*Article*

# Joint Beam-Forming, User Clustering and Power Allocation for MIMO-NOMA Systems

## Jiayin Wang, Yafeng Wang * and Jiarun Yu

School of Information and Communication Engineering, Beijing University of Posts and Telecommunications, Beijing 100876, China; jiayinwang@bupt.edu.cn (J.W.); yujiarun@bupt.edu.cn (J.Y.)
* Correspondence: wangyf@bupt.edu.cn; Tel.: +86-188-1103-1862

**Abstract:** In this paper, we consider the optimal resource allocation problem for multiple-input multiple-output non-orthogonal multiple access (MIMO-NOMA) systems, which consists of beam-forming, user clustering and power allocation, respectively. Users can be divided into different clusters, and the users in the same cluster are served by the same beam vector. Inter-cluster orthogonality can be guaranteed based on multi-user detection (MUD). In this paper, we propose a three-step framework to solve the multi-dimensional resource allocation problem. In step 1, we propose a beam-forming algorithm for a given user cluster. Specifically, fractional transmitting power control (FTPC) is applied for intra-cluster power allocation. The considered beam-forming problem can be transformed into a non-constrained one and the limited-memory Broyden–Fletcher–Goldfarb–Shanno (L-BFGS) method is applied to obtain the optimal solution. In step 2, optimal user clustering is further considered. Channel differences and correlations are both involved in the design of user clustering. By assigning different weights to the two factors, we can produce multiple candidate clustering schemes. Based on the proposed beam-forming algorithm, beam-forming can be done for each candidate clustering scheme to compare their performances. Moreover, based on the optimal user clustering and beam-forming schemes, in step 3, power allocation can be further optimized. Specifically, it can be formalized as a difference of convex (DC) programming problem, which is solved by successive convex approximation (SCA) with strong robustness. Simulations results show that the proposed scheme can effectively improve spectral efficiency (SE) and edge users' data rates.

**Keywords:** NOMA; MIMO; user clustering; power allocation; analog beam forming

## 1. Introduction

Traditional orthogonal multiple access (OMA) has met a bottleneck, since the limited spectrum resources cannot meet the ever-growing demand for mobile data traffic. As an alternative, non-orthogonal multiple access (NOMA) has attracted considerable attention since it allows multiple users to occupy the same spectrum resource simultaneously. According to NOMA protocols, users can be divided into different clusters based on their channel characteristics. The signals of the users in the same cluster will be further transmitted utilizing the same time-frequency resource [1]. In each cluster, the channel differences among different users should be large enough to perform successive interference cancellation (SIC) successfully [2–4]. Moreover, weak users can be compensated in the power allocation process, which not only improves edge users' performances, but helps to better identify multiplexed users in the power domain [5–8].

Moreover, multiple-input multiple-output (MIMO) also serves as a promising technique by which to multiply the spectrum efficiency (SE) gain [9–11]. In massive MIMO systems, beam-forming can effectively improve SE based on spacial diversity [12]. Conventionally, a specific beam vector can be designed for each user. The interference among multiple users can be eliminated when the number of antennas is greater than that of users. Specifically, the beam vector of each user can be set orthogonal to the channel vectors of others based on the zero-forcing beam-forming (ZF-BF) algorithm [13,14].

In this passage, SE can be further improved by exploiting both the spacial and power domain, i.e., MIMO-NOMA. We divide users into different clusters and further design a beam vector for each cluster. For each user, the inter-cluster interference can be transformed into the inter-beam interference, which is further eliminated based on ZF-BF [14–16]. Moreover, the beam vector of each cluster mainly depends on the channel characteristics of strong intra-cluster users [17]. The gaps between strong users and weak users will continue to widen, which is favored by NOMA.

## 1.1. Related Works

MIMO-NOMA has received considerable research interest for its ability to improve SE. Optimal user clustering for a downlink NOMA system was considered in [4], where the users were divided into different clusters based on an improved sorting algorithm. The authors of [5] applied NOMA to a MIMO system and demonstrated that the combined application can bring extra SE improvements only when the channel correlations among multiplexed users were sufficiently high. In [9], beam-forming and power allocation were jointly optimized for MIMO-NOMA based on semi-definite programming (SDP). In [12], simultaneous wireless information and power transfer (SWIPT) was applied in cooperation with NOMA, with the aim of enhancing edge users' data rates. The angle domain was exploited in [16] to identify those users occupying the same spectrum resources, and the beam-forming problem was further considered based on an estimation of users' angle information. In [17], receiver antenna selection (RAS) was applied in an uplink MIMO-NOMA system to ensure that cell-edge users could be more likely to participate in the communication process. The authors of [18] proposed a beam-forming algorithm for MIMO-NOMA. An effective channel vector was obtained for each cluster to describe the channel characteristics of intra-cluster users, which provided a compatible dimension for ZF-BF. Moreover, high-speed beam-forming for MIMO-NOMA was studied in [19]. The authors of [20] integrated device-to-device (D2D) communications with MIMO-NOMA to further improve SE. Consequently, in [20], a novel resource allocation scheme was proposed for the integrated system to overcome interference. Research on the problem of resource allocation in NOMA has also been expanded to a multi-cell scenario. [21] investigated the resource allocation problem for multi-cell MIMO-NOMA-based internet of things (IoT) networks. Moreover, [22] investigated the energy efficiency (EE) maximization problem for multi-cell, massive MIMO-NOMA networks with wireless power transfer. The authors of [22] proposed a novel joint power, time, antenna and subcarrier allocation scheme, which could properly allocate the time for energy harvesting and data transmission.

## 1.2. Our Contributions

In MIMO-NOMA systems, SIC is of great significance in reducing intra-cluster interference. While a user is decoding the signals of others based on SIC, past research has tended to set a lower bound for the received signal to interference and noise ratio (SINR) to ensure the decoding process goes smoothly. In [23], the optimal power allocation scheme for the downlink of NOMA system was obtained based on Karush-Kuhn-Tucker (KKT) conditions with a SINR bound of 0.3. Additionally, [9] jointly optimized power allocation and beam-forming for MIMO-NOMA with a SINR bound less than 0.5. Unfortunately, most related works fail to obtain a feasible solution when the SINR bound is greater than 1, which makes the received SINR for users relatively lower and, in turn, decreases system reliability. One explanation for this is that, most related works usually consider the joint optimization of power allocation and beam-forming. The scale of the considered problem is relatively large, which makes it challenging for optimization tools to obtain a feasible solution. To address this issue, we decompose the multi-dimensional resource allocation problem into three sub-problems. The scope of each sub-problem is relatively small, which helps to obtain a feasible solution with strong robustness.

Moreover, most existing works only consider channel difference characteristics while determining the clustering scheme for MIMO-NOMA. However, since the users in the same cluster are served by the same beam vector, their channel correlations should be relatively high to bring the advantages of MIMO into full play. The clustering criterion in [9–12] was to make the channel differences among multiplexed users as big as possible, which neglected channel correlations characteristics and was not directly related to the ultimate system performance. In this paper, channel correlations and differences are both involved in the design of user clustering. By assigning different weights to the two factors, we can produce multiple possible clustering schemes. Beam-forming can be done for each possible clustering scheme to compare their performance, which ensures the ultimate clustering scheme achieves the maximum SE performance.

In addition, some related literature only considers the resource allocation problem for a two-user-cluster. In this paper, the size of each cluster is not fixed, which makes the proposed scheme more practical. The main contributions of this paper are summarized as below:

1. We present a system model for MIMO-NOMA. Multiple users can be divided into different clusters and the size of each cluster is not fixed. The users in the same cluster are served by the same beam vector. Each user is assumed to detect signals based on a specific receiving coefficient to ensure inter-cluster orthogonality. Moreover, SIC is applied to users to alleviate intra-cluster interference.

2. We propose a three-step framework to solve the multi-dimensional resource allocation problem. In step 1, a beam-forming algorithm is proposed to obtain the optimal beam vector for a given user cluster. Specifically, fractional transmitting power control (FTPC) is applied to perform intra-cluster power allocation. The considered beam-forming problem can be transformed into a non-constrained one and the limited-memory Broyden–Fletcher–Goldfarb–Shanno (L-BFGS) method is applied to obtain a local optimal with less complexity.

3. In step 2, user clustering is further considered, based on the proposed beam-forming algorithm. For each user $k$, we define a utility function to describe its preference on each cluster $n$. The utility function consists of two terms, which depict the channel differences and correlations between user $k$ and the existing users in cluster $n$, respectively. A relative weight is introduced for the two factors to balance the tradeoff between channel differences and correlations. Based on the utility function, user $k$ can be further assigned to its favorite cluster. In this paper, the relative weight is obtained by particle swarm optimization (PSO). In PSO, we can simultaneously produce multiple possible solutions for the relative weight, each corresponding to a possible clustering scheme. Based on the proposed beam-forming algorithm, beam-forming can be done for each possible clustering scheme to compare their performance, which ensures the ultimate clustering scheme achieves the maximum SE performance.

4. In step 3, power allocation is further optimized based on the optimal user clustering and beam-forming schemes. As mentioned before, it can be formalized as a difference of convex (DC) programming problem utilizing the specific characteristic of the objective function, which can be solved by successive convex approximation (SCA) through limited iterations. We evaluate the performance of the proposed scheme and some other existing schemes to illustrate the significance of the proposed scheme.

The rest of the paper is organized as follows: Section 2 presents the system model for MIMO-NOMA and further provides a mathematical expression of the optimal resource allocation problem. Section 3 introduces more details about the proposed beam-forming algorithm. Sections 4 and 5 introduce the user clustering and power allocation schemes, respectively. The performance of the proposed scheme is evaluated in Section 6. Section 7 concludes this paper.

## 2. System Model

Consider a single-cell downlink MIMO-NOMA system, in which there is one base station (BS) equipped with $N$ antennas and $K$ single-antenna users. Let $\mathbf{U} = \{1, 2, \ldots, K\}$ denote the set of users. Without loss of generality, the users are indexed by the descending order of channel gains, i.e., $|\mathbf{h_1}|^2 > |\mathbf{h_2}|^2 > \ldots > |\mathbf{h_K}|^2$, where $\mathbf{h_k} \in \mathbb{C}^{N \times 1}$ ($k \in \{1, 2, \ldots, K\}$) denotes the channel vector of user $k$. All the $K$ users will be further divided into $S$ different clusters. Let $\mathbf{U_n} = \{i_n(1), i_n(2), \ldots, i_n(m_n)\}$ ($n \in \{1, 2, \ldots, S\}$) denote the set of users assigned to cluster $n$, where $m_n$ denotes the size of $\mathbf{U_n}$, and $i_n(l)$ ($l \in \{1, 2, \ldots, m_n\}$) denotes the index of the $l$-th user in $\mathbf{U_n}$. Specifically, the users in $\mathbf{U_n}$ are sorted in the ascending order of their indexes, i.e., $i_n(1) < i_n(2) < \ldots < i_n(m_n)$. The superposed signal at the BS is given by

$$\mathbf{ss} = \sum_{n=1}^{S} \mathbf{w_n} \sum_{l=1}^{m_n} \sqrt{p_{i_n(l)}} s_{i_n(l)} \tag{1}$$

where $\mathbf{w_n} \in \mathbb{C}^{N \times 1}$ denotes the beam vector of cluster $n$, $s_{i_n(l)}$ and $p_{i_n(l)}$ denote the signal and power of user $i_n(l)$, respectively. Assume the beam vector of each cluster has constant modulus (CM) elements. The received signal at user $i_n(l)$ is given by

$$
\begin{aligned}
r_{i_n(l)} = \; & \mathbf{h_{i_n(1)}}^H \mathbf{w_n} \sqrt{p_{i_n(l)}} s_{i_n(l)} \\
& + \mathbf{h_{i_n(1)}}^H \underbrace{\sum_{q \in \{1,2,\ldots,S\} \backslash \{n\}} \mathbf{w_q} \sum_{a=1}^{m_q} \sqrt{p_{i_q(a)}} s_{i_q(a)}}_{\text{inter-cluster-interference}} \\
& + \mathbf{h_{i_n(1)}}^H \mathbf{w_n} \underbrace{\sum_{b \in \{1,2,\ldots,m_n\} \backslash \{l\}} \sqrt{p_{i_n(b)}} s_{i_n(b)}}_{\text{intra-cluster-interference}} + \omega_{i_n(l)}
\end{aligned} \tag{2}
$$

where $\mathbf{h_{i_n(1)}}$ denotes the channel vector of user $i_n(l)$. The first term in (2) represents the received desired signal. The second and third term represent the inter-cluster and intra-cluster interference, respectively. The noise term $\omega_{i_n(l)}$ is a zero-mean complex additive white Gaussian noise (AWGN) with variance $\sigma^2$. One can observe that the received interference is significantly larger. To solve this problem, each user is assumed to detect signals via a specific receiving coefficient, given by

$$\alpha_{i_n(l)} = \mathbf{v}_{\mathbf{i_n(1)}}^H \mathbf{h_{i_n(1)}} \tag{3}$$

where $\alpha_{i_n(l)}$ denotes the receiving coefficient of user $i_n(l)$, $\mathbf{v_{i_n(1)}} \in \mathbb{C}^{N \times 1}$. Then, the received signal at user $i_n(l)$ can be re-written as

$$
\begin{aligned}
\bar{r}_{i_n(l)} = \; & \alpha_{i_n(l)} r_{i_n(l)} \\
= \; & \mathbf{v}_{\mathbf{i_n(1)}}^H \mathbf{H_{i_n(1)}} \mathbf{w_n} \sqrt{p_{i_n(l)}} s_{i_n(l)} \\
& + \mathbf{v}_{\mathbf{i_n(1)}}^H \mathbf{H_{i_n(1)}} \sum_{q \in \{1,2,\ldots,S\} \backslash \{n\}} \mathbf{w_q} \sum_{a=1}^{m_q} \sqrt{p_{i_q(a)}} s_{i_q(a)} \\
& + \mathbf{v}_{\mathbf{i_n(1)}}^H \mathbf{H_{i_n(1)}} \mathbf{w_n} \sum_{b \in \{1,2,\ldots,m_n\} \backslash \{l\}} \sqrt{p_{i_n(b)}} s_{i_n(b)} + \bar{\omega}_{i_n(l)}
\end{aligned} \tag{4}
$$

where $\mathbf{H_{i_n(1)}} = \mathbf{h_{i_n(1)}} \mathbf{h}_{\mathbf{i_n(1)}}^H$. For user $i_n(l)$, the interference from cluster $q$ ($q \neq n$) can be eliminated when the following condition satisfies:

$$\mathbf{v}_{\mathbf{i_n}(1)}^{H}\mathbf{H}_{\mathbf{i_n}(1)}\mathbf{w_q} = 0 \tag{5}$$

Let $\tilde{\mathbf{h}}_{\mathbf{n,l}}^{\mathbf{q}} = \mathbf{H}_{\mathbf{i_n}(1)}\mathbf{w_q}$, and let $\tilde{\mathbf{H}}_{\mathbf{n,l}} = [\tilde{\mathbf{h}}_{\mathbf{n,l}}^{\mathbf{1}}, \tilde{\mathbf{h}}_{\mathbf{n,l}}^{\mathbf{2}}, \dots, \tilde{\mathbf{h}}_{\mathbf{n,l}}^{\mathbf{n-1}}, \tilde{\mathbf{h}}_{\mathbf{n,l}}^{\mathbf{n+1}}, \dots, \tilde{\mathbf{h}}_{\mathbf{n,l}}^{\mathbf{S}}]$. For user $i_n(l)$, the inter-cluster interference can be totally eliminated by setting $\mathbf{v}_{\mathbf{i_n}(1)}$ as the left singular vector of $\tilde{\mathbf{H}}_{\mathbf{n,l}}$ corresponding to the zero singular value. It is worth noting that there is a constraint for this operation, i.e., $N \geq S - 1$. In addition, $\alpha_{i_n(l)}$ should be normalized to ensure that it will not bring extra SE gains, i.e., $\left|\mathbf{v}_{\mathbf{i_n}(1)}^{H}\mathbf{h}_{\mathbf{i_n}(1)}\right| = 1$.

Accordingly, the received signal at user $i_n(l)$ can be transformed into

$$\begin{aligned} \bar{r}_{i_n(l)} = \ & \mathbf{v}_{\mathbf{i_n}(1)}^{H}\mathbf{H}_{\mathbf{i_n}(1)}\mathbf{w_n}\sqrt{p_{i_n(l)}}s_{i_n(l)} \\ & + \mathbf{v}_{\mathbf{i_n}(1)}^{H}\mathbf{H}_{\mathbf{i_n}(1)}\mathbf{w_n}\sum_{b\in\{1,2,\dots,m_n\}\setminus\{l\}}\sqrt{p_{i_n(b)}}s_{i_n(b)} + \bar{\omega}_{i_n(l)} \end{aligned} \tag{6}$$

Moreover, SIC is performed at users to further reduce intra-cluster interference. According to SIC, in each cluster, a user can decode the signals of the others with poorer channel conditions. Conventionally, since $\mathbf{U_n}$ is a sorted sequence, user $i_n(b)$ ($b \in \{1, 2, \dots, m_n\}\setminus\{l\}$) can decode the signals of user $i_n(l)$ if and only if $b < l$. However, in MIMO-NOMA system, users' channel gains depend not only on the physical environments but on the beams, i.e., the decoding priority may not be fixed and is subject to the beam-forming scheme. Specifically, user $i_n(b)$ can decode the signals of user $i_n(l)$ when the following condition satisfies:

$$\mathbf{w_n}^{H}\mathbf{h}_{\mathbf{i_n}(\mathbf{b})}\mathbf{h}_{\mathbf{i_n}(\mathbf{b})}^{H}\mathbf{w_n} - \mathbf{w_n}^{H}\mathbf{h}_{\mathbf{i_n}(1)}\mathbf{h}_{\mathbf{i_n}(1)}^{H}\mathbf{w_n} > 0 \tag{7}$$

Obviously, beam-forming affects the decoding order by adjusting users' effective channel gains. Accordingly, we introduce a decoding indicator $\lambda_n^{b,l}$ to depict whether or not user $i_n(b)$ can decode the signals of user $i_n(l)$, given by

$$\lambda_n^{b,l} = \frac{1}{2}(1 + \text{sgn}(\mathbf{w_n}^{H}\mathbf{h}_{\mathbf{i_n}(\mathbf{b})}\mathbf{h}_{\mathbf{i_n}(\mathbf{b})}^{H}\mathbf{w_n} - \mathbf{w_n}^{H}\mathbf{h}_{\mathbf{i_n}(1)}\mathbf{h}_{\mathbf{i_n}(1)}^{H}\mathbf{w_n})) \tag{8}$$

Here, the sign function is introduced to denote the decoding priority, which returns 1 when its input is positive, and $-1$ otherwise. From (8), if user $i_n(b)$ can decode the signals of user $i_n(l)$, $\lambda_n^{b,l} = 1$; otherwise, $\lambda_n^{b,l} = 0$. When $\lambda_n^{b,l} = 1$, there is an implicit power constraint, given by

$$p_{i_n(l)} - p_{i_n(b)} > 0 \tag{9}$$

From (9), when $\lambda_n^{b,l} = 1$, the power of user $i_n(l)$ should be larger than that of user $i_n(b)$ to make $i_n(l)$ more easily detected. The received SINR at user $i_n(l)$ can be expressed as

$$SINR_{i_n(l)} = \frac{p_{i_n(l)}g_{i_n(l)}}{1 + \sum_{b\in\{1,2,\dots,m_n\}\setminus\{l\}}\lambda_n^{b,l}p_{i_n(b)}g_{i_n(l)}} \tag{10}$$

where $g_{i_n(l)} = \frac{\mathbf{w_n}^{H}\mathbf{h}_{\mathbf{i_n}(1)}\mathbf{h}_{\mathbf{i_n}(1)}^{H}\mathbf{w_n}}{\left|\bar{\omega}_{i_n(l)}\right|^2}$ denotes the normalized channel gain of user $i_n(l)$. Based on the discussion above, the considered problem can be mathematically expressed as below:

$$\max_{\mathbf{P,W,I}} : \sum_{n=1}^{S}\sum_{l=1}^{m_n}\log(1 + \frac{p_{i_n(l)}g_{i_n(l)}}{1 + \sum_{b\in\{1,2,\dots,m_n\}\setminus\{l\}}\lambda_n^{b,l}p_{i_n(b)}g_{i_n(l)}}) \tag{11a}$$

$$s.t. \frac{p_{i_n(l)}g_{i_n(b)}}{1 + \sum_{j\in\{1,2,\dots,m_n\}\setminus\{l\}}\lambda_n^{j,l}p_{i_n(j)}g_{i_n(b)}} > \Gamma, \lambda_n^{b,l} = 1 \tag{11b}$$

$$p_{i_n(l)} - p_{i_n(b)} > 0, \lambda_n^{b,l} = 1 \tag{11c}$$

$$|[\mathbf{w_n}]_c| = \frac{1}{\sqrt{N}}, \forall n, \forall c = 1, 2, \dots, N \tag{11d}$$

$$\sum_{n=1}^{S} \sum_{l=1}^{m_n} p_{i_n(l)} \le P_{tot} \tag{11e}$$

$$m_n \le M, \forall n \tag{11f}$$

where $\mathbf{P} = \{p_k, k = 1, 2, \dots, K\}$, $\mathbf{W} = \{\mathbf{w_n}, \forall n\}$ and $\mathbf{I} = \{\mathbf{U_n}, \forall n\}$ denote the power allocation, beam-forming and user clustering schemes for MIMO-NOMA, respectively. (11b) denotes the constraint on SIC, where $\Gamma$ denotes the SINR threshold for a successful decoding. (11c) denotes the implicit power constraint. (11d) denotes the CM constraint, where $[\mathbf{w_n}]_c$ denotes the $c$-th element in $\mathbf{w_n}$. Constraint (11e) provides power budget $P_{tot}$ for the considered system. Constraint (11f) indicates that each NOMA cluster can serve at most $M$ users.

## 3. Beam-Forming Algorithm for a Given User Cluster

Problem (11) considers the joint optimization of user clustering, beam-forming and power allocation for MIMO-NOMA, which is challenging to be solved in a polynomial time. Due to the orthogonality among different clusters, in this section, we first consider the beam-forming problem for a given user cluster (the optimal user clustering and power allocation schemes will be further discussed in Sections 4 and 5, respectively). Without loss of generality, we assign the first $m$ users of $\mathbf{U}$ to cluster $n$, i.e., $i_n(l) = l, \forall l = 1, 2, \dots, m$. The beam-forming problem for $n$ can be mathematically expressed as below:

$$\max_{\mathbf{w_n}, \mathbf{P_n}} : \sum_{l=1}^{m} \log(1 + \frac{p_l g_l}{1 + \sum_{b=\{1,2,\dots,m\}\setminus\{l\}} \lambda_n^{b,l} p_b g_l}) \tag{12a}$$

$$s.t. \lambda_n^{b,l} = \frac{1}{2}(1 + \text{sgn}(\mathbf{w_n}^H \mathbf{h_b} \mathbf{h_b}^H \mathbf{w_n} - \mathbf{w_n}^H \mathbf{h_l} \mathbf{h_l}^H \mathbf{w_n})), \forall b, l \in \{1, 2, \dots, m\}, b \ne l \tag{12b}$$

$$p_l - p_b > 0, \lambda_n^{b,l} = 1 \tag{12c}$$

$$\sum_{l=1}^{m} p_l \le \frac{P_{tot}}{S} \tag{12d}$$

$$|[\mathbf{w_n}]_c| = \frac{1}{\sqrt{N}}, \forall c = 1, 2, \dots, N \tag{12e}$$

where $\mathbf{P_n} = \{p_l, \forall l = 1, 2, \dots, m\}$ denotes the power allocation scheme for the $m$ considered users. For the sake of simplify, the SINR constraint is omitted here and will be further considered in Section 5. Each cluster is assumed to have the same power budget, denoted by $\frac{P_{tot}}{S}$. Due to (12e), the beam vector can be represented as $\mathbf{w_n} = \frac{1}{\sqrt{N}}(e^{j\phi_1}, e^{j\phi_2}, \dots, e^{j\phi_N})^T$, where $\phi_c$ denotes the phase of the $c$-th element in $\mathbf{w_n}$. The beam vector is obtained once the phases of its elements are determined. Inspired by this observation, we treat $\mathbf{\Phi} = [\phi_1, \phi_2, \dots, \phi_N]$ as variables. Based on perfect square formula, users' normalized channel gains can be expressed in terms of $\mathbf{\Phi}$ (for the details of derivation, see Appendix A).

$$g_l = \frac{1}{|\bar{\omega}_l|^2 N} \|\mathbf{h_l}\|_2^2 + \frac{2}{|\bar{\omega}_l|^2 N} \sum_{c=1}^{N} \sum_{d=c+1}^{N} \kappa_{l,c} \kappa_{l,d} \cos(\phi_c - \phi_d - (\varphi_{l,c} - \varphi_{l,d})) \tag{13}$$

where $\kappa_{l,c}$ and $\varphi_{l,c}$ denote the amplitude and phase of the $c$-th element in $\mathbf{h_l}$, respectively. However, problem (12) is still difficult to solve due to (12c) and (12d). To predigest the scope of (12), we first produce a feasible solution for $\mathbf{P_n}$ and then maximize (12a) by optimizing $\mathbf{\Phi}$.

Specifically, the power allocation scheme can be obtained based on FTPC, i.e., the transmit power of user $l$ can be represented by:

$$p_l = \frac{P_{tot}}{S} \frac{g_l^{-\gamma}}{\sum\limits_{j=1}^{m} g_j^{-\gamma}} \tag{14}$$

where $\gamma$ denotes the decay factor. With FTPC, constraint (12d) always holds since $\sum\limits_{l} p_l = \frac{P_{tot}}{S}$. Moreover, $\gamma$ determines the correlation between users' channel gains and transmitting power. When $\gamma = 0$, transmitting power is totally unrelated to normalized gains, i.e., each user has the same transmitting power. Moreover, $p_l$ and $g_l$ will be negatively correlated as $\gamma$ increases, which is consistent with (12c) and thus makes (14) a feasible solution.

Accordingly, problem (12) can be transformed into

$$\max_{\boldsymbol{\Phi}} : \sum_{l=1}^{m} \log(1 + \frac{\frac{P_{tot}}{S} g_l^{1-\gamma}}{\sum\limits_{j=1}^{m} g_j^{-\gamma} + \frac{P_{tot}}{S} \sum_{b \in \{1,2,\dots,m\} \backslash \{l\}} \lambda_n^{b,l} g_b^{-\gamma} g_l}) \tag{15}$$

However, it is still challenging for us to solve (15) since the sign function in (12b) is non-differentiable. To solve this problem, we produce an approximation of $\lambda_n^{b,l}$, given by

$$\bar{\lambda}_n^{b,l} = \frac{1}{1 + \exp(g_l - g_b)} \tag{16}$$

The Sigmoid function is introduced which is first-order differentiable. From (16), when $g_b - g_l \to +\infty$, $\bar{\lambda}_n^{b,l} \to 1$; when $g_b - g_l \to -\infty$, $\bar{\lambda}_n^{b,l} \to 0$. Since the output of (16) ranges from zero to one, we consider (16) as the probability that user $b$ successfully decodes the signals of user $l$.

Then, (15) can be re-written as

$$\max_{\boldsymbol{\Phi}} : f(\boldsymbol{\Phi}) = \sum_{l=1}^{m} \log(1 + \frac{\frac{P_{tot}}{S} g_l^{1-\gamma}}{\sum\limits_{j=1}^{m} g_j^{-\gamma} + \frac{P_{tot}}{S} \sum_{b \in \{1,2,\dots,m\} \backslash \{l\}} \frac{g_b^{-\gamma} g_l}{1+\exp(g_l - g_b)}}) \tag{17}$$

Consider the partial derivatives:

$$\frac{\partial f}{\partial \phi_c} = \sum_{l=1}^{m} \frac{\partial f}{\partial g_l} \frac{\partial g_l}{\partial \phi_c} \tag{18}$$

$$\frac{\partial g_l}{\partial \phi_c} = \frac{2}{\bar{\omega}_l N} \sum_{d=1}^{N} \kappa_{l,d} \kappa_{l,c} \sin(\phi_d - \varphi_{l,d} - (\phi_c - \varphi_{l,c})) \tag{19}$$

Since problem (17) is a differentiable non-constrained problem, a quasi-Newton method named L-BFGS can be applied to solve it in limited iterations. In each iteration, L-BFGS produces an updating direction for $\boldsymbol{\Phi}$ based on the information from the last $T$ iterations. Once the update direction is determined, the Armijo rule is applied to obtain a proper step size. More details about the proposed algorithm are as shown in Algorithm 1.

---

**Algorithm 1** Beam-forming Algorithm for A Given Cluster

---

**Require: $\mathbf{U_n}$**
**Ensure: $\boldsymbol{\Phi}$**
 1: Initialize $T$, $\eta$.
 2: $\mathbf{Y_n} = \varnothing$, $\mathbf{S_n} = \varnothing$, $\mathbf{R_n} = \varnothing$.
 3: Randomly initialize $\boldsymbol{\Phi}$.
 4: $\mathbf{g_{pre}} \leftarrow$ the gradient of $f$ at $\boldsymbol{\Phi}$.
 5: Calculate the updating direction $\mathbf{y} = -\mathbf{g_{pre}}$.
 6: Obtain the optimal step size $\mu$ based on the Armijo rule.
 7: $\mathbf{y} \leftarrow \mu\mathbf{y}$, $\boldsymbol{\Phi} \leftarrow \boldsymbol{\Phi} + \mathbf{y}$.
 8: $\mathbf{g_{cur}} \leftarrow$ the gradient of $f$ at $\boldsymbol{\Phi}$.
 9: $\mathbf{s} \leftarrow \mathbf{g_{cur}} - \mathbf{g_{pre}}$.
10: $\rho \leftarrow \mathbf{y}^H \mathbf{s}$.
11: **while** $|\mathbf{g_{cur}}| \geq \eta$ **do**
12:     $\mathbf{g_{pre}} \leftarrow \mathbf{g_{cur}}$.
13:     Insert $\mathbf{y}$ to $\mathbf{Y_n}$.
14:     Insert $\mathbf{s}$ to $\mathbf{S_n}$.
15:     Insert $\rho$ to $\mathbf{R_n}$.
16:     $L \leftarrow$ the number of the elements in $\mathbf{Y_n}$.
17:     **if** $L > T$ **then**
18:       Pop the first element in $\mathbf{Y_n}$.
19:       Pop the first element in $\mathbf{S_n}$.
20:       Pop the first element in $\mathbf{R_n}$.
21:       $L \leftarrow L - 1$.
22:     **end if**
23:     (Back Propagating)
24:     **for** $i = L:-1:1$ **do**
25:       $\mathbf{s} \leftarrow$ the $i$-th element in $\mathbf{S_n}$.
26:       $\mathbf{y} \leftarrow$ the $i$-th element in $\mathbf{Y_n}$.
27:       $\rho \leftarrow$ the $i$-th element in $\mathbf{R_n}$.
28:       $\chi_i \leftarrow \rho\mathbf{s}^H\mathbf{g_{cur}}$.
29:       $\mathbf{g_{cur}} \leftarrow \mathbf{g_{cur}} - \chi_i\mathbf{y}$.
30:     **end for**
31:     (Forward Propagating)
32:     $\mathbf{res} \leftarrow \mathbf{g_{cur}}$.
33:     **for** $i = 1:L$ **do**
34:       $\mathbf{s} \leftarrow$ the $i$-th element in $\mathbf{S_n}$.
35:       $\mathbf{y} \leftarrow$ the $i$-th element in $\mathbf{Y_n}$.
36:       $\rho \leftarrow$ the $i$-th element in $\mathbf{R_n}$.
37:       $\beta_i \leftarrow \rho\mathbf{y}^H\mathbf{res}$.
38:       $\mathbf{res} \leftarrow \mathbf{res} + (\chi_i - \beta_i)\mathbf{s}$.
39:     **end for**
40:     $\mathbf{y} = -\mathbf{res}$.
41:     steps (6)–(10)
42: **end while**

---

## 4. User Clustering for MIMO-NOMA System

In this section, optimal user clustering is further considered based on Algorithm 1. In each cluster, the channel differences among multiplexed users should be large enough to perform SIC successfully. Moreover, since the users in the same cluster are served by the same beam vector, their channel correlations should also be emphasized to bring the advantages of MIMO into full play. Accordingly, the optimal clustering scheme will be obtained with consideration for both the two factors.

Due to SIC, the strongest user in each cluster is in fact served by OMA, which can achieve good performance with less power when its channel gain is relatively large. Therefore, the first $S$ users of $\mathbf{U}$ will be assigned to $S$ different clusters, respectively. Due to the high channel gains, these users could achieve good performances with less power, which can in turn enable more power budget for others. After initializing $S$ clusters, the remaining users in $\mathbf{U}$ will successively select a suitable cluster to join. For each user $k$, we define a utility function to assess its preference on different clusters, given by

$$u_k(n) = \frac{1}{m_n} \sum_{l=1}^{m_n} \frac{\left| \mathbf{h}_\mathbf{k}^H \mathbf{h}_{\mathbf{i_n}(1)} \right|}{\left| \mathbf{h_k} \right| \left| \mathbf{h}_{\mathbf{i_n}(1)} \right|} - \theta \frac{\left( \left| \mathbf{h_k} \right| + \sum_{l=1}^{m_n} \left| \mathbf{h}_{\mathbf{i_n}(1)} \right| \right)^2}{(m_n + 1)\left( \left| \mathbf{h_k} \right|^2 + \sum_{l=1}^{m_n} \left| \mathbf{h}_{\mathbf{i_n}(1)} \right|^2 \right)} \tag{20}$$

where $u_k(n)$ describes user $k$'s preference for cluster $n$. The first term depicts the channel correlations between user $k$ and the existing users in cluster $n$. The second term is the Jain's fairness index, which measures the channel difference between user $k$ and the existing users in cluster $n$. Specifically, the second term ranges from $\frac{1}{m_n+1}$ to 1 and will decrease as the channel difference gets larger. $\theta$ denotes the relative weight for the two aspects. Based on the utility function, user $k$ will be further assigned to its favorite cluster $n_k$, given by

$$n_k = \underset{n \in \{1,2,\dots,S\}}{\arg\max} (u_k(n)) \tag{21}$$

Clusters will reject users only when condition (11f) is not met. For each cluster $n$, when $m_n = M + 1$, $n$ should reject a user to satisfy the size-constraint. Accordingly, we can produce multiple possible user set for cluster $n$ by removing any single user from $\mathbf{U_n}$. Based on Algorithm 1, beam-forming can be done for each possible user set to compare their performances, and the optimal user set for cluster $n$ is obtained accordingly.

One can observe that the relative weight is of great significance in steering the ultimate clustering scheme. When $\theta$ is relatively small, channel correlations play a decisive role in the clustering process. As $\theta$ increases, channel differences, in turn, become the controlling factor of the ultimate clustering scheme. With any given $\theta$, user clustering can be performed based on Algorithm 2. Then, Algorithm 1 can be applied to obtain a beam-forming scheme, and the corresponding achievable SE can be denoted by $w(\theta)$.

---

**Algorithm 2** User Clustering Scheme for MIMO-NOMA with A Given Relative Weight

---

1: Assign the first $S$ users of $\mathbf{U}$ to $S$ different clusters.
2: Construct the utility function as (20) based on the given relative weight.
3: **for** $j = S + 1{:}K$ **do**
4:     Sort multiple clusters based on user $j$'s preference.
5:     Denote the sorted sequence by $\mathbf{\Omega_j}$.
6:     **while** $\mathbf{\Omega_j} \neq \varnothing$ **do**
7:         $n_j \leftarrow$ the first cluster in $\mathbf{\Omega_j}$.
8:         Insert user $j$ to $\mathbf{U_{n_j}}$.
9:         $NUM \leftarrow$ the number of the users in $\mathbf{U_{n_j}}$.
10:        **if** $NUM \leq M$ **then**
11:            Break.
12:        **else**
13:            **for** $i = 1{:}M + 1$ **do**
14:                Remove the $i$-th user from $\mathbf{U_{n_j}}$.
15:                Obtain the optimal beam vector for cluster $n_j$ by Algorithm 1.
16:                $\varepsilon_i \leftarrow$ the sum rate of the users in $\mathbf{U_{n_j}}$.
17:                Insert the removed user to its original position.
18:            **end for**
19:            $\tilde{i} \leftarrow$ the position of the maximum in $[\varepsilon_1, \varepsilon_2, \dots, \varepsilon_{M+1}]$.
20:            $x \leftarrow$ the index of the $\tilde{i}$-th user in $\mathbf{U_{n_j}}$.
21:            Remove the $\tilde{i}$-th user from $\mathbf{U_{n_j}}$.
22:        **end if**
23:        **if** $x = j$ **then**
24:            Remove cluster $n_j$ from $\mathbf{\Omega_j}$.
25:        **else**
26:            Break.
27:        **end if**
28:    **end while**
29: **end for**

---

In this section, the optimal $\theta$ is obtained by PSO. In PSO, the optimal $\theta$ can be obtained through numerous iterations. In each iteration, PSO produces multiple possible solutions for $\theta$, each corresponding to a possible clustering scheme. Based on Algorithm 1, beam-forming can be done for each possible clustering scheme to compare their performance. We will further select the one with the maximum SE performance as the optimal clustering scheme, and its corresponding relative weight is exactly the optimal $\theta$ obtained by PSO. More details are as described in Algorithm 3. Note that the random variable $\delta$ in step (13) denotes the step size, which is a real number ranging from 0 to 1.

---

**Algorithm 3** PSO-based Optimal User Clustering

---

1: Initialize group size $G$.
2: Initialize the total number of iterations $D$.
3: **for** $g = 1:G$ **do**
4:   Initialize position $\theta_g$ for particle $g$.
5:   $e_{best,g} = w(\theta_g)$.
6:   $p_{best,g} = \theta_g$.
7: **end for**
8: $\bar{g} \leftarrow$ the position of the maximum in $\{e_{best,g}, \forall g\}$.
9: $G_{best} = p_{best,\bar{g}}$.
10: $t = 1$.
11: **while** $t \leq D$ **do**
12:   **for** $g = 1:G$ **do**
13:     $\theta_g = 0.5(p_{best,g} + G_{best}) + \delta(p_{best,g} - G_{best})$.
14:     $e_g = w(\theta_g)$.
15:     **if** $e_g > e_{best,g}$ **then**
16:       $p_{best,g} = \theta_g$.
17:       $e_{best,g} = e_g$.
18:     **end if**
19:   **end for**
20:   $\bar{g} \leftarrow$ the position of the maximum in $\{e_{best,g}, \forall g\}$.
21:   $G_{best} = p_{best,\bar{g}}$.
22:   $t = t + 1$.
23: **end while**
24: $\theta = G_{best}$.
25: Perform user clustereing with the obtained $\theta$ based on Algorithm 2.

---

## 5. Power Allocation for MIMO-NOMA

User clustering and beam-forming are jointly solved in Section 4. However, FTPC is still applied for intra-cluster power allocation, which needs further improvements. In this section, power allocation is optimized based on the optimal user clustering and beam-forming schemes. Without loss of generality, the users in $\mathbf{U_n}$ are re-ordered in the descending order of effective channel gains. The re-ordered sequence can be denoted by $\tilde{\mathbf{U}}_{\mathbf{n}} = \{\tilde{i}_n(1), \tilde{i}_n(2), \ldots, \tilde{i}_n(m_n)\}$, where $\tilde{i}_n(l)$ ($l \in \{1, 2, \ldots, m_n\}$) denotes the index of the $l$-th user in $\tilde{\mathbf{U}}_{\mathbf{n}}$. Moreover, we have $\left|\mathbf{w_n}^H \mathbf{h}_{\tilde{\mathbf{i}}_{\mathbf{n}}(1)}\right|^2 > \left|\mathbf{w_n}^H \mathbf{h}_{\tilde{\mathbf{i}}_{\mathbf{n}}(2)}\right|^2 > \ldots > \left|\mathbf{w_n}^H \mathbf{h}_{\tilde{\mathbf{i}}_{\mathbf{n}}(m_n)}\right|^2$, where $\mathbf{h}_{\tilde{\mathbf{i}}_{\mathbf{n}}(1)}$ denotes the channel vector of user $\tilde{i}_n(l)$. The achievable rate of user $\tilde{i}_n(l)$ with normalized bandwidth can be represented as:

$$
\begin{aligned}
R_{\tilde{i}_n(l)} &= \log(1 + \frac{p_{\tilde{i}_n(l)} g_{\tilde{i}_n(l)}}{1 + \sum\limits_{b=1}^{l-1} p_{\tilde{i}_n(b)} g_{\tilde{i}_n(l)}}) \\
&= \log(1 + \sum_{b=1}^{l} p_{\tilde{i}_n(b)} g_{\tilde{i}_n(l)}) - \log(1 + \sum_{b=1}^{l-1} p_{\tilde{i}_n(b)} g_{\tilde{i}_n(l)})
\end{aligned}
\tag{22}
$$

$$g_{\tilde{i}_n(l)} = \frac{\left|\mathbf{w_n}^H \mathbf{h}_{\tilde{\mathbf{i}}_{\mathbf{n}}(l)}\right|^2}{\left|\omega_{\tilde{\mathbf{i}}_{\mathbf{n}}(l)}\right|^2} \tag{23}$$

The power allocation problem can be mathematically expressed as below:

$$\max_{\mathbf{P}} : \sum_{n=1}^{S}\sum_{l=1}^{m_n} R_{\tilde{i}_n(l)} \tag{24a}$$

$$s.t. \frac{p_{\tilde{i}_n(l)}g_{\tilde{i}_n(b)}}{1 + \sum\limits_{j=1}^{l-1} p_{\tilde{i}_n(j)}g_{\tilde{i}_n(b)}} > \Gamma, \forall b, l \in \{1, 2, \ldots, m_n\}, b < l \tag{24b}$$

$$\sum_{n=1}^{S}\sum_{l=1}^{m_n} p_{\tilde{i}_n(l)} \leq P_{tot} \tag{24c}$$

As mentioned in Section 2, $\mathbf{P} = \{p_k, k = 1, 2, \ldots, K\}$ denotes the power allocation scheme. (24b) denotes the SINR constraint for decoding. In each cluster $n$, the signals of user $\tilde{i}_n(l)$ should be decoded from the others with higher channel gains. The series SINR constraints can be represented as below:

$$\frac{p_{\tilde{i}_n(l)}g_{\tilde{i}_n(l-1)}}{1 + \sum\limits_{j=1}^{l-1} p_{\tilde{i}_n(j)}g_{\tilde{i}_n(l-1)}} > \Gamma \tag{25}$$

$$\frac{p_{\tilde{i}_n(l)}g_{\tilde{i}_n(l-2)}}{1 + \sum\limits_{j=1}^{l-1} p_{\tilde{i}_n(j)}g_{\tilde{i}_n(l-2)}} > \Gamma \tag{26}$$

$$\vdots$$

$$\frac{p_{\tilde{i}_n(l)}g_{\tilde{i}_n(1)}}{1 + \sum\limits_{j=1}^{l-1} p_{\tilde{i}_n(j)}g_{\tilde{i}_n(1)}} > \Gamma \tag{27}$$

Since $g_{\tilde{i}_n(1)} > g_{\tilde{i}_n(2)} > \ldots > g_{\tilde{i}_n(m_n)}$, (25)–(27) will all hold when (25) holds. The considered problem can be further simplified as:

$$\max_{\mathbf{P}} : \sum_{n=1}^{S}\sum_{l=1}^{m_n} R_{\tilde{i}_n(l)} \tag{28a}$$

$$s.t. \frac{p_{\tilde{i}_n(l)}g_{\tilde{i}_n(l-1)}}{1 + \sum\limits_{j=1}^{l-1} p_{\tilde{i}_n(j)}g_{\tilde{i}_n(l-1)}} > \Gamma, \forall l = 2, 3, \ldots, m_n \tag{28b}$$

$$\sum_{n=1}^{S}\sum_{l=1}^{m_n} p_{\tilde{i}_n(l)} \leq P_{tot} \tag{28c}$$

The power budget of each cluster can be auto-adjusted based on the channel characteristics of intra-cluster users. To solve (28), we first introduce an auxiliary variable $t$ to bound (28a) from below and then optimize (28) by maximizing $t$. The equivalence problem is given by

$$\max_{t,\mathbf{P}} : t \tag{29a}$$

$$s.t. \sum_{n=1}^{S} \sum_{l=1}^{m_n} \log(1 + \sum_{b=1}^{l} p_{\tilde{i}_n(b)} g_{\tilde{i}_n(l)}) - \log(1 + \sum_{b=1}^{l-1} p_{\tilde{i}_n(b)} g_{\tilde{i}_n(l)}) > t \tag{29b}$$

$$p_{\tilde{i}_n(l)} g_{\tilde{i}_n(l-1)} - \Gamma(1 + \sum_{j=1}^{l-1} p_{\tilde{i}_n(j)} g_{\tilde{i}_n(l-1)}) > 0, \forall l = 2, 3, \ldots, m_n \tag{29c}$$

$$\sum_{n=1}^{S} \sum_{l=1}^{m_n} p_{\tilde{i}_n(l)} \leq P_{tot} \tag{29d}$$

However, problem (29) is non-convex since (29b) is a non-convex constraint. To address this issue, we produce a convex relaxation of (29b) based on SCA. Accordingly, (29) is transformed into a convex problem, which can be efficiently solved with a polynomial time. To relax (29b), we first consider a DC function, given by

$$\xi_{n,l} = \log(1 + \sum_{b=1}^{l} p_{\tilde{i}_n(b)} g_{\tilde{i}_n(l)}) - \log(1 + \sum_{b=1}^{l-1} p_{\tilde{i}_n(b)} g_{\tilde{i}_n(l)}) \tag{30}$$

The first and second term in (30) are both logarithmic functions, which makes (30) a DC function. Due to the concavity of logarithmic functions, the second term in (30) can be tightly bounded from above with its first-order Taylor expansion, i.e., with any given $\{\bar{p}_{\tilde{i}_n(1)}, \bar{p}_{\tilde{i}_n(2)}, \ldots, \bar{p}_{\tilde{i}_n(l-1)}\}$, we have

$$\log(1 + \sum_{b=1}^{l-1} p_{\tilde{i}_n(b)} g_{\tilde{i}_n(l)}) < \log(1 + \sum_{b=1}^{l-1} \bar{p}_{\tilde{i}_n(b)} g_{\tilde{i}_n(l)}) \\ + \frac{g_{\tilde{i}_n(l)}}{1 + \sum\limits_{b=1}^{l-1} \bar{p}_{\tilde{i}_n(b)} g_{\tilde{i}_n(l)}} \sum_{b=1}^{l-1} (p_{\tilde{i}_n(b)} - \bar{p}_{\tilde{i}_n(b)}) \tag{31}$$

Substitute (31) into (30), we obtain

$$\xi_{n,l} > \log(1 + \sum_{b=1}^{l} p_{\tilde{i}_n(b)} g_{\tilde{i}_n(l)}) - \log(1 + \sum_{b=1}^{l-1} \bar{p}_{\tilde{i}_n(b)} g_{\tilde{i}_n(l)}) \\ - \frac{g_{\tilde{i}_n(l)}}{1 + \sum\limits_{b=1}^{l-1} \bar{p}_{\tilde{i}_n(b)} g_{\tilde{i}_n(l)}} \sum_{b=1}^{l-1} (p_{\tilde{i}_n(b)} - \bar{p}_{\tilde{i}_n(b)}) \tag{32}$$

The left-hand side (LHS) of (29b) can be represented as $\sum\limits_{n=1}^{S} \sum\limits_{l=1}^{m_n} \xi_{n,l}$. Accordingly, with any given $\{\bar{p}_k, k = 1, 2, \ldots, K\}$, we can derive a lower bound $B$ for the LHS of (29b), represented as (33) from the top of next page. Obviously, $B$ is convex in **P**, and (29b) can be further relaxed by restricting $B$ to be greater than $t$. The equivalence convex problem is given by (34)–(37).

$$\sum_{n=1}^{S} \sum_{l=1}^{m_n} \log(1 + \sum_{b=1}^{l} p_{\tilde{i}_n(b)} g_{\tilde{i}_n(l)}) - \log(1 + \sum_{b=1}^{l-1} p_{\tilde{i}_n(b)} g_{\tilde{i}_n(l)}) > \underbrace{\sum_{n=1}^{S} \sum_{l=1}^{m_n} (\log(1 + \sum_{b=1}^{l} p_{\tilde{i}_n(b)} g_{\tilde{i}_n(l)}) - \log(1 + \sum_{b=1}^{l-1} \bar{p}_{\tilde{i}_n(b)} g_{\tilde{i}_n(l)}) - \frac{g_{\tilde{i}_n(l)}}{1 + \sum\limits_{b=1}^{l-1} \bar{p}_{\tilde{i}_n(b)} g_{\tilde{i}_n(l)}} \sum_{b=1}^{l-1} (p_{\tilde{i}_n(b)} - \bar{p}_{\tilde{i}_n(b)}))}_{B} \tag{33}$$

$$\max_{t, \mathbf{P}} : t \tag{34}$$

$$s.t. B > t \tag{35}$$

$$p_{\tilde{i}_n(l)} g_{\tilde{i}_n(l-1)} - \Gamma(1 + \sum_{j=1}^{l-1} p_{\tilde{i}_n(j)} g_{\tilde{i}_n(l-1)}) > 0, \forall l = 2, 3, \ldots, m_n \tag{36}$$

$$P_{tot} - \sum_{n=1}^{S} \sum_{l=1}^{m_n} p_{\tilde{i}_n(l)} \geq 0 \tag{37}$$

According to the principles of SCA, the solution of problem (29) should be obtained through multiple iterations. In each iteration, we produce an equivalence problem of (29) as (34)–(37), which is further solved by some effective optimization tools. Specifically, in the *i*-th iteration, $\{\bar{p}_k, k = 1, 2, \ldots, K\}$ should be set as the solution obtained in the previous iteration. More details are as described in Algorithm 4.

---

**Algorithm 4** Power Allocation Scheme for MIMO-NOMA

---

1: $\bar{p}_{\tilde{i}_n(l)} = \dfrac{P_{tot}}{S} \dfrac{g_{\tilde{i}_n(l)}^{-\gamma}}{\sum\limits_{j=1}^{m_n} g_{\tilde{i}_n(j)}^{-\gamma}}, \forall n, l.$

2: Initialize $\eta$.

3: $r_0 \leftarrow \sum\limits_{n=1}^{S} \sum\limits_{l=1}^{m_n} \log(1 + \dfrac{\bar{p}_{\tilde{i}_n(l)} g_{\tilde{i}_n(l)}}{1 + \sum\limits_{b=1}^{l-1} \bar{p}_{\tilde{i}_n(b)} g_{\tilde{i}_n(l)}})$ .

4: $i = 0$.

5: **repeat**

6:    $i = i + 1$.

7:    Produce an equivalence problem of (29) as (34)–(37) based on the Taylor expansion at $\{\bar{p}_k, \forall k\}$.

8:    Solve the obtained convex problem to get $\{\tilde{p}_k, \forall k\}$.

9:    $r_i \leftarrow \sum\limits_{n=1}^{S} \sum\limits_{l=1}^{m_n} \log(1 + \dfrac{\tilde{p}_{\tilde{i}_n(l)} g_{\tilde{i}_n(l)}}{1 + \sum\limits_{b=1}^{l-1} \tilde{p}_{\tilde{i}_n(b)} g_{\tilde{i}_n(l)}})$.

10:    $\bar{p}_k \leftarrow \tilde{p}_k, \forall k$.

11: **until** $\left| \dfrac{r_i - r_{i-1}}{r_{i-1}} \right| < \eta$

12: $p_k = \bar{p}_k, \forall k$.

---

## 6. Simulations Results

In this section, the performance of the proposed scheme is evaluated by multiple simulations. The distance from users to the BS is uniformly distributed in the range of 0 to 500 m. The channel vector of each user is assumed to be the product of large-scale path loss and Rayleigh fading. We also evaluate the performances of two existing schemes to illustrate the significance of the proposed scheme [9,21]. Some key parameters are as summarized in Table 1. The effects of multiple factors will be discussed in more details.

**Table 1.** Simulation Parameters.

| Parameter | Value |
|-----------|-------|
| number of antennas ($N$) | 2 |
| noise power spectral density ($N_0$) | $-169$ dBm |
| system bandwidth ($B$) | 360 kHz |
| large-scale path loss model | free-space path loss model |
| $T$ | 20 |
| $\eta$ | $10^{-3}$ |
| $G$ | 30 |
| $P_{tot}$ | 30 dBm |
| $D$ | 100 |

Figure 1 plots SE versus *M* with $\gamma = 0.2$ and $S = 2$. As shown in Figure 1, SE increases with *M* since a larger *M* allows a cluster to serve more users. However, such effect gets saturated as *M* increases to a certain degree, subject to the total power budget.

The effect of user diversity is also considered. Figure 2 plots SE versus *K* with $\gamma = 0.2$ and $S = 2$. From Figure 2, one can observe that user diversity plays a key role in increasing SE. Moreover, the proposed scheme outperforms the two existing schemes in terms of SE.



**Figure 1.** Spectrum efficiency (SE) versus *M*.



**Figure 2.** SE versus *K*.

Next, the impact of the decay parameter $\gamma$ is further considered. In Algorithm 1, $\gamma$ determines the correlation between users' channel gains and transmitting power. As $\gamma$ increases, weak users have made notable gains at the expense of strong users. However, the achievable SE mainly depends on strong users due to their high channel gains. Figure 3 plots SE versus $\gamma$ with $S = 2$ and $M = 2$. From Figure 3, SE decreases with the increase of $\gamma$ due to the performances degradation of strong users.

**Figure 3.** SE versus $\gamma$.

In Algorithm 2, channel differences and correlations are both involved in the design of user clustering. A relative weight $\theta$ is introduced for the two aspects, which is obtained based on group hunting strategy. As discussed before, $\theta$ is of great significance in steering the ultimate clustering scheme. With $M = 2$, $\gamma = 0.2$ and $S = 2$, Figure 4 plots SE versus $K$ under different $\theta$. When $\theta = 0.1$, the achievable SE is relatively small because of neglect of channel difference characteristics. As $\theta$ increases, we can achieve a better balance between the two contributing factors and SE will increase accordingly. However, when $\theta$ increases to a certain degree, e.g., $\theta \geq 10$, the achievable SE will further decrease for overlooking channel correlations characteristics. Moreover, the performance upper bound is also considered. Exhaustive user search can be done to find the optimal clustering scheme. Based on Algorithms 1 and 4, beam-forming and power allocation can be done for each possible clustering scheme to compare their performance. From Figure 4, the performance of the proposed scheme can approach to the upper bound due to the optimization strategy of PSO.
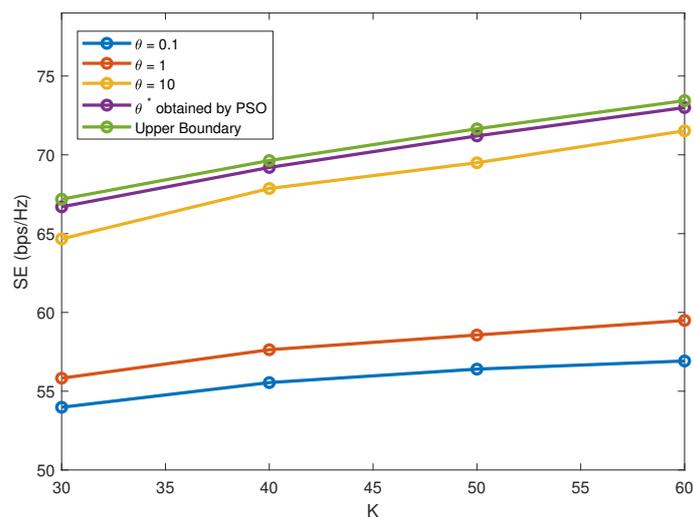


**Figure 4.** SE versus $K$ under different $\theta$.

Figure 5 plots SE versus $S$ with $M = 2$ and $\gamma = 0.2$. The number of users that can be served simultaneously will increase with the increase of $S$. Moreover, when $N \geq S - 1$, there is no interference among different clusters. As shown in Figure 5, SE has a nearly linear increase with $S$ due to the orthogonality among different clusters.
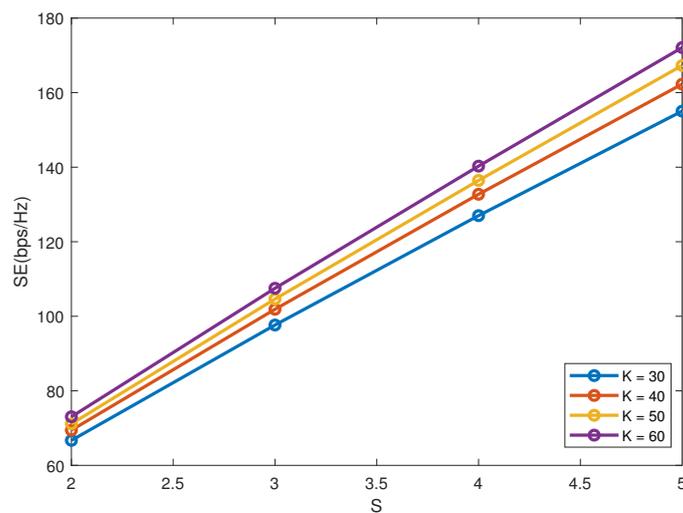
**Figure 5.** SE versus *S*.

The effect of Γ is also considered. In the SIC-based decoding process, the received SINR at users will increase with the increase of Γ, which not only improves system reliability, but also enables more power budget for edge users. Accordingly, as Γ increases, there is a corresponding increase in edge users' data rates. With $\gamma = 0.2$, Figure 6 plots SE versus *K* under different Γ. From Figure 6, the achievable SE is almost unaffected by Γ, i.e., as Γ increases, we sacrifice strong users' data rates in exchange for weak users' rates to ensure all multiplexed users can achieve satisfactory performances. However, the existing schemes fail to obtain a feasible solution when Γ is greater than 0.5. By contrast, the proposed scheme is more robust which can obtain a feasible solution with a larger Γ. One explanation for this is that, the proposed scheme solves beam-forming and power allocation separately, which predigests the scope of the considered problem and helps to achieve a better performance with strong robustness.
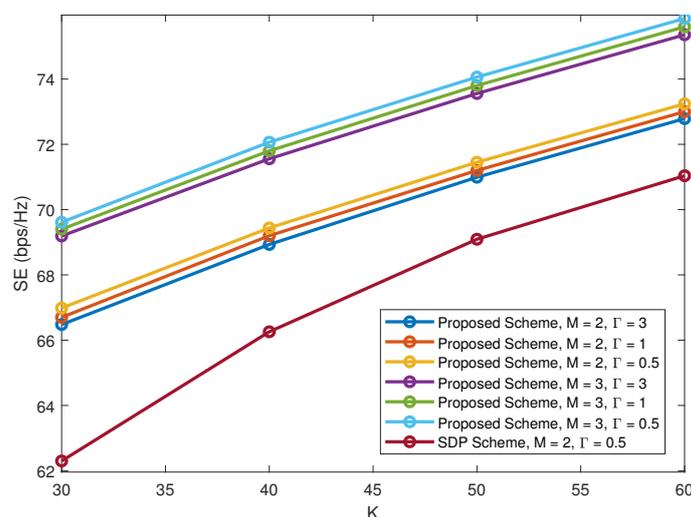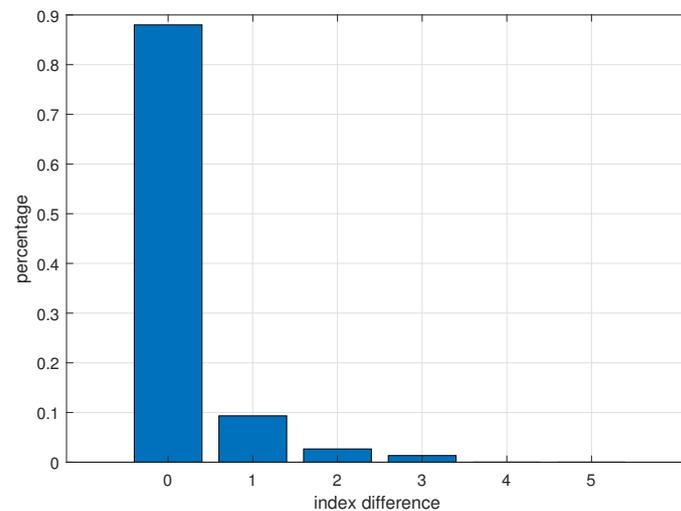


**Figure 6.** SE versus *K* under different Γ.

The effect of beam-forming on the optimal decoding order is also investigated. We generate 1000 instances with $\gamma = 0.2$, $K = 60$, $S = 5$ and $M = 6$. The proposed scheme is applied for the realization of each instance. In each cluster, the intra-cluster users are sorted in the descending order of channel gains and normalized channel gains, respectively. The positions of each user in the two sorted sequences are recorded, and their difference can be utilized to describe how often beam-forming changes the optimal SIC order. Figure 7 plots the distribution of the position differences. From Figure 7, the decoding order

generally remains unchanged. However, there are circumstances where the optimal SIC order is slightly adjusted.



**Figure 7.** Position difference distribution.

*Complexity Analysis*

With any given relative weight, the corresponding clustering scheme can be obtained by Algorithm 2. For each user $k$, Algorithm 2 measures $k'$ s preference for different clusters.

User $k$ will be first assigned to its favorite cluster $n_k$. After assigning user k to cluster $n_k$, two cases can occur:

- Case 1: The number of the users in cluster $n_k$ is no greater than M;
- Case 2: The number of the users in cluster $n_k$ is greater than M.

In case 1, user $k$ can be directly assigned to cluster $n_k$. In case 2, cluster $n_k$ should reject a user to meet the size constraint. Algorithm 2 produces $(M + 1)$ possible user set for cluster $n_k$. Based on Algorithm 1, beam-forming can be done for each possible user set to compare their performance, and the rejected user is obtained accordingly. If the rejected user is user $k$, $k$ will be further assigned to its second-favorite cluster. The above process will be repeatedly executed until either user $k$ is successfully assigned to a single cluster or all the clusters are processed.

Accordingly, Algorithm 2 consists of two parts: part 1 measures each user's preference for different clusters; part 2 helps each user select a suitable cluster to join. The complexity of part 1 is O($SK$), and the complexity of part 2 is O($SK$). The complexity of Algorithm 2 is O($SK$).

The optimal relative weight can be obtained by Algorithm 3 through $D$ iterations. In each iteration, Algorithm 3 produces $G$ possible relative weights, each corresponding to a possible clustering scheme. Based on Algorithm 1, beam-forming can be done for each possible clustering scheme to compare their performance. The complexity of Algorithm 3 is O($SK$).

## 7. Conclusions

In this passage, we consider the multi-dimensional resource allocation problem for MIMO-NOMA, which consists of power allocation, user clustering and beam-forming, respectively. A three-step resource allocation framework is proposed to solve the considered problem: step 1 solves the beam-forming problem for a given user cluster; step 2 obtains the optimal clustering scheme based on the proposed beam-forming algorithm; step 3 further optimizes power allocation based on the optimal user clustering and beam-forming schemes. Simulation results show that the proposed scheme can effectively increase the

received SINR at users. Additionally, the performance of the proposed scheme can approach the performance upper bound in terms of SE.

## Appendix A

The normalized channel gain of user $l$ is given by

$$g_l = \frac{1}{|\bar{\omega}_l|^2}\left|\mathbf{h_l}^H \mathbf{w_n}\right|^2 = \frac{1}{|\bar{\omega}_l|^2}\left|\frac{1}{\sqrt{N}}\sum_{c=1}^{N}\kappa_{l,c}e^{j(\phi_c - \varphi_{l,c})}\right|^2 \tag{A1}$$

$$= \frac{1}{|\bar{\omega}_l|^2}\left|\frac{1}{\sqrt{N}}\sum_{c=1}^{N}\kappa_{l,c}\cos(\phi_c - \varphi_{l,c}) + j\frac{1}{\sqrt{N}}\sum_{c=1}^{N}\kappa_{l,c}\sin(\phi_c - \varphi_{l,c})\right|^2 \tag{A2}$$

$$= \frac{1}{|\bar{\omega}_l|^2 N}(\sum_{c=1}^{N}\kappa_{l,c}\cos(\phi_c - \varphi_{l,c}))^2 + \frac{1}{|\bar{\omega}_l|^2 N}(\sum_{c=1}^{N}\kappa_{l,c}\sin(\phi_c - \varphi_{l,c}))^2 \tag{A3}$$

$$= \frac{1}{|\bar{\omega}_l|^2 N}\sum_{c=1}^{N}\kappa_{l,c}^2\cos^2(\phi_c - \varphi_{l,c}) + \frac{1}{|\bar{\omega}_l|^2 N}\sum_{c=1}^{N}\kappa_{l,c}^2\sin^2(\phi_c - \varphi_{l,c}) \tag{A4}$$

$$+ \frac{2}{|\bar{\omega}_l|^2 N}\sum_{c=1}^{N}\sum_{d=c+1}^{N}\kappa_{l,c}\kappa_{l,d}\cos(\phi_c - \varphi_{l,c})\cos(\phi_d - \varphi_{l,d}) \tag{A5}$$

$$+ \frac{2}{|\bar{\omega}_l|^2 N}\sum_{c=1}^{N}\sum_{d=c+1}^{N}\kappa_{l,c}\kappa_{l,d}\sin(\phi_c - \varphi_{l,c})\sin(\phi_d - \varphi_{l,d}) \tag{A6}$$

$$= \frac{1}{|\bar{\omega}_l|^2 N}\sum_{c=1}^{N}\kappa_{l,c}^2 + \frac{2}{|\bar{\omega}_l|^2 N}\sum_{c=1}^{N}\sum_{d=c+1}^{N}\kappa_{l,c}\kappa_{l,d}\cos(\phi_c - \varphi_{l,c} - (\phi_d - \varphi_{l,d})) \tag{A7}$$

$$= \frac{1}{|\bar{\omega}_l|^2 N}\|\mathbf{h_l}\|_2^2 + \frac{2}{|\bar{\omega}_l|^2 N}\sum_{c=1}^{N}\sum_{d=c+1}^{N}\kappa_{l,c}\kappa_{l,d}\cos(\phi_c - \varphi_{l,c} - (\phi_d - \varphi_{l,d})) \tag{A8}$$

where $\kappa_{l,c}$ and $\varphi_{l,c}$ denote the amplitude and phase of the $c$-th element in $\mathbf{h_l}$, respectively.

## References

1. Lei, L.; Yuan, D.; Ho, C.K.; Sun, S. Power and Channel Allocation for Non-orthogonal Multiple Access in 5G Systems: Tractability and Computation. *IEEE Trans. Wirel. Commun.* **2016**, *15*, 8580–8594. [CrossRef]
2. Xiao, Z.; Zhu, L.; Choi, J.; Xia, P.; Xia, X.G. Joint Power Allocation and Beamforming for Non-Orthogonal Multiple Access (NOMA) in 5G Millimeter-Wave Communications. *IEEE Trans. Wirel. Commun.* **2018**, *17*, 2961–2974. [CrossRef]
3. Fang, F.; Cheng, J.; Ding, Z. Joint Energy Efficient Subchannel and Power Optimization for a Downlink NOMA Heterogeneous Network. *IEEE Trans. Veh. Technol.* **2018**, *62*, 1351–1364. [CrossRef]
4. Saito, Y.; Benjebbour, A.; Kishiyama, Y.; Nakamura, T. System-level performance evaluation of downlink non-orthogonal multiple access (NOMA). In Proceedings of the 2013 IEEE 24th International Symposium on Personal Indoor and Mobile Radio Communications (PIMRC), London, UK, 8–11 September 2013.
5. Choi, J. On generalized downlink beamforming with NOMA. *J. Commun. Netw.* **2017**, *19*, 319–328. [CrossRef]
6. Zhao, J.; Liu, Y.; Chai, K.K.; Nallanathan, A.; Yue, C.; Zhu, H. Spectrum Allocation and Power Control for Non-Orthogonal Multiple Access in HetNets. *IEEE Trans. Wirel. Commun.* **2017**, *16*, 5825–5837. [CrossRef]
7. Li, Q.; Zhang, Q.; Qin, J. A Special Class of Fractional QCQP and Its Applications on Cognitive Collaborative Beamforming. *IEEE Trans. Signal Process.* **2014**, *62*, 2151–2164. [CrossRef]

8.  Ali, M.S.; Tabassum, H.; Hossain, E. Dynamic User Clustering and Power Allocation for Uplink and Downlink Non-Orthogonal Multiple Access (NOMA) Systems. *IEEE Access* **2017**, *4*, 6325–6343. [CrossRef]
9.  Sun, X.; Yang, N.; Yan, S.; Ding, Z.; Ng, D.; Shen, C.; Zhong, Z. Joint Beamforming and Power Allocation in Downlink NOMA Multiuser MIMO Networks. *IEEE Trans. Wirel. Commun.* **2018**, *17*, 5367–5381. [CrossRef]
10.  Pan, Y.; Chen, M.; Yang, Z.; Huang, N.; Shikh-Bahaei, M. Energy-Efficient NOMA-Based Mobile Edge Computing Offloading. *IEEE Commun. Lett.* **2018**, *23*, 310–313. [CrossRef]
11.  Guo, S.; Zhou, X. Robust power allocation for NOMA in heterogeneous vehicular communications with imperfect channel estimation. In Proceedings of the 2017 IEEE 28th Annual International Symposium on Personal, Indoor, and Mobile Radio Communications (PIMRC), Montreal, QC, Canada, 8–13 October 2017.
12.  Xu, Y.; Chao, S.; Ding, Z.; Sun, X.; Shi, Y.; Gang, Z.; Zhong, Z. Joint Beamforming and Power Splitting Control in Downlink Cooperative SWIPT NOMA Systems. *IEEE Trans. Signal Process.* **2017**, *65*, 4874–4886. [CrossRef]
13.  Zhang, H.; Fang, F.; Cheng, J.; Long, K.; Wang, W.; Leung, V. Energy-Efficient Resource Allocation in NOMA Heterogeneous Networks. *IEEE Wirel. Commun.* **2018**, *25*, 48–53. [CrossRef]
14.  Nibedita, N.; Sudhan, M.; Wu, H.C. Secure Beamforming for MIMO-NOMA Based Cognitive Radio Network. *IEEE Commun. Lett.* **2018**, *22*, 1708–1711.
15.  Qin, Z.; Yue, X.; Liu, Y.; Ding, Z.; Arumugam, N. User Association and Resource Allocation in Unified NOMA Enabled Heterogeneous Ultra Dense Networks. *IEEE Commun. Mag.* **2018**, *56*, 86–92. [CrossRef]
16.  Hai, L.; Gao, F.; Shi, J.; Li, G.Y. A New View of Multi-User Hybrid Massive MIMO: Non-Orthogonal Angle Division Multiple Access. *IEEE J. Sel. Areas Commun.* **2017**, *35*, 2268–2280.
17.  Al-Hussaibi, W.A.; Ali, F.H. Efficient User Clustering, Receive Antenna Selection, and Power Allocation Algorithms for Massive MIMO-NOMA Systems. *IEEE Access* **2019**, *7*, 31865–31882. [CrossRef]
18.  Ding, Z.; Yang, Z.; Fan, P.; Poor, H.V. On the Performance of Non-Orthogonal Multiple Access in 5G Systems with Randomly Deployed Users. *Signal Process. Lett. IEEE* **2014**, *21*, 1501–1505. [CrossRef]
19.  Ding, J.; Cai, J.; Yi, C. An Improved Coalition Game Approach for MIMO-NOMA Clustering Integrating Beamforming and Power Allocation. *IEEE Trans. Veh. Technol.* **2019**, *68*, 1672–1687. [CrossRef]
20.  Solaiman, S.; Nassef, L.; Fadel, E. User Clustering and Optimized Power Allocation for D2D Communications at mmWave Underlaying MIMO-NOMA Cellular Networks. *IEEE Access* **2021**, *9*, 57726–57742. [CrossRef]
21.  Wang, Q.; Wu, Z. Beamforming Optimization and Power Allocation for User-Centric MIMO-NOMA IoT Networks. *IEEE Access* **2020**, *9*, 339–348. [CrossRef]
22.  Wang, Z.; Lin, Z.; Lv, T.; Ni, W. Energy-Efficient Resource Allocation in Massive MIMO-NOMA Networks with Wireless Power Transfer: A Distributed ADMM Approach. *IEEE Internet Things J.* **2021**, *8*, 14232–14247. [CrossRef]
23.  Ali, M.S.; Hossain, E.; Kim, D.I. Non-Orthogonal Multiple Access (NOMA) for Downlink Multiuser MIMO Systems: User Clustering, Beamforming, and Power Allocation. *IEEE Access* **2016**, *5*, 565–577. [CrossRef]