

Joint Channel Bandwidth and Power Allocations for Downlink Non-orthogonal Multiple Access Systems

Yuan Wu^{1,2}, Liping Qian¹, Haowei Mao¹, Weidang Lu¹, Haibo Zhou³, Changsheng Yu⁴

¹College of Information Engineering, Zhejiang University of Technology, Hangzhou, China

²State Key Laboratory of Integrated Services Networks, Xidian University, Xian, 710071, China

³Department of Electrical and Computer Engineering, University of Waterloo, ON, Canada

⁴Nokia Solutions and Networks System Technology, Co. Ltd, Hangzhou, China

Abstract—The advanced non-orthogonal multiple access (NOMA) has been considered as a promising scheme to satisfy the ultimate goals of future 5G cellular networks for providing ultra-high throughput and ultra-dense connections. By enabling a group of mobile users (MUs) to simultaneously share a same frequency channel and adopting successive interference cancellation to mitigate the co-channel interference, the NOMA can significantly improve the spectrum efficiency compared with the conventional orthogonal multiple access (OMA). However, due to cellular operators' limited and crowded spectrum resources, a critical question is how to properly size the channel bandwidth for the NOMA-enabled transmission to satisfy all MUs' traffic demands. In this paper, we propose a joint optimization scheme of bandwidth and power allocations for the NOMA-enabled downlink transmission, with the objective of minimizing the overall resource consumption cost that accounts for both the spectrum consumption cost and power consumption cost. In spite of the non-convexity nature of the joint optimization problem, we propose an efficient algorithm to compute the optimal bandwidth allocation and power allocation. Numerical results validate the proposed algorithm and the performance advantage of the proposed NOMA-enabled transmission in saving the overall resource consumption cost.

I. INTRODUCTION

Non-orthogonal multiple access (NOMA) has been considered as a promising multiple access scheme to achieve ultimate goals of providing ultra-high throughput, ultra-low latency, and ultra-dense connectivity in future fifth generation (5G) cellular systems [1] as well as the promising practical applications of 5G systems such as future vehicular networks and Internet of Things [2] [3]. The key advantage of NOMA lies in that it enables a group of mobile users (MUs) to adopt the power domain division which allows the MUs to share a same frequency channel/time-slot/code simultaneously and improves spectrum-efficiency compared with the conventional orthogonal multiple access (OMA). By further adopting the successive interference cancellation (SIC) to mitigate the co-channel interference, the NOMA further enhances the MUs' throughput and power-efficiency. Due to its potential advantages, the NOMA has attracted lots of research interests.

In [4], Ding *et al.* analyzed the performance of downlink NOMA by focusing on the scenario of MUs being uniformly

and randomly deployed within the cell. Further taking into account the partial channel information, the performance NOMA has been analyzed in [5]. Both analysis in [4] and [5] have demonstrated that NOMA can achieve superior performance compared to the conventional OMA. Due to allowing the co-channel interference among MUs in NOMA, both the user-scheduling (or pairing) in NOMA and the corresponding power allocation are critical to the performance of NOMA. For instance, the impact of user-pairing in NOMA has been investigated in [7]. Power allocation schemes have been proposed in [10] and [11] to improve the energy efficiency in downlink NOMA systems. The joint user-association and power allocation scheme for uplink NOMA has been studied in [8]. In [9], Wu *et al.* proposed a NOMA-enabled traffic offloading through small-cell networks. The SIC in NOMA requires the MUs with weak channel power gains to suffer from the interference the MUs with strong gains, which leads to a fairness in NOMA. In [6], Timotheou *et al.* proposed a power allocation scheme that ensures the fairness among users in downlink NOMA. Further taking into account the multi-carrier NOMA, it is a critical issue about how to assign MUs to share different sub-carriers in NOMA [12]–[14]. In [12], Sun *et al.* studied the joint power and subcarrier allocations for downlink multi-carrier NOMA with the objective of maximizing the weighted system throughput. In [13], Di *et al.* proposed a joint sub-channel assignment and power allocation for downlink NOMA to achieve the balance between the number of served users and the total throughput. In [14], Lei *et al.* proposed a joint channel and power allocation for downlink NOMA for maximizing the overall served throughput.

However, in all the aforementioned studies, the (sub)channel bandwidth used for NOMA has been treated as fixed. It is an open question about how to properly size the channel bandwidth used for NOMA, which is an important issue due to the limited even crowded spectrum resource of cellular operators. In this paper, we thus investigate the joint bandwidth and power allocation for the downlink NOMA, with the objective of minimizing a system-wise resource usage cost comprised of the power consumption cost and the bandwidth consumption cost, while satisfying all MUs' traffic demands. Despite the non-convexity of the formulated joint optimization problem, we identify its hidden convexity (after invoking some equivalent transformations) and further propose an efficient algorithm to compute the optimal bandwidth and power al-

This work was supported, in partial, by the National Natural Science Foundation of China (61572440), the Zhejiang Provincial Natural Science Foundation of China (LR17F010002 and LR16F010003), and the Young Talent Cultivation Project of Zhejiang Association for Science and Technology (2016YCGC011).

locations. Numerical results have validated the effectiveness of the proposed algorithm and the advantage of the proposed joint optimization scheme.

II. SYSTEM MODEL AND PROBLEM FORMULATION

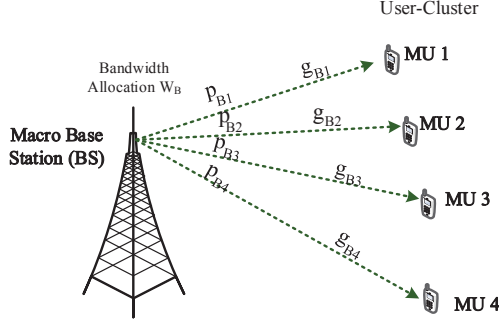


Fig. 1. Considered system model in this paper (assuming that $g_{B1} > g_{B2} > g_{B3} > g_{B4}$ in the user-cluster)

We consider a downlink offloading scenario in which a group of T MUs are under the coverage of a macro base station (BS). The BS uses NOMA to send data to the MUs. Since in NOMA, the MUs' successive interference cancellation (SIC) requires an ordering of the MUs according to their channel power gains with respect to the BS. We introduce the index-set \mathcal{I} , in which the group of T MUs follow the following ordering:

$$g_{B1} > g_{B2} > \dots > g_{Bi} > g_{Bj} > \dots > g_{BT}, \quad (1)$$

where g_{Bi} denotes the channel power gain from the BS to the i -th MU (we will illustrate the detailed modeling of channel gain in the numerical section).

We consider that the BS knows the instantaneous channel power gains $\{g_{Bi}\}_{i \in \mathcal{I}}$ of all MUs. Based on NOMA, the BS broadcasts the superposition of signals to all the MUs via power domain division. At the MU-side, SIC is used to eliminate the MUs' interference. Specifically, for MU i , it firstly decodes the message of MU j (with $j > i$) and then removes the decoded message from the received signal (in the order of $j = T, T-1, T-2, \dots, i+1$). Meanwhile, for MU i , it treats the message of MU j (with $j < i$) as noise. With the above decoding scheme, the throughput from the BS to MU i (i.e., the i -th MU in \mathcal{I}) can be given by:

$$R_{Bi} = W_B \log_2 \left(1 + \frac{g_{Bi} p_{Bi}}{g_{Bi} \sum_{j=1}^{i-1} p_{Bj} + W_B n_0} \right), \forall i \in \mathcal{I}. \quad (2)$$

We use p_{Bi} to denote the BS's transmit-power for MU i , and W_B to denote the BS's bandwidth allocation for serving the group of MUs. Parameter n_0 denotes the power density of the background noise, and the subscript "B" denotes the BS.

Our objective is to minimize the BS's system-wise resource consumption cost comprised of the bandwidth usage cost and power consumption cost to satisfy all MUs' traffic demands,

which corresponds to the following Total resource consumption Cost Minimization (TCM) problem:

$$(TCM): \quad \min \alpha \sum_{i \in \mathcal{I}} p_{Bi} + \beta W_B$$

$$\text{subject to:} \quad \sum_{i \in \mathcal{I}} p_{Bi} \leq P_B^{\text{tot}}, \quad (3)$$

$$W_B \leq W_B^{\text{tot}}, \quad (4)$$

$$R_{Bi} \geq R_i^{\text{req}}, \forall i \in \mathcal{I}, \quad (5)$$

$$\text{variables:} \quad p_{Bi} > 0, \forall i \in \mathcal{I}, \text{ and } W_B > 0.$$

In the objective function, we use parameters α and β to denote the respective emphasis on the BS's transmit-power consumption and the BS's bandwidth consumption. A larger ratio of β/α means that we put more emphasis on reducing the BS's bandwidth usage, while a smaller ratio of β/α means that we put more emphasis on reducing the BS's power consumption. We notice that similar expressions of the linear combination have also been used to minimize system-wise resource utilizations in several studies (e.g., [15] [16]). Constraint (3) means that the BS's total transmit-power cannot exceed the BS's maximum power-capacity P_B^{tot} . Constraint (4) means that the BS's bandwidth allocation cannot exceed the maximum bandwidth capacity W_B^{tot} . Finally, we use parameter R_i^{req} in (5) to denote MU i 's throughput requirement to achieve. Directly solving Problem (TCM) is challenging. Our key idea to solve Problem (TCM) is to transform it into a bandwidth allocation problem as shown in the next section.

III. EQUIVALENT TRANSFORMATION INTO A BANDWIDTH ALLOCATION PROBLEM

To transform Problem (TCM) into an equivalent bandwidth allocation problem, we introduce β_{Bi} to denote the received signal to interference plus noise ratio (SINR) from the BS to the i -th MU in \mathcal{I} , namely,

$$\beta_{Bi} = \frac{g_{Bi} p_{Bi}}{g_{Bi} \sum_{j=1}^{i-1} p_{Bj} + W_B n_0}, \forall i \in \mathcal{I}. \quad (6)$$

Given $\{\beta_{Bi}\}_{i \in \mathcal{I}}$, we can use the following recursive computing to derive the required minimum transmit-power for the BS to transmit to the i -th MU in \mathcal{I} :

$$p_{Bi} = \beta_{Bi} \sum_{j=1}^{i-1} p_{Bj} + \beta_{Bi} \frac{W_B n_0}{g_{Bi}}, \forall i \in \mathcal{I}. \quad (7)$$

The recursive computing (7) shows that each MU i 's power allocation is increasing in $\{\beta_{Bj}\}_{j \leq i}$. Thus, taking into account constraint (5), we have the following lemma.

Lemma 1: The global optimum of Problem (TCM) is achieved, when each MU i achieves $\beta_{Bi} = \beta_{Bi}^{\min} = 2^{\frac{R_i^{\text{req}}}{W_B}} - 1$.

Proof: Suppose that the optimal solution of Problem (TCM) leads to $\beta_{Bi}^* > \beta_{Bi}^{\min}$ for any MU i , we can always slightly reduce β_{Bi}^* without violating any constraint but reducing the objective function, which leads to a contradiction. ■

With Lemma 1, we have the following important result that characterizes the BS's minimum total power consumption to satisfy constraint (5).

Proposition 1: Given the MUs' traffic demands $\{R_i^{\text{req}}\}_{i \in \mathcal{I}}$, the overall minimum transmit-power for the BS to transmit to the MUs, as a function of the BS's bandwidth allocation W_B , can be compactly given by:

$$p_B^{\min}(W_B) = W_B \sum_{i=1}^T \left(\frac{n_0}{g_{Bi}} - \frac{n_0}{g_{Bi-1}} \right) (2^{\frac{1}{W_B} \sum_{j=i}^T R_j^{\text{req}}} - 1), \quad (8)$$

where for simplicity, we set the auxiliary parameter g_{B0} as a sufficiently large number such that $\frac{n_0}{g_{B0}} = 0$.

Proof: The proof of eq. (8) is similar to the proof of Proposition 1 in [9], except that we now explicitly specify the BS's bandwidth allocation W_B (which is treated as a fixed value in the proof). We thus skip the detailed procedures. ■

With Proposition 1, we can thus transform Problem (TCM) as an equivalent Bandwidth Allocation (BA) problem:

$$\begin{aligned} \text{(BA):} \quad & \min \alpha p_B^{\min}(W_B) + \beta W_B \\ \text{subject to:} \quad & p_B^{\min}(W_B) \leq P_B^{\text{tot}}, \end{aligned} \quad (9)$$

$$W_B \leq W_B^{\text{tot}}, \quad (10)$$

$$\text{variables:} \quad W_B > 0.$$

Compared with Problem (TCM), Problem (BA) only involves a single variable W_B which thus easier to solve.

IV. PROPOSED ALGORITHM TO SOLVE PROBLEM (BA)

However, Problem (BA) is still a non-convex optimization, i.e., no general existing algorithm that can directly solve it. To efficiently solve Problem (BA), we introduce the following variable-change as:

$$x = \frac{1}{W_B}. \quad (11)$$

With x and Proposition 1, we can transform Problem (BA) into the following Equivalent form (here, "E" denotes the "Equivalent"):

(BA-E):

$$\begin{aligned} \min & \alpha \frac{1}{x} \sum_{i=1}^T \left(\frac{n_0}{g_{Bi}} - \frac{n_0}{g_{Bi-1}} \right) (2^{x \sum_{j=i}^T R_j^{\text{req}}} - 1) + \beta \frac{1}{x} \\ \text{subject to:} \quad & \frac{1}{x} \sum_{i=1}^T \left(\frac{n_0}{g_{Bi}} - \frac{n_0}{g_{Bi-1}} \right) (2^{x \sum_{j=i}^T R_j^{\text{req}}} - 1) \leq P_B^{\text{tot}}, \\ \text{variables:} \quad & x \geq \frac{1}{W_B^{\text{tot}}} \end{aligned} \quad (12)$$

Although Problem (BA-E) is still nonconvex (due to the non-convexity of the objective function and constraint (12)), we can efficiently solve Problem (BA-E) as follows. Specifically, we introduce an additional variable v and consider the following optimization problem with respect to (v, x)

(BA-EV): $\min v$

$$\begin{aligned} \text{subject to:} \quad & \sum_{i=1}^T \left(\frac{n_0}{g_{Bi}} - \frac{n_0}{g_{Bi-1}} \right) (2^{x \sum_{j=i}^T R_j^{\text{req}}} - 1) \leq \frac{v}{\alpha} x - \frac{\beta}{\alpha}, \\ & \sum_{i=1}^T \left(\frac{n_0}{g_{Bi}} - \frac{n_0}{g_{Bi-1}} \right) (2^{x \sum_{j=i}^T R_j^{\text{req}}} - 1) \leq x P_B^{\text{tot}}, \\ \text{variables:} \quad & x \geq \frac{1}{W_B^{\text{tot}}} \text{ and } v \geq 0. \end{aligned}$$

Notice that Problem (BA-EV) is essentially equivalent to Problem (BA-E), and the optimal value of the objective function v^* corresponds to the minimum total resource consumption cost of the very original Problem (TCM).

An important observation of Problem (BA-EV) is as follows. If the value of v is fixed, Problem (BA-EV) turns to a convex feasibility-check problem. Moreover, given the value of v , this feasibility-check problem can be equivalently expressed as the following optimization problem:

(BA-EVsub):

$$\begin{aligned} H_v^* &= \min \sum_{i=1}^T \left(\frac{n_0}{g_{Bi}} - \frac{n_0}{g_{Bi-1}} \right) (2^{x \sum_{j=i}^T R_j^{\text{req}}} - 1) - \\ & \quad \min \left\{ \frac{v}{\alpha} x - \frac{\beta}{\alpha}, x P_B^{\text{tot}} \right\} \\ \text{variables:} \quad & x \geq \frac{1}{W_B^{\text{tot}}}. \end{aligned}$$

If $H_v^* < 0$, then Problem (BA-EV) is feasible under the given v , which implies that we can further reduce the value of v . On the other hand, if $H_v^* \geq 0$, then Problem (BA-EV) is infeasible, which implies that we need to increase the value of v . Such a updating continues until H_v^* sufficiently approaches to zero.

To compute H_v^* , we have the following important property regarding Problem (BA-EVsub).

Proposition 2: Given v , Problem (BA-EVsub) is a convex optimization problem with respect to x .

Proof: Due to $v > 0$, the piece-wise minimum function $\min \left\{ \frac{v}{\alpha} x - \frac{\beta}{\alpha}, x P_B^{\text{tot}} \right\}$ is concave with respect to x . Therefore, Problem (BA-EVsub) is a convex optimization problem [17]. ■

Furthermore, we have the following property regarding H_v^* .

Proposition 3: The optimal value H_v^* of Problem (BA-EVsub) is a non-increasing function of v .

Proof: It can be verified that the objective function of Problem (BA-EVsub) is non-increasing in v . Moreover, the feasible region does not depend on v . Thus, the optimal value of the objective function, i.e., H_v^* , is non-increasing in v . ■

Proposition 2 indicates that we can efficiently solve Problem (BA-EVsub) (e.g., by using the interior point method [17]) and compute the value of H_v^* under a given v . With H_v^* and further based on Proposition 3, we propose the following algorithm (i.e., Algorithm (Sol-BA)) to solve Problem (BA-EV). Specifically, Algorithm (Sol-BA) executes the bisection search on v until we reach $H_v^* = 0$ with the specified accuracy. As we have described below Problem (BA-EV), v^* output by Algorithm (Sol-BA) corresponds to the minimum total resource consumption cost of the original Problem (TCM).

By using the output of Algorithm (Sol-BA), we can derive the optimal bandwidth allocation of Problem (TCM) as:

$$W_B^* = \frac{1}{x^*}, \quad (13)$$

and we can recursively derive the optimal transmit-power

Algorithm (Sol-BA): to solve Problem (BA-EV) and compute x^* and v^*

```

1: Input:  $v^{\max}, v^{\min}$ , and  $\text{tol}$ .
2: while  $|v^{\max} - v^{\min}| \geq \text{tol}$  do
3:   Set  $v = \frac{v^{\max} + v^{\min}}{2}$ 
4:   Given  $v$ , solve Problem (BA-EVSub) to obtain  $H_v^*$  and the corresponding  $x_v^*$ .
5:   if  $H_v^* < 0$  then
6:     Update  $v^{\min} = v$ .
7:   else
8:     Update  $v^{\max} = v$ .
9:   end if
10: end while
11: Output:  $x^* = x_v^*$  and  $v^* = v$  for Problem (BA-EV).

```

allocation for each MU i as

$$p_{Bi}^* = (2^{x^* R_i^{\text{req}}} - 1) \left(\sum_{j=1}^{i-1} p_{Bj}^* + \frac{n_0}{g_{Bi}} \frac{1}{x^*} \right), \quad (14)$$

according to (7). We thus finish solving the original Problem (TCM) completely.

V. NUMERICAL RESULTS

This section evaluates the performance of the proposed Algorithm (Sol-BA) (to solve Problem (TCM)) and the proposed joint optimization of the bandwidth and power allocations for NOMA. We setup a scenario in which the BS is located at the origin, and all MUs are randomly distributed within a disk whose central is the BS, and the radius is 100. To evaluate the performance, we use a 10-MU scenario, i.e., ten different MUs are randomly distributed as described. We first set $\alpha = 1$ and $\beta = 0.5$ (we will vary α and β later on).

Figure 2 illustrates the rationale of the proposed Algorithm (Sol-BA) to solve Problem (TCM). The top-subplot shows the convergence of auxiliary variable v (to the minimum overall bandwidth consumption cost and power consumption cost). For the purpose of comparison, we also use the enumeration method, in which we enumerate W_B with a small step-size, to solve Problem (TCM-E). The horizontal line in red denotes the result obtained by the enumeration method. The top-subplot shows that the proposed Algorithm (Sol-BA) quickly converges to the minimum overall cost, which validates the effectiveness of Algorithm (Sol-BA). The bottom-subplot further illustrates the rationale of Algorithm (Sol-BA) that executes the bisection search on v according to the value of H_v^* (which is the output of Problem (BAEV)). Algorithm (Sol-BA) converges (i.e., the value of v approaches to the minimum overall cost) until H_v^* approach to zero with the specified accuracy.

Figure 3 further validates the effectiveness of Algorithm (Sol-BA). Specifically, we plot the total resource consumption cost versus different BS's bandwidth allocation. For the case of $R_i^{\text{req}} = 10\text{Mbps}$ for each MU, we use star to mark out v^* (i.e., the minimum total resource consumption cost) obtained by Algorithm (Sol-BA). Meanwhile, for the case of $R_i^{\text{req}} = 9\text{Mbps}$, we use circle to mark out v^* obtained by Algorithm (Sol-BA). Both cases show that the output of Algorithm (Sol-BA) correspond to the minimum of the total resource consumption cost versus different W_B , which thus validates the effectiveness of Algorithm (Sol-BA).

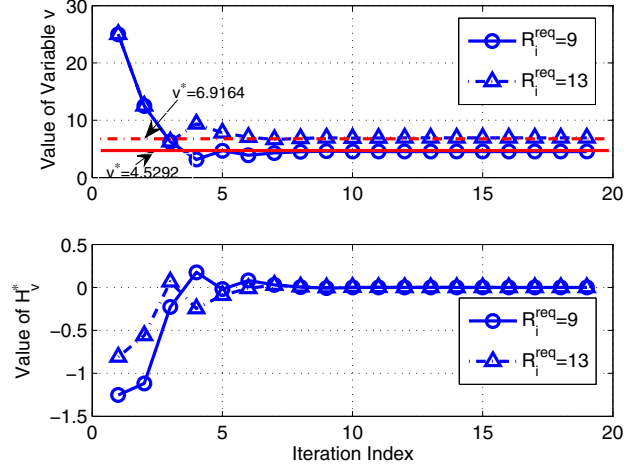


Fig. 2. Convergence of Proposed Algorithm (Sol-BA). Top-subplot: Convergence of variable v ; Bottom-subplot: Convergence of H_v^* .

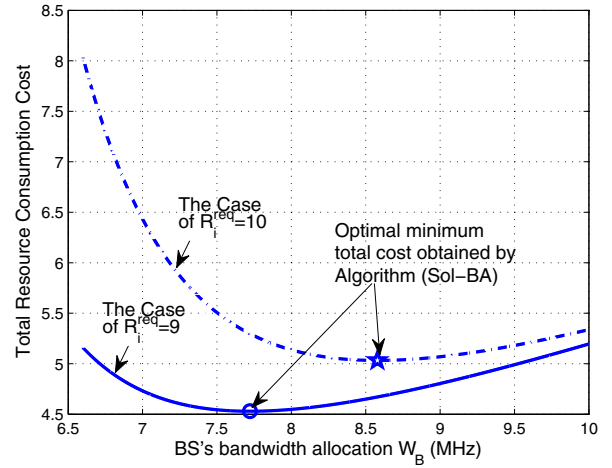


Fig. 3. Validation of proposed Algorithm (Sol-BA).

Figure 4 shows the advantage of the proposed joint optimization scheme in saving the overall resource consumption cost. For the purpose of comparison, we consider three other heuristic schemes, namely, the BS's allocated 30%, 50%, and 70% of the total bandwidth W_B^{tot} for NOMA. In particular, we consider two cases of $W_B^{\text{tot}} = 10\text{MHz}$ (in the left-subplot) and $W_B^{\text{tot}} = 15\text{MHz}$ (in the right-subplot). For each case, we plot the relative saving-gain, i.e., the value of $\frac{v(q) - v^*}{v^*}$ (where $v(q)$ denotes the overall resource consumption cost by setting $W_B = W_B^{\text{tot}} \times q$ with $q = 30\%$, 50% , and 70% , respectively). In all the tested cases, the relative saving-gains are positive, meaning that the proposed joint optimization scheme always outperforms the other three heuristic schemes. This result verifies the importance of taking into account the channel bandwidth allocation in NOMA. In particular, as shown in Figure 4, compared with the heuristic scheme with fixed bandwidth allocation $W_B^{\text{tot}} \times 30\%$, the relative saving-gain first decreases and then increases, when the MUs' traffic demand R_i^{req} increases. Such a curve of U-shape is consistent

with the intuition. Let us take the case of $W_B^{\text{tot}} \times 30\%$ as an example. Specifically, $W_B^{\text{tot}} \times 30\%$ might correspond to the optimal bandwidth allocation of Problem (TCM) (namely, $W_B^* = W_B^{\text{tot}} \times 30\%$) at a particular R_i^{req} , when we enumerate R_i^{req} . This leads to that the relative saving-gain is zero at the particular R_i^{req} . However, the larger the difference from such a particular R_i^{req} , the larger the relative saving-gain achieved by our proposed joint optimization scheme.

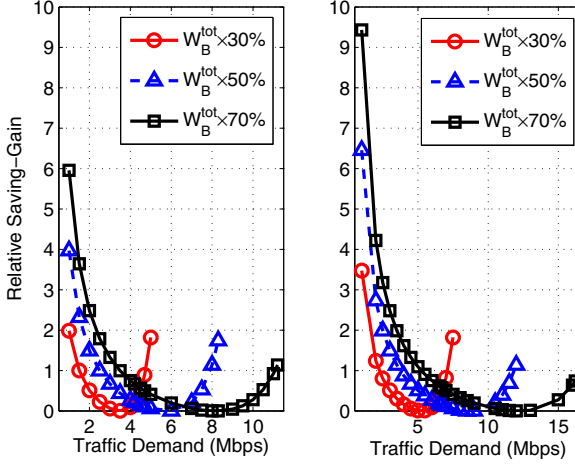


Fig. 4. Relative saving-gain of the proposed joint optimization scheme. Left-subplot: $W_B^{\text{tot}} = 10\text{MHz}$; Right-subplot: $W_B^{\text{tot}} = 15\text{MHz}$.

Figure 5 shows the impact of weighting factor (α, β) on the optimal bandwidth usage and the corresponding minimum total power consumption. Specifically, we vary the ratio of β/α from 0.4 to 2 (with fixed $\alpha = 1$) and plot the corresponding optimal bandwidth allocation W_B^* (in the left-subplot) and the minimum total power consumption $p_B^{\text{min}}(W_B^*)$ (in the right-subplot). As shown in Figure 5, when the ratio of β/α increases, the optimal bandwidth usage decreases while the minimum total power consumption increases, which thus validates the tradeoff between the bandwidth usage and power consumption in NOMA to satisfy the MUs' demands.

VI. CONCLUSION

In this paper, we have proposed the joint optimization scheme of bandwidth and power allocations for the NOMA-enabled downlink transmission, with the objective of minimizing the overall spectrum consumption cost and power consumption cost. In spite of the non-convexity nature of the joint optimization problem, we have proposed the efficient algorithm to compute the optimal bandwidth allocation and power allocation. Numerical results validate the proposed algorithm and the performance advantage of the proposed NOMA-enabled transmission in saving the overall spectrum consumption and power consumption cost.

REFERENCES

[1] L. Dai, B. Wang, Y. Yuan, S. Han, C.-L.I., and Z. Wang, "Non-orthogonal multiple access for 5G: solutions, challenges, opportunities, and future research trends," *IEEE Communications Magazine*, vol. 53, no. 9, pp. 74-81, Sep. 2015.

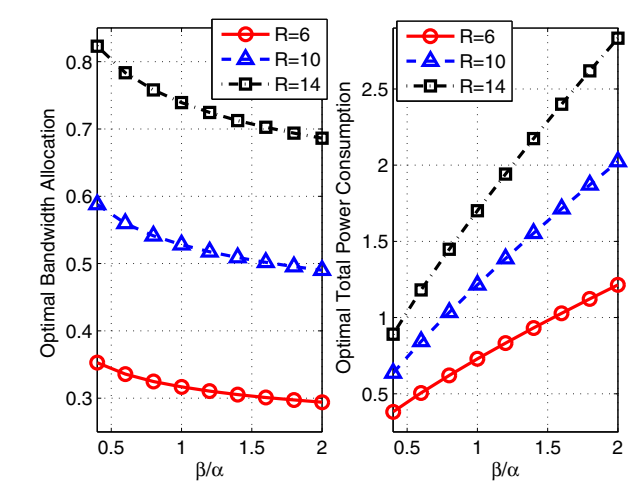


Fig. 5. Influence of weighting factors on the optimal resource consumption. Left-subplot: Optimal bandwidth usage; Right-subplot: Minimum total power consumption.

[2] Z. Su, Y. Hui, and Q. Yang, "The Next Generation Vehicular Networks: A Content Centric Framework", *IEEE Wireless Communications*, vol.24, no.1, pp.60-66, Feb. 2017.

[3] Y. Hui, et. al., "Utility Based Data Computing Scheme to Provide Sensing Service in Internet of Things", to appear in *IEEE Transactions on Emerging Topics in Computing*, June 2017, DOI:10.1109/TETC.2017.2674023.

[4] Z. Ding, et. al., "On the performance of nonorthogonal multiple access in 5G systems with randomly deployed users," *IEEE Signal Processing Letter*, vol. 21, no. 12, pp. 1501-1505, Dec. 2014.

[5] Z. Yang, Z. Ding, P. Fan, and G.K. Karagiannidis, "On the performance of non-orthogonal multiple access systems with partial channel information," *IEEE Transactions on Communications*, vol. 64, no. 2, Feb. 2016.

[6] S. Timotheou, I. Kriridis, "Fairness for non-orthogonal multiple access in 5G systems," *IEEE Signal Processing Letters*, vol. 22, no. 10, pp. 1647-1651, Oct. 2015.

[7] Z. Ding, P. Fan, and H. V. Poor, "Impact of user pairing on 5G nonorthogonal multiple access," *IEEE Transactions on Vehicle Technology*, vol. 65, no. 8, pp. 6010-6023, Aug. 2016.

[8] L. Qian, Y. Wu, H. Zhou, and X. Shen "Joint uplink base station association and power control for small-cell networks with non-orthogonal multiple access," to appear in *IEEE Transactions on Wireless Communications*, DOI:10.1109/TWC.2017.2664832.

[9] Y. Wu, L. Qian, "Energy-efficient NOMA-enabled traffic offloading via dual-connectivity in small-cell networks", to appear in *IEEE Communications Letters*, DOI:10.1109/LCOMM.2017.2685384.

[10] Q. Sun, S. Han, C.-L.I and Z. Pan, "Energy efficiency optimization for fading MIMO non-orthogonal multiple access systems," in *Proc. of IEEE ICC'2015*.

[11] Y. Zhang, H. Wang, T. Zheng, and Q. Yang, "Energy-efficient transmission design in non-orthogonal multiple access," *IEEE Transactions on Vehicle Technology*, to be published.

[12] Y. Sun, D.W.K Ng, Z. Ding, and R. Schober, "Optimal Joint Power and Subcarrier Allocation for MC-NOMA Systems," in *Proc. of IEEE Globe Communication Conference (GlobeCom'2016)*.

[13] B. Di, S. Bayat, L. Song, and Y. Li, "Radio Resource Allocation for Downlink Non-Orthogonal Multiple Access (NOMA) Networks using Matching Theory," in *Proc. of IEEE Globecom'2015*.

[14] L. Lei, D. Yuan, C. K. Ho, and S. Sun, "Joint Optimization of Power and Channel Allocation with Non-Orthogonal Multiple Access for 5G Cellular Systems," *Proc. of IEEE Globecom'2014*.

[15] Y. Wu, et. al., "Energy-Aware Cooperative Traffic Offloading via Device-to-Device Cooperations: An Analytical Approach," *IEEE Transactions on Mobile Computing*, vol. 16, no. 1, pp. 97-114, Jan. 2017.

[16] Y. Wu, et. al., "Secrecy-based Energy-Efficient Data Offloading via Dual-Connectivity over Unlicensed Spectrums," *IEEE Journal on Selected Areas in Communications*, vol. 34, no. 12, pp. 3252-3270, Dec. 2016.

[17] S. Boyd, and L. Vandenberghe, "Convex Optimization," Cambridge University Press, 2004.