hand, if δ is too large, the non-convergence issue still exists. On the other hand, if δ is too small, the convergence speed becomes very slow, so the number of N needs to be very large to offer a good MIA, increasing the burden of the

hardware implementation.

Instead of (7), another Θ was applied in [7], [8], which achieves a better MIA [7]. Specifically, \mathbf{G} is decomposed as

$$\mathbf{G} = \mathbf{D} + \mathbf{E}, \quad (8)$$

where \mathbf{D} is a diagonal matrix including the diagonal elements of \mathbf{G} , and \mathbf{E} is a hollow matrix including the off-diagonal elements of \mathbf{G} . Then, Θ is chosen as

$$\Theta = \mathbf{D}^{-1}. \quad (9)$$

To achieve a good MIA with quick convergence, (9) requires that \mathbf{G} is a Diagonally Dominant Matrix (DDM) [7], [8], i.e.,

$$|g_{ii}| > \sum_{j, j \neq i} |g_{ij}|, i, j = 1, \dots, K. \quad (10)$$

The performance-complexity trade-off and hardware implementation of the NS-based MIA employing (9) have been discussed for the downlink and uplink in [7] and [8] respectively. In both cases, the NS-based MIA employing (9) was considered as a promising and practical method for massive MIMO systems. As a result, the analysis carried out in this paper is based on the choice of (9).

As mentioned in Section I, there still some issues on the application of (9) for finite M and K values. Firstly, it is unclear that what α can offer high convergence probability. Secondly, it is unclear that what α can achieve high probability for \mathbf{G} to be diagonal dominant. Moreover, more accurate approximation error analysis for practical N values is needed. In the next sections, the aforementioned issues are addressed.

III. CONVERGENCE AND α

According to the theory of matrix power series [9], for a $K \times K$ matrix \mathbf{B} , the product \mathbf{B}^N converges to $\mathbf{0}_K$ only when the spectral radius of \mathbf{B} , denoted by $\rho(\mathbf{B})$, i.e., the maximum modulus of eigenvalues of \mathbf{B} , is less than 1. Then, for the choice of (9), a good MIA of (2) requires

$$\rho(\mathbf{I}_K - \mathbf{D}^{-1}\mathbf{G}) < 1. \quad (11)$$

Since the elements of \mathbf{H} are i.i.d. zero-mean unit-variance complex Gaussian random variables, when the number of M is large, the diagonal elements of \mathbf{D} approach to $ME\{|h_{kk}|^2\} = M$ by the law of large

numbers [1], [2], [4]. Therefore, the diagonal matrix \mathbf{D} can be replaced by $M\mathbf{I}_K$, Then, the condition (11) changes to

$$|M - \lambda(\mathbf{G})| < M \Rightarrow 0 < \lambda(\mathbf{G}) < 2M. \quad (12)$$

As \mathbf{G} is a positive-definite matrix [9], its eigenvalues are all larger than 0 [9]. As a result, (12) is equivalent to

$$\lambda_{\max}(\mathbf{G}) < 2M. \quad (13)$$

Note that $\mathbf{G} = \mathbf{H}^H\mathbf{H}$ is a complex central Wishart matrix [9], and the distribution of $\lambda_{\max}(\mathbf{G})$ is provided in [11] as

$$P(\lambda_{\max}(\mathbf{G}) < x) = \frac{\mathcal{C}\Gamma_K(K)}{\mathcal{C}\Gamma_K(M+K)} x^{KM} \times {}_1F_1(M; M+K; -x\mathbf{I}), \quad (14)$$

where x is a non-negative number. The complex multivariate gamma function $\mathcal{C}\Gamma_p(a)$ is defined as

$$\mathcal{C}\Gamma_p(a) = \pi^{p(p-1)/2} \prod_{i=1}^p \Gamma[a - i + 1], \quad (15)$$

where p is a positive integer, a is a complex-valued number, and $\Gamma[a]$ is the gamma function. The hypergeometric function ${}_1F_1(M; M+K; -x\mathbf{I})$ is

$${}_1F_1(M; M+K; -x\mathbf{I}) = \sum_{k=0}^{\infty} \sum_{\kappa} \frac{[M]_{\kappa}}{[M+K]_{\kappa}} \frac{C_{\kappa}(-x\mathbf{I})}{k!}. \quad (16)$$

The details of $[M]_{\kappa}$ and $C_{\kappa}(-x\mathbf{I})$ in (16) can be found in [11]. Based on (14), the probability of (13) can be directly derived. However, (14) includes the summation of infinite terms in (16) which has extreme complexity, so it cannot provide a closed-form convergence condition of (13) in terms of α .

Fortunately, based on (4), the condition of (13) changes to

$$M \left(1 + \frac{1}{\sqrt{\alpha}}\right)^2 < 2M. \quad (17)$$

Based on (17), a high probability convergence condition in terms of α is derived as

$$\alpha > \frac{1}{(\sqrt{2} - 1)^2} \approx 5.83. \quad (18)$$

With (18), the maximum possible number of K can be found for a specific number of M to achieve a

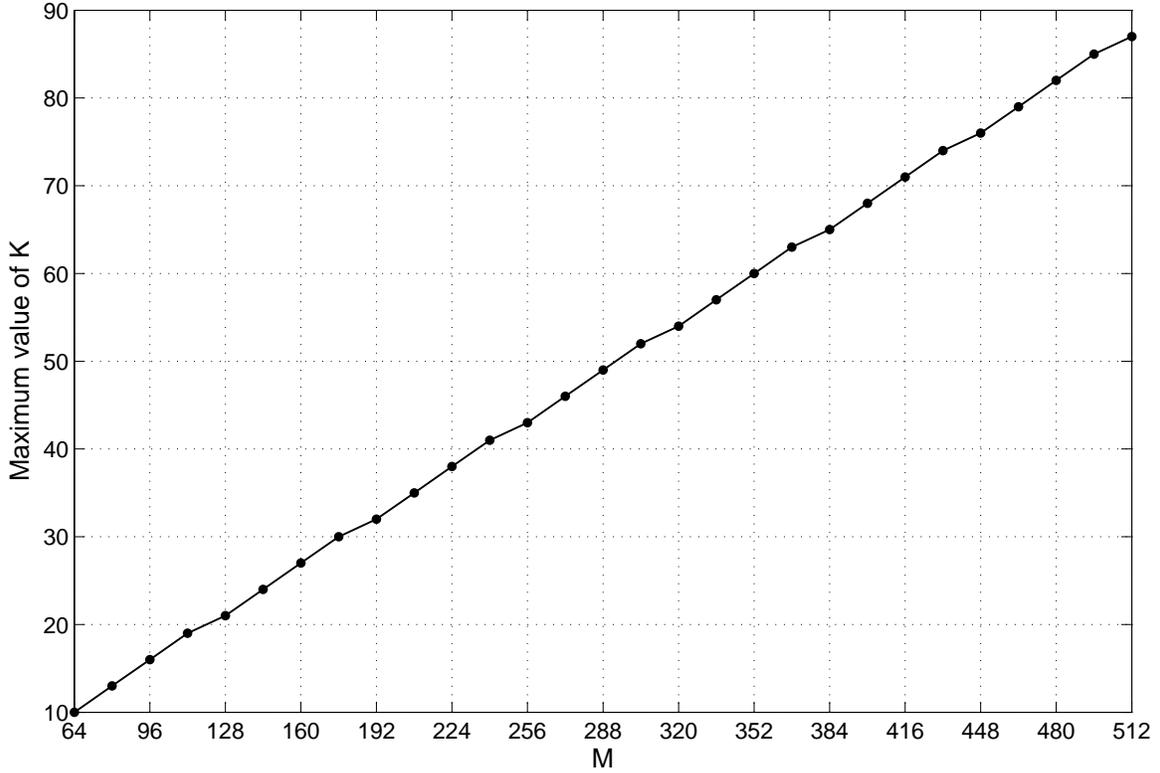


Fig. 1. The maximum K values for different M values that satisfy (18)

TABLE I
TYPICAL M VALUES WITH THEIR ASSOCIATED MAXIMUM VALUES OF K AND CONVERGENCE PROBABILITY VALUES

M	64	128	256	512
K	10	21	43	87
Probability of (3)	0.999	0.998	0.995	0.991

very high probability of convergence for (3).

Fig. 1 illustrates the maximum values of K corresponding to M values that vary from 64 to 512 based on the convergence condition (18). With these K values, the simulated convergence probability values of (3) based on the accurate condition (11) and the approximated condition (13) are shown in Fig. 2. The results indicate that they provide close probability with (11) being slightly better in massive MIMO systems with large M . The results verify that (13) is an acceptable approximation of (11). Furthermore, the results show that the condition (18) in terms of α can offer high convergence probability for (3). Table I summarizes the typical M values of massive MIMO systems with their corresponding maximum values of K and the convergence probability values of (3). Note that (18) does not ensure fast convergence of (3), so a more strict α condition for \mathbf{G} being a DDM is studied in the next section.

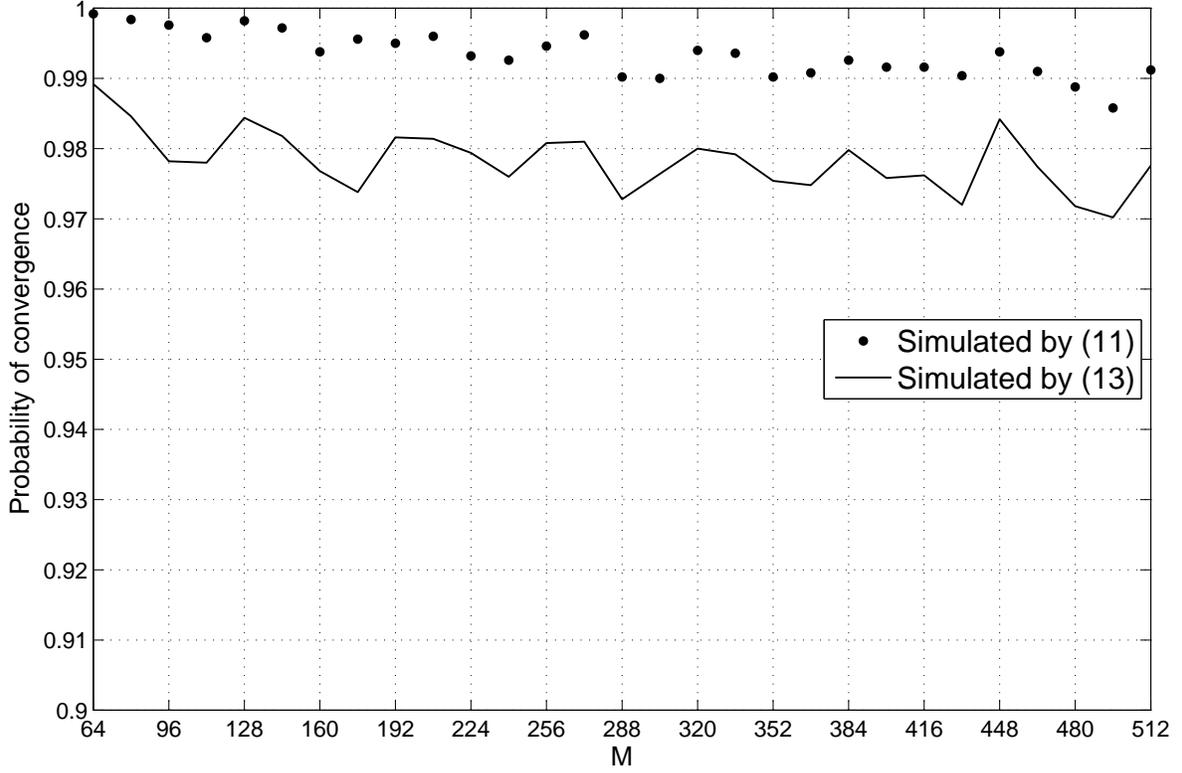


Fig. 2. Convergence probability values of (3) for the K values in Fig. 1

IV. DIAGONALLY DOMINANT AND α

Let \mathbf{h}_k denote the k th column vector of the $M \times K$ channel matrix \mathbf{H} . Then, \mathbf{h}_k represents the M -dimensional channel vector for the k th user. Hence, the elements of the $K \times K$ matrix $\mathbf{G} = \mathbf{H}^H \mathbf{H}$ is calculated as

$$\begin{cases} g_{ii} = \|\mathbf{h}_i\|_2^2, & i = 1, \dots, K, \\ g_{ij} = \mathbf{h}_i^H \mathbf{h}_j, & j = 1, \dots, K, j \neq i. \end{cases} \quad (19)$$

As mentioned in Section III, the diagonal elements g_{ii} approach to M when the number of M is large. As a result, the requirement (10) in Section II for \mathbf{G} being a DDM can be approximated as

$$\Delta_i = \sum_{j \neq i} |r_{ij}| < 1, \forall i, \quad (20)$$

where r_{ij} is the normalized correlation coefficient between \mathbf{h}_i and \mathbf{h}_j defined as

$$r_{ij} = \frac{\mathbf{h}_i^H \mathbf{h}_j}{\|\mathbf{h}_i\|_2 \|\mathbf{h}_j\|_2} \approx \frac{\mathbf{h}_i^H \mathbf{h}_j}{M}. \quad (21)$$

Let $x = |r_{ij}|$. Note that the Probability Density Function (PDF) of x was derived in [12] as

$$f(x) = 2(M-1)x(1-x^2)^{M-2}, \quad 0 \leq x \leq 1. \quad (22)$$

Hence, the mean of x is

$$E(x) = \int_0^1 xf(x) dx = (M-1)B(1.5, M-1), \quad (23)$$

where $B(a, b)$ with a and b being complex-valued numbers is the beta function defined as

$$B(a, b) = \int_0^1 t^{a-1} (1-t)^{b-1} dt, \quad \Re\{a\}, \Re\{b\} > 0. \quad (24)$$

Although (23) provides the values of $E(x)$, since the number of K is not large enough, Δ_i in (20) can be larger than $(K-1)E(x)$. However, Δ_i has a high probability being smaller than $(K-1)[E(x) + \delta(x)]$ where $\delta(x)$ denotes the standard deviation of x , which is

$$\delta(x) = \sqrt{E(x^2) - E(x)^2}, \quad (25)$$

with

$$E(x^2) = \int_0^1 x^2 f(x) dx = (M-1)B(2, M-1). \quad (26)$$

Therefore, the condition (20) can be approximated as

$$(K-1)[E(x) + \delta(x)] < 1. \quad (27)$$

Based on (27), a high probability condition for the \mathbf{G} matrix being a DDM in terms of α is derived as

$$\alpha > \frac{M[E(x) + \delta(x)]}{E(x) + \delta(x) + 1}. \quad (28)$$

With (28), the maximum possible number of K can be found for a specific number of M to achieve a very high probability for \mathbf{G} being a DDM.

Fig. 3 shows the maximum values of K corresponding to M values that vary from 64 to 512 based on the diagonally dominant condition (28). With these K values, the simulated DDM probability based on the definition (10) and the approximated condition (20) are illustrated in Fig. 4. The results show that they achieve close probability in massive MIMO systems with large M . The results verify that (20) is a good

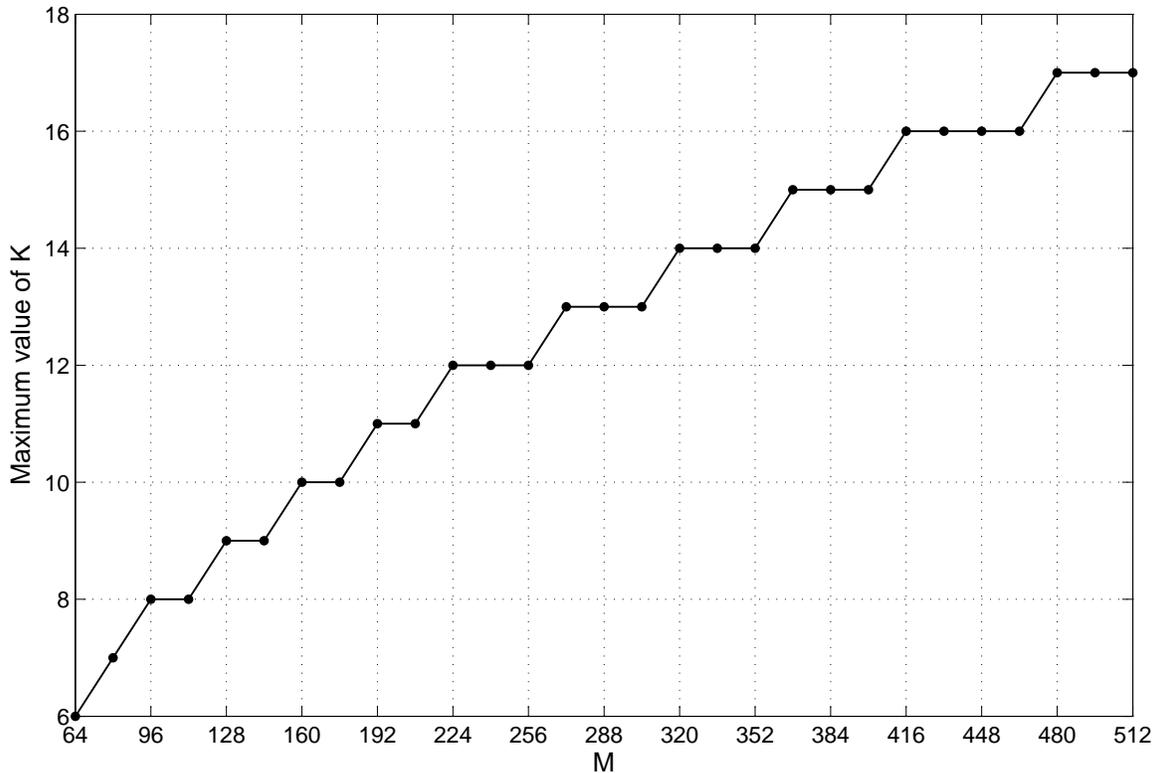


Fig. 3. The maximum K values for different M values that satisfy (28)

TABLE II
TYPICAL M VALUES WITH THEIR ASSOCIATED MAXIMUM VALUES OF K AND DIAGONALLY DOMINANT PROBABILITY VALUES

M	64	128	256	512
K	6	9	12	17
Probability of (10)	0.990	0.977	0.998	0.999

approximation of (10), especially when M is very large. Moreover, the results show that the condition (28) in terms of α can offer high DDM probability. Table II summarizes the typical M values of massive MIMO systems with their corresponding maximum values of K and the diagonally dominant probability values of (10). Note that the DDM condition (28) is sufficient for the convergence condition (18) and leads to quicker convergence, so it is more useful in practice.

V. ERROR ANALYSIS

Based on (8) and (9), the NS-based MIA of (2) changes to

$$\mathbf{G}_N^{-1} = \sum_n^{N-1} (-\mathbf{D}^{-1}\mathbf{E})^n \mathbf{D}^{-1}. \quad (29)$$

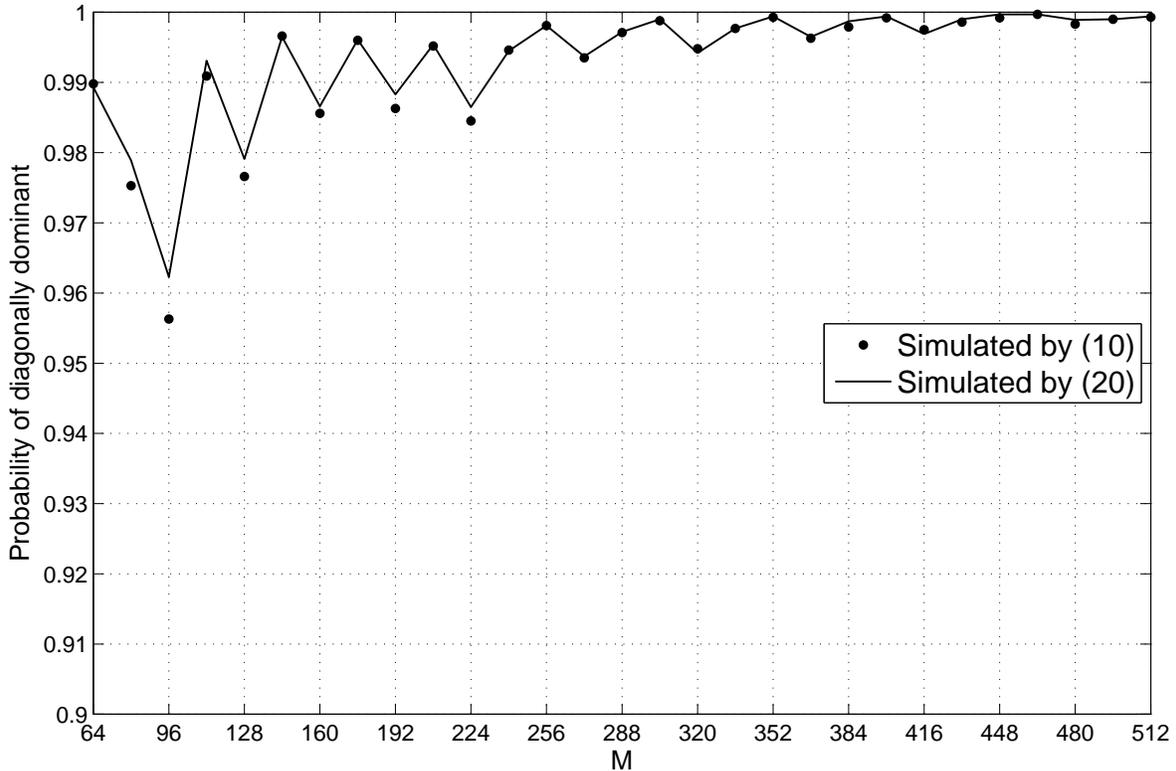


Fig. 4. Diagonally dominant probability values for the K values in Fig. 3

Note that if the convergence condition (3) is satisfied, \mathbf{G}_∞^{-1} is the exact matrix inverse of \mathbf{G} . However, in practice, the number of N cannot be very large. Otherwise, it would cause excessive burden for hardware implementation. In this case, residual error resulted from the NS-based MIA \mathbf{G}_N^{-1} exists. Let the K -dimensional vector \mathbf{s} denote the transmitted symbols for the uplink or the downlink. Without loss of generality, $E(|s_k|^2) = 1$ is assumed, with $k = 1, \dots, K$. Let $\mathbf{Z} = \mathbf{D}^{-1}\mathbf{E}$. Then, the Mean Square Error (MSE) of the NS-based MIA \mathbf{G}_N^{-1} for the uplink is derived as

$$\begin{aligned}
\epsilon_N^{\text{ul}} &= E \left\{ \left\| (\mathbf{G}_\infty^{-1} - \mathbf{G}_N^{-1}) \mathbf{H}^H \mathbf{H} \mathbf{s} \right\|_2^2 \right\} \\
&= E \left\{ \left\| \mathbf{Z}^N \sum_{n=0}^{\infty} (-\mathbf{Z})^n \mathbf{D}^{-1} \mathbf{H}^H \mathbf{H} \mathbf{s} \right\|_2^2 \right\} \\
&= E \left\{ \left\| \mathbf{Z}^N \mathbf{G}_\infty^{-1} \mathbf{G} \mathbf{s} \right\|_2^2 \right\} \\
&= E \left\{ \left\| \mathbf{Z}^N \mathbf{s} \right\|_2^2 \right\} \\
&= E \left\{ \text{Tr} \left[\mathbf{Z}^N \mathbf{s} \mathbf{s}^H (\mathbf{Z}^N)^H \right] \right\}
\end{aligned}$$

$$\begin{aligned}
&= \mathbb{E} \left\{ \text{Tr} \left[\mathbf{s} \mathbf{s}^H (\mathbf{Z}^N)^H \mathbf{Z}^N \right] \right\} \\
&= \text{Tr} \left\{ \mathbb{E} \left[\mathbf{s} \mathbf{s}^H \right] \mathbb{E} \left[(\mathbf{Z}^N)^H \mathbf{Z}^N \right] \right\} \\
&= \text{Tr} \left\{ \mathbf{I}_K \mathbb{E} \left[(\mathbf{Z}^N)^H \mathbf{Z}^N \right] \right\} \\
&= \mathbb{E} \left\{ \text{Tr} \left[(\mathbf{Z}^N)^H \mathbf{Z}^N \right] \right\} \\
&= \mathbb{E} \left\{ \|\mathbf{Z}^N\|_F^2 \right\}. \tag{30}
\end{aligned}$$

Note that for the downlink case, the MSE result is

$$\begin{aligned}
\epsilon_N^{\text{dl}} &= \mathbb{E} \left\{ \|\mathbf{s}^T (\mathbf{G}_\infty^{-1} - \mathbf{G}_N^{-1}) \mathbf{H}^H \mathbf{H}\|_2^2 \right\} \\
&= \mathbb{E} \left\{ \|\mathbf{s}^T \mathbf{Z}^N\|_2^2 \right\} \\
&= \mathbb{E} \left\{ \|\mathbf{Z}^N\|_F^2 \right\}, \tag{31}
\end{aligned}$$

which is the same as (30). Hence, ϵ_N is used instead of ϵ_N^{ul} and ϵ_N^{dl} in this section below. Since ϵ_N can be interpreted as the power of the residual interference of ZF precoding or detection, the average Signal-to-Interference Ratio (SIR) for each user is calculated as

$$\gamma_N = \frac{\frac{\|\mathbf{s}\|^2}{K}}{\frac{\epsilon_N}{K}} = \frac{K}{\epsilon_N}. \tag{32}$$

In [8], the MSE ϵ_N in (30) and (31) is upper bounded as

$$\epsilon_N \leq \left[(K^2 - K) \sqrt{\frac{2M(M+1)}{(M-1)(M-2)(M-3)(M-4)}} \right]^N, \tag{33}$$

with $M > 4$. Unfortunately, (33) is a very loose upper bound, resulting in a very loose lower bound of γ_N in (32). In Fig. 5, the exact SIR values and the lower bound values are compared with $M = 128$ for different K and M values. The results show substantial differences when $N > 1$, which cannot provide sufficient insight for the residual error of the NS-based MIA for $N > 1$. Due to this reason, we seek to derive a more accurate approximation of ϵ_N in this section below.

When M is large, because \mathbf{D} can be approximated as $M\mathbf{I}_K$ as mentioned in Section III, according to (19) and (21), the elements of \mathbf{Z} is approximated as

$$\begin{cases} z_{ii} = 0 & i = 1, \dots, K, \\ z_{ij} = z_{ji}^* \approx \frac{\mathbf{h}_i^H \mathbf{h}_j}{M} \approx r_{ij}, & j = 1, \dots, K, j \neq i. \end{cases} \tag{34}$$

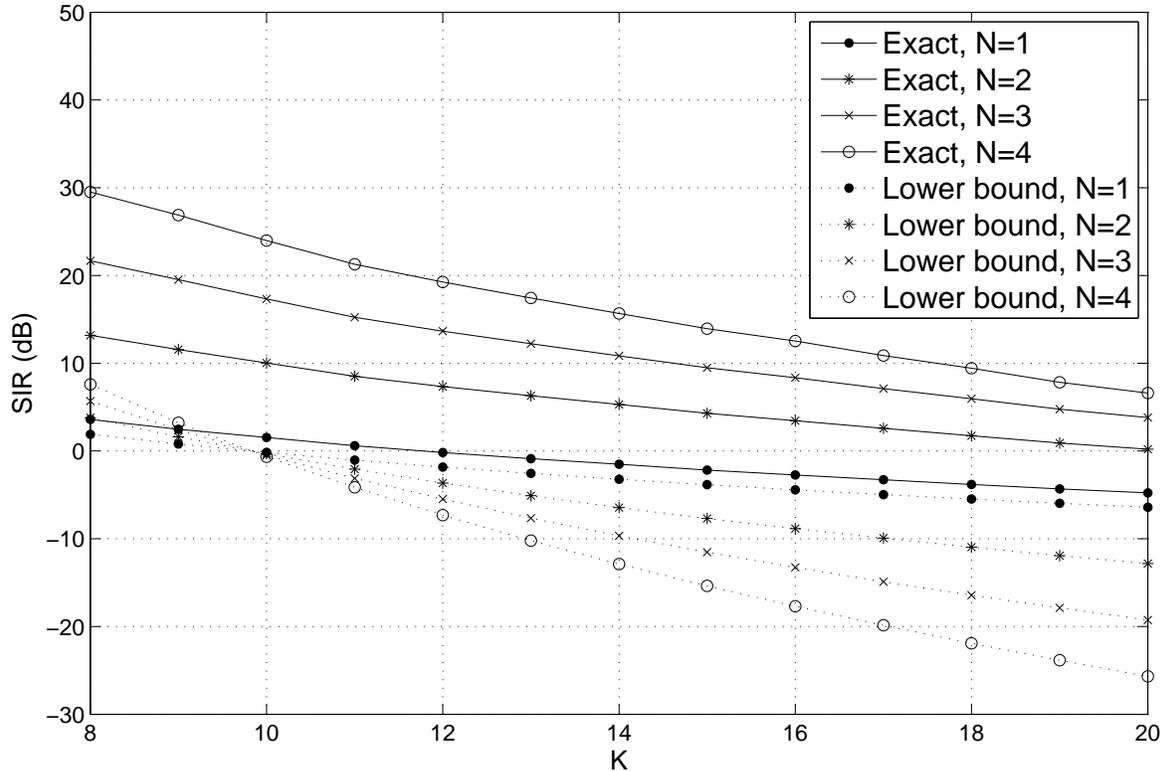


Fig. 5. Comparison between the exact SIR values and the lower bound values resulting from (33) with $M = 128$ for different K and N values.

As a result, the PDF of $x = |z_{ij}|$ can be approximated as (22). Then, a more accurate approximation of γ_N can be derived based on (22).

When $N = 1$, the MSE ϵ_N in (30) and (31) changes to

$$\epsilon_1 = \|\mathbf{Z}\|_F^2 = \sum_{i=1}^K \sum_{j=1, j \neq i}^K |z_{ij}|^2 \approx K(K-1)E(x^2). \quad (35)$$

Since $E(x^2)$ has been derived as (26), the term ϵ_1 in (35) is rewritten as

$$\epsilon_1 \approx K(K-1)B_{2,M}, \quad (36)$$

where $B_{a,M}$ is defined as

$$B_{a,M} = (M-1)B(a, M-1). \quad (37)$$

When $N = 2$, the MSE ϵ_N in (30) and (31) changes to

$$\epsilon_2 = \|\mathbf{Z}^2\|_F^2, \quad (38)$$

where the elements in $\mathbf{Y} = \mathbf{Z}^2$ is

$$\begin{cases} y_{ii} = \sum_{k=1, k \neq i}^K |z_{ik}|^2, & i = 1, \dots, K, \\ y_{ij} = \sum_{k=1, k \neq i, j}^K z_{ik} z_{jk}^*, & j = 1, \dots, K, j \neq i. \end{cases} \quad (39)$$

Note that $\|\mathbf{Z}^2\|_{\mathbb{F}}^2$ can be written as a summation of polynomial terms, which can be classified into three categories. The first category includes $(K-1)K$ terms of $|z_{ik}|^4$ with $i \neq k$. The second category includes $(K-2)(K-1)K$ terms of $|z_{ik}|^2|z_{il}|^2$ with $i \neq k \neq l$, as well as $(K-2)(K-1)K$ terms of $|z_{ik}|^2|z_{jk}|^2$ with $i \neq j \neq k$. Hence, the total number of terms for the second category is $2(K-2)(K-1)K$. Finally, the third category includes $(K-3)(K-2)(K-1)K$ terms of $z_{ik}z_{jk}^*z_{il}^*z_{jl}$ with $i \neq j \neq k \neq l$. Because the elements of \mathbf{H} are i.i.d zero-mean unit-variance complex Gaussian random variables, based on (34), the elements of z_{ij} are i.i.d. zero-mean random variables. As a result, the terms of the third category are also i.i.d. zero-mean random variable. Therefore, the sum of the terms of the third category can be approximated as zero. For the terms of the first category, the mean can be calculated based on (22) as

$$\mathbb{E}(x^4) = \int_0^1 x^4 f(x) dx = B_{3,M}. \quad (40)$$

Similarly, the mean of the terms of the second category is can be approximated as

$$\mathbb{E}(x_1^2 x_2^2) = \mathbb{E}(x^2)^2 = B_{2,M}^2. \quad (41)$$

Due to (40) and (41), the term ϵ_2 in (38) is approximated as

$$\epsilon_2 \approx K(K-1)B_{3,M} + 2(K-2)(K-1)KB_{2,M}^2. \quad (42)$$

When $N > 2$, the MSE ϵ_N in (30) and (31) can be derived with the similar method applied by $N = 2$. The results of $N = 3$ and $N = 4$ are directly provided below as

$$\begin{aligned} \epsilon_3 &\approx (K-2)(K-1)K(5K-8)B_{2,M}^3 \\ &\quad + (2K-3)(K-1)KB_{3,M}B_{2,M}, \end{aligned} \quad (43)$$

and

$$\epsilon_4 \approx (2K-3)(K-1)KB_{3,M}^2$$

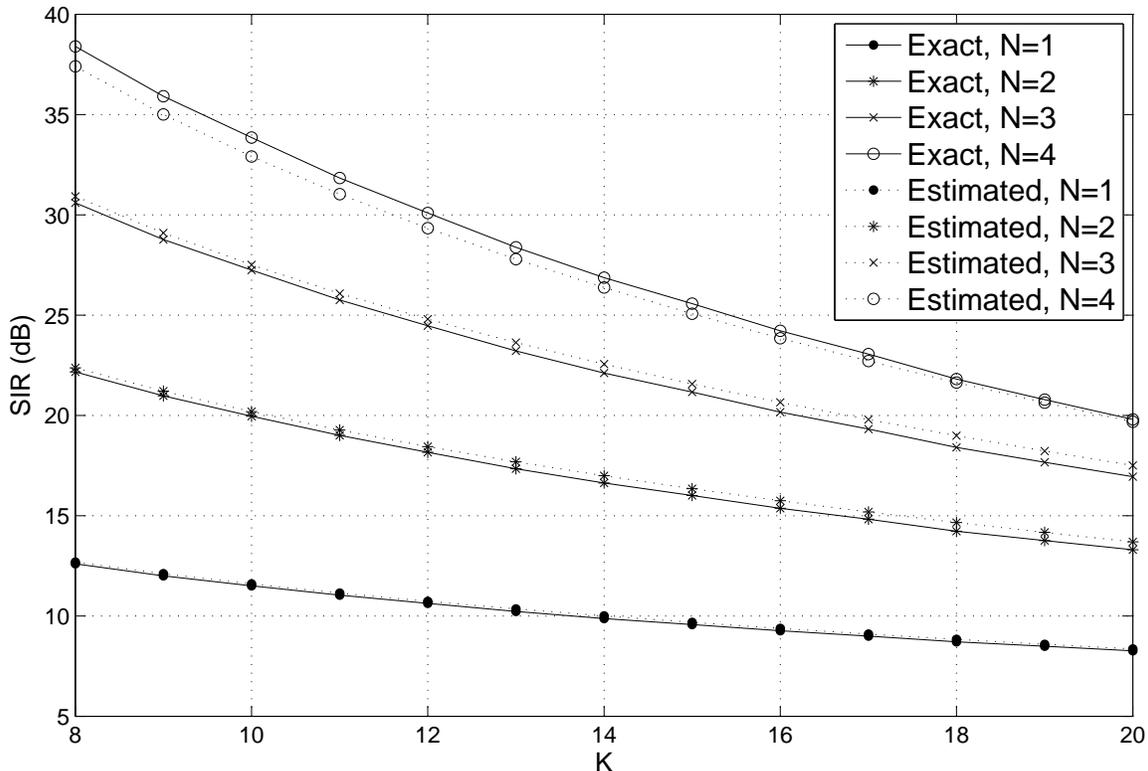


Fig. 6. Comparison between the exact and estimated SIR values with $M = 128$ for different K and N values.

$$\begin{aligned}
& + (2K - 3)^2 (K - 1)^2 K B_{4,M} B_{2,M} \\
& + (K - 2) (K - 1)^2 K^2 B_{2,M}^4.
\end{aligned} \tag{44}$$

With the estimated residual error formulas (36), (42)-(44), the estimated SIR formulas can be easily derived according to (32). Fig. 6-8 compare the exact and estimated SIR values for different K and N values, with $M = 128$, $M = 256$, and $M = 512$ respectively. The results show that the estimated SIR values are very close to the exact SIR values, which verifies the high accuracy of SIR estimation formulas based on (36), (42)-(44).

VI. DISCUSSIONS

In massive MIMO systems, α is commonly considered to be very large to offer good performance [2], [4], e.g., $\alpha > 10$. Hence, the convergence condition (18), i.e., $\alpha > 5.83$, derived in Section III is generally satisfied for massive MIMO systems. Note that the convergence probability values provided in Fig. 1 and Table I, which are already close to 1, correspond to the smallest α values that satisfy (18). Hence, the convergence probability values for massive MIMO systems are not lower than the values provided in Fig.

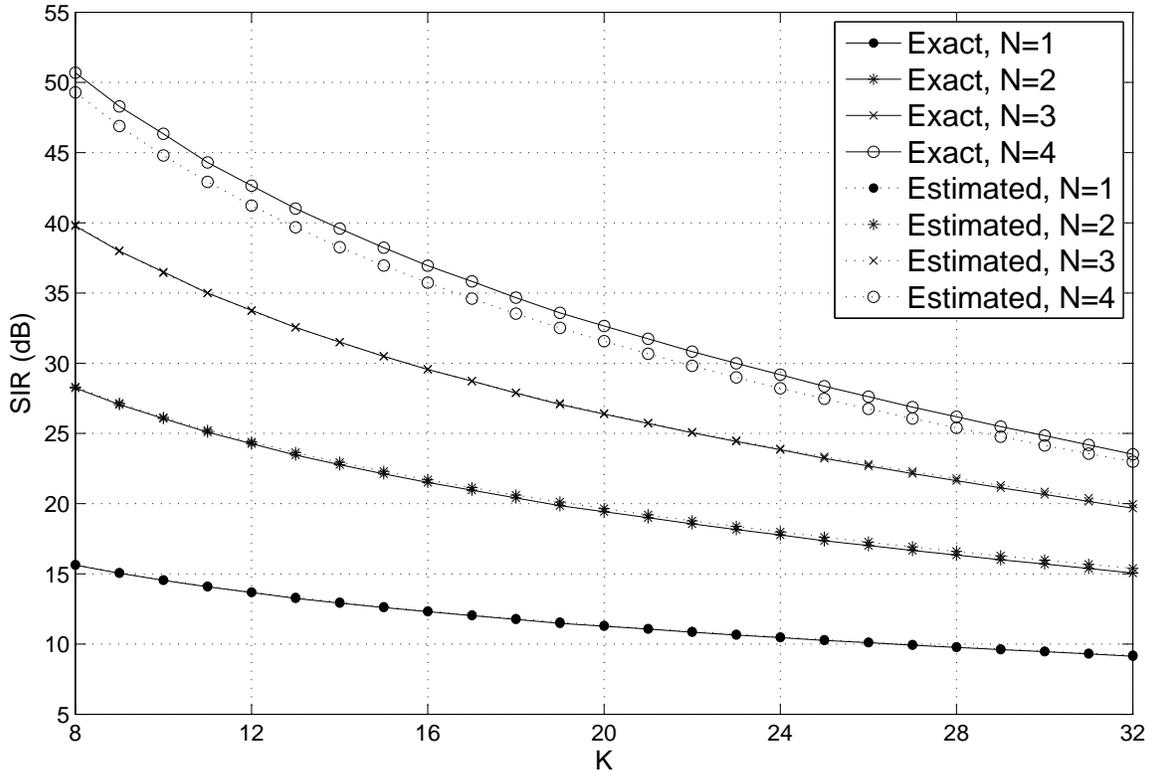


Fig. 7. Comparison between the exact and estimated SIR values with $M = 256$ for different K and N values.

1 and Table I. Therefore, the convergence of NS-based MIA is guaranteed so that it is a valid method for massive MIMO systems, and its accuracy can be improved by increasing N .

As mentioned at the end of Section III, the convergence condition (18) does not guarantee quick convergence of (3). With the diagonally dominant condition (28) derived in Section IV, however, the NS-based MIA can achieve good accuracy with quick convergence, i.e., a small N can offer a sufficiently good MIA. Otherwise, with the same N value, violating (28) results in performance loss for the ZF decoding or detection employing the NS-based MIA. Take the simulation results provided in [8] as examples, with $M = 128$ and $N = 3$, the choice of $K = 4$ satisfying (28) achieves close performance to the exact inverse, while the choice of $K = 12$ violating (28) suffers huge performance loss. However, (28) requires very small α values, and α becomes smaller as M increases, which can be seen from Table IV. The strict requirement of α may reduce the spatial multiplexing advantage of massive MIMO systems, i.e., at most $K = 17$ users can be served by $M = 512$ antennas. To relieve this issue, one comprised choice is to apply an α slightly higher than (28) with slightly larger N of the NS-based MIA, depending on the hardware capability.

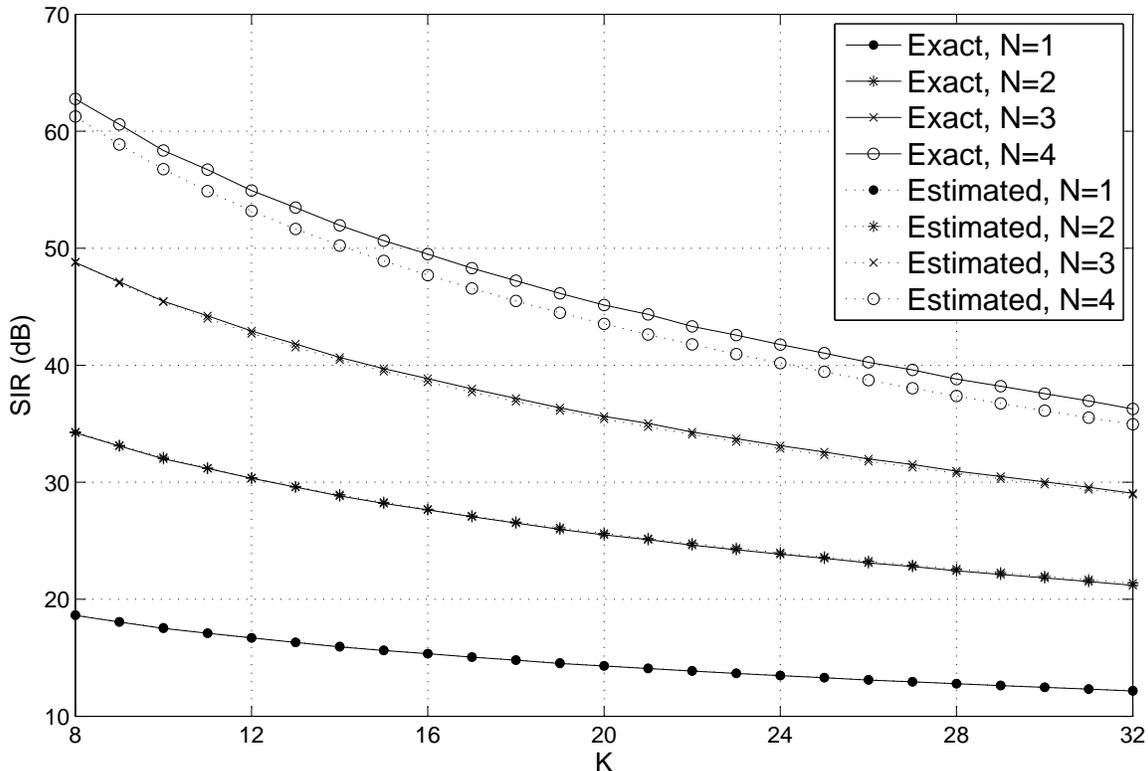


Fig. 8. Comparison between the exact and estimated SIR values with $M = 512$ for different K and N values.

The SIR discussed in Section V reflects the performance error floor for ZF precoding or detection employing practical NS-based MIA in massive MIMO systems. The performance error floor decides the best performance that the ZF precoding or detection employing the NS-based MIA can achieve. As a result, with M , K , and N , the best achievable performance can be easily estimated based on (36), (42)-(44). In addition, since larger N causes higher hardware implementation complexity, with M , K , and the target performance, the smallest choice of N that can offer sufficiently good performance can be determined to relieve the complexity. Note that a revised form of (2) was provided in [7] as

$$\begin{aligned} \mathbf{G}_N^{-1} &\approx \sum_{n=0}^{N-1} (\mathbf{I}_K - \Theta \mathbf{G})^n \Theta \\ &= \prod_{l=0}^{L-1} \left[\mathbf{I}_K + (\mathbf{I}_K - \Theta \mathbf{G})^{2^l} \right] \Theta, \end{aligned} \quad (45)$$

where L is a positive integer with $N = 2^L$. Hence, $L = 1$, $L = 2$, and $L = 3$ of the alternative expression (45) correspond to $N = 2$, $N = 4$, and $N = 8$ of the regular expression (2) respectively. As a result, with the alternative expression (45), after the choice of $N = 4$, the NS-based MIA with the choice of $N = 8$

can be quickly calculated. Therefore, if the choice of $N = 4$ is not good enough based on the estimation formula (44), the choice of $N = 8$ can be directly selected based on (45). Furthermore, note that the complexity of the NS-based MIA with the choice of $N > 3$ is considered to be $O(K^3)$ in [8], which loses the complexity advantage over the exact matrix inverse of $O(K^3)$. In fact, however, the NS-based MIA can be implemented as a series of cascaded matched filter so that the complexity can be reduced to $O(K^2)$, as discussed in [2]. In this way, the NS-based MIA still has the complexity advantage over the exact inverse even with the choice of $N = 8$.

VII. CONCLUSIONS

In this paper, three issues related to the practical application of the NS-based MIA in massive MIMO systems are addressed. Firstly, $\alpha > 5.83$ as in (18) is offered for the NS-based MIA to achieve very high convergence probability. In other words, with the number of BS antennas M , the maximum number of served users K for the NS-based MIA to be a valid method in massive MIMO systems can be determined. Then, a tighter condition (28) is provided for \mathbf{G} to be a DDM in very high probability, resulting in a good NS-based MIA with a small number of N . This means that given the number of BS antennas M , the maximum number of served users K for the NS-based MIA to achieve good performance and quick convergence for ZF decoding or detection can be determined. Finally, by approximation error analysis, residual error estimation formulas (36), (42)-(44) with very high accuracy are derived for practical N values, which can be applied to estimate the error floor caused by the NS-based MIA. Thus, given the number of BS antennas M , the number of served users K , and the number of terms employed by the NS-based MIA N , highly accurate estimation of the SIR caused by the NS-based MIA can be obtained. These results offer useful guidelines for practical application of the NS-based MIA in massive MIMO systems.

REFERENCES

- [1] T. L. Marzetta, "Noncooperative Cellular Wireless with Unlimited Numbers of Base Station Antennas," *IEEE Trans. Wireless Commun.*, vol. 9, no. 11, pp. 3590–3600, Nov. 2010.
- [2] F. Rusek, D. Persson, B. K. Lau, E. G. Larsson, T. L. Marzetta, O. Edfors, and F. Tufvesson, "Scaling up MIMO: Opportunities and Challenges with Very Large Arrays," *IEEE Signal Process. Mag.*, vol. 30, no. 1, pp. 40–46, Jan. 2013.
- [3] E. G. Larsson, F. Tufvesson, O. Edfors, and T. L. Marzetta, "Massive MIMO for Next Generation Wireless Systems," *IEEE Commun. Mag.*, vol. 52, no. 2, pp. 186–195, Feb. 2014.
- [4] J. Hoydis, S. Brink, and M. Debbah, "Massive MIMO in the UL/DL of Cellular Networks: How Many Antennas Do We Need?" *IEEE Sel. Areas Commun.*, vol. 31, no. 2, pp. 160–171, Feb. 2013.

- [5] C. Shepard, H. Yu, N. Anand, L. E. Li, T. Marzetta, R. Yang, and L. Zhong, "Argos: Practical Many-Antenna Base Stations," in *Proc. MobiCom 12*, Istanbul, Turkey, Aug. 2012.
- [6] H. Yang and T. L. Marzetta, "Performance of Conjugate and Zero-Forcing Beamforming in Large-scale Antenna Systems," *IEEE J. Sel. Areas Commun.*, vol. 31, no. 2, pp. 172–179, Feb. 2013.
- [7] H. Prabhu, J. Rodrigues, O. Edfors, and F. Rusek, "Approximative Matrix Inverse Computations for Very-large MIMO and Applications to Linear Pre-coding Systems," in *Proc. IEEE WCNC 13*, Shanghai, China, Apr. 2013.
- [8] M. Wu, B. Yin, G. Wang, C. Dick, J. Cavallaro, and C. Studer, "Large-Scale MIMO Detection for 3GPP LTE: Algorithms and FPGA Implementations," *IEEE J. Sel. Topics Signal Process.*, vol. 8, no. 5, pp. 916–929, Oct. 2014.
- [9] R. A. Horn and C. R. Johnson, *Matrix Analysis*. Cambridge Univeristy Press, 1990.
- [10] A. Edelman, "Eigenvalues and Condition Numbers of Random Matrices," Ph.D. dissertation, MIT, 1989.
- [11] T. Ratnarajah, R. Vaillancourt, and M. Alvo, "Eigenvalues and Condition Numbers of Complex Random Matrices," *SIAM Journal on Matrix Analysis and Applications*, vol. 26, no. 2, pp. 441–456, 2004.
- [12] D. Zhu, B. Li, and P. Liang. (2014) Normalized Volume of Hyperball in Complex Grassmann Manifold and Its Application in Large-Scale MU-MIMO Communication Systems. arXiv:1402.4543. [Online]. Available: <http://arxiv.org>