

Distributed Reinforcement Learning Algorithm for Dynamic Economic Dispatch with Unknown Generation Cost Functions

Dai, Pengcheng; Yu, Wenwu; Wen, Guanghui; Baldi, Simone

DOI

[10.1109/TII.2019.2933443](https://doi.org/10.1109/TII.2019.2933443)

Publication date

2020

Document Version

Final published version

Published in

IEEE Transactions on Industrial Informatics

Citation (APA)

Dai, P., Yu, W., Wen, G., & Baldi, S. (2020). Distributed Reinforcement Learning Algorithm for Dynamic Economic Dispatch with Unknown Generation Cost Functions. *IEEE Transactions on Industrial Informatics*, 16(4), 2258-2267. <https://doi.org/10.1109/TII.2019.2933443>

Important note

To cite this publication, please use the final published version (if applicable). Please check the document version above.

Copyright

Other than for strictly personal use, it is not permitted to download, forward or distribute the text or part of it, without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license such as Creative Commons.

Takedown policy

Please contact us and provide details if you believe this document breaches copyrights. We will remove access to the work immediately and investigate your claim.

Distributed Reinforcement Learning Algorithm for Dynamic Economic Dispatch With Unknown Generation Cost Functions

Pengcheng Dai, *Student Member, IEEE*, Wenwu Yu , *Senior Member, IEEE*, Guanghui Wen , *Senior Member, IEEE*, and Simone Baldi , *Member, IEEE*

Abstract—In this article, the dynamic economic dispatch (DED) problem for smart grid is solved under the assumption that no knowledge of the mathematical formulation of the actual generation cost functions is available. The objective of the DED problem is to find the optimal power output of each unit at each time so as to minimize the total generation cost. To address the lack of *a priori* knowledge, a new distributed reinforcement learning optimization algorithm is proposed. The algorithm combines the state-action-value function approximation with a distributed optimization based on multiplier splitting. Theoretical analysis of the proposed algorithm is provided to prove the feasibility of the algorithm, and several case studies are presented to demonstrate its effectiveness.

Index Terms—Distributed reinforcement learning, dynamic economic dispatch (DED), multiplier splitting, state-action-value function approximation.

I. INTRODUCTION

THE POWER grid is undergoing significant changes due to the integration of distributed energy resources, devel-

opment of smart technologies, high demand of transactions and energy management, and so on [1], [2]. Within this context, smart grids have received increasing attention [3]. The smart grid technology makes full use of communication and sensing in an effort to attain safe, efficient, stable, and sustainable power services [4]–[6]. In smart grids, the dynamic economic dispatch (DED) problem has attracted much attention. The aim of DED is to find the optimal power output of each generator at each time to minimize the total generation cost in a given time horizon. In most practical cases, the DED problem needs to be solved in a distributed way. It has been learned from existing literature that multiagent systems theory [7]–[9] is an appealing framework to solve such a problem. The static economic dispatch (SED) problem is a special case of DED which has also been studied in the framework of multiagent systems [10]–[20]. Specifically, a fully distributed λ -consensus algorithm was proposed in [10] for smart grids with a directed topology. The authors of [11] proposed a distributed discrete-time consensus algorithm under a jointly connected switching undirected topology. In [12], under a uniformly jointly strong connected directed graph with time-varying delays, some distributed gradient push-sum algorithms were discussed for SED. A distributed Laplacian-gradient algorithm was proposed in [13] with feasible initial point. Yi *et al.* [14] solved the SED problem via an initialization-free distributed algorithm based on the multiplier splitting method. Guo *et al.* [15] proposed an average consensus algorithm and the distributed projection gradient algorithm to solve SED with consideration of wind turbines and energy storage systems. A distributed auction-based algorithm was proposed in [16] to solve a nonconvex SED. In the presence of communication uncertainties, an adaptive incremental cost consensus-based algorithm was proposed in [18]. In contrast, few results on the DED problem are reported in the literature due to the complexity of this problem [21]–[23]. A distributed primal–dual dynamic algorithm was proposed in [21]. Zhao *et al.* [22] deal with a fully decentralized optimization for the multiarea DED through the cutting plane consensus algorithm. More recently, by integrating the average consensus protocol and alternating direction method of multipliers (ADMM), a distributed coordination algorithm has been proposed in [24] to solve the dynamic social welfare problem. In practice, the accurate mathematical expression of the cost functions in a DED problem may be unavailable as the

Manuscript received March 20, 2019; revised July 2, 2019; accepted July 22, 2019. Date of publication August 6, 2019; date of current version January 17, 2020. This work was supported in part by the National Natural Science Foundation of China under Grant 61673107 and Grant 61673104, in part by the National Ten Thousand Talent Program for Young Top-Notch Talents under Grant W2070082, in part by the General Joint Fund of the Equipment Advance Research Program of the Ministry of Education under Grant 6141A020223, in part by the Six Talent Peaks of Jiangsu Province under Grant 2019-DZXX-006, in part by the Fundamental Research Funds for the Central Universities under Grant 4007019109 (RECON-STRUCT), in part by the Special Guiding Funds for Double First-Class under Grant 4007019201, and in part by the Jiangsu Provincial Key Laboratory of Networked Collective Intelligence under Grant BM2017002. Paper no. TII-19-0995. (*Corresponding author: Wenwu Yu.*)

P. Dai and G. W. Wen are with the School of Mathematics, Southeast University, Nanjing 210096, China (e-mail: jldaipc@163.com; wenguanghui@gmail.com).

W. Yu is with the School of Mathematics, Southeast University, Nanjing 210096, China, and also with the Department of Electrical Engineering, Nantong University, Nantong 226019, China (e-mail: wwyu@seu.edu.cn).

S. Baldi is with the School of Mathematics, Southeast University, Nanjing 210096, China, and also with the Delft Center for Systems and Control, Delft University of Technology 2628CD Delft, The Netherlands (e-mail: s.baldi@tudelft.nl).

Color versions of one or more of the figures in this article are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TII.2019.2933443

cost functions are affected by various factors, such as operating conditions and aging of the generator. Most of the aforementioned algorithms no longer work when the accurate mathematical formation of the cost function is unavailable. Hence, it is of both theoretical and practical interest to design an algorithm to solve the DED problem with little information of the actual cost functions.

Reinforcement learning [25] is a method through which an agent can find the optimal policy by interacting with the environment. This has motivated the application of reinforcement learning algorithms in control and optimization problems, sometimes in the context of multiagent systems [26]–[31]. The reinforcement learning-based approach is used to investigate the optimal tracking control problem in [26]. Data-driven optimal control based on reinforcement learning was proposed in [27] for discrete-time multiagent systems with unknown dynamics. Wang *et al.* [28] proposed a dual heuristic dynamic programming algorithm for a class of nonlinear discrete-time systems affected by time-varying delay. The method of policy iteration in reinforcement learning was used in [29] to find the optimal control for zero-sum games. Exciting applications of deep reinforcement learning are [30] and [31], which show that an agent can learn to play Atari better than humans. In this article, we draw inspiration from reinforcement learning techniques, especially from state-action-value function approximation and from nonlinear programming theories to solve the DED problem with little information of actual cost functions.

The contributions of this article are as follows.

- 1) The techniques of state-action-value function approximation based on semigradient Q-learning and distributed optimization algorithm based on multiplier splitting are successfully combined in the proposed algorithm. This algorithm can deal with the situations in which the mathematical expression of the cost functions is not available.
- 2) The update of the operating policy depends not only on the optimal solution of the approximate state-action-value function but also on the last operating policy. This means that the cost can be proven to be monotonically nonincreasing at each iteration.
- 3) Time-varying parameters in approximate state-action-value function are proposed. As compared to the use of time-invariant parameters, they enable us to reduce the error and preserve convexity of approximate state-action-value function. To the best of our knowledge, this is the first attempt to employ time-varying parameters in the approximation of the state-action-value function.

The rest of this article is organized as follows. The DED problem is formulated in Section II. The distributed reinforcement learning optimization algorithm is proposed in Section III. Section IV confirms the feasibility of the distributed reinforcement learning optimization algorithm. Simulation results to demonstrate the effectiveness of the algorithm are provided in Section V. Finally, Section VI concludes this article. The Appendix gives preliminaries about convex analysis, algebraic graph theory, and reinforcement learning.

II. PROBLEM STATEMENT

A. Dynamic Economic Dispatch

We consider a smart grid setting where N units must make their electricity generation equal to the total power demand at each time slot t . The objective of the DED problem is to find the optimal electricity allocation such that the total generation cost of all units is minimized. The mathematical expression of this problem is

$$\begin{aligned} \min \quad & \sum_{t=1}^T \sum_{i=1}^N F_i(p_{i,t}) \\ \text{s.t.} \quad & \sum_{i=1}^N p_{i,t} = D_t, \quad t = 1, 2, \dots, T \\ & \underline{p}_i \leq p_{i,t} \leq \bar{p}_i, \quad i = 1, 2, \dots, N \\ & |p_{i,t} - p_{i,t-1}| \leq p_i^R, \quad i = 1, \dots, N, t = 1, \dots, T \end{aligned} \quad (1)$$

where $F_i(\cdot)$ is the generation cost function of unit i , $p_{i,t}$ is the power output of unit i at time t , D_t is the total power demand at time t , p_i^R is the ramp-rate limit of unit i . \underline{p}_i and \bar{p}_i are the minimum and maximum power output of unit i , respectively. For notational brevity, set $p_{i,0} + p_i^R = \bar{p}_i$, $p_{i,0} - p_i^R = \underline{p}_i$, and $D_t - \sum_{i=1}^N p_i^R \leq D_{t+1} \leq D_t + \sum_{i=1}^N p_i^R$, $t = 1, 2, \dots, T-1$. We denote $\mathcal{P}_i = [\underline{p}_i, \bar{p}_i]$ as the set of admissible power output of unit i .

Various forms of the generation cost function have been proposed in the literature. The most common generation cost function is $F_i(p_{i,t}) = a_i p_{i,t}^2 + b_i p_{i,t} + c_i$, where a_i , b_i , and c_i are some coefficients for unit i [19]. The cost function considered in this article is a more general sinusoidal cost function inspired by [33]

$$F_i(p_{i,t}) = a_i p_{i,t}^2 + b_i p_{i,t} + c_i + |e_i \cdot \sin(f_i \cdot (\underline{p}_i - p_{i,t}))|$$

where the additional coefficients e_i and f_i are related to the capacity of unit i . The mathematical expression of this cost function is known for simulation purposes, but it is unknown for the purpose of controller design.

When considering the above cost function, the following challenges should be taken into account: 1) the nonconvex objective function invalidates existing algorithms based on convex optimization problems and 2) only the value of the generation cost is known while the mathematical formulation of the cost function is unknown. Fortunately, reinforcement learning algorithm can be applied to tackle such challenges.

Remark 1: In the DED problem, the total demand D_t , the feasible power output combination (FPOC) of units and the generation cost at each time slot can be seen as the state, action, and reward in the mind of reinforcement learning. Furthermore, the generation cost at each time slot is also important and should be fully considered while dealing with the DED problem. Hence, the discount factor γ introduced in the step of reinforcement learning (cf. Appendix C) is set as 1 in the DED problem.

Two standard assumptions are made to guarantee the existence of an optimal distributed solution to (1).

Assumption 1: There exists at least one FPOC $(p_{1,1}, \dots, p_{N,1}, \dots, p_{1,T}, \dots, p_{N,T})^T$ at all times such that

$\sum_{i=1}^N p_{i,t} = D_t, p_{i,t} \in \mathcal{P}_i, |p_{i,t} - p_{i,t-1}| \leq p_i^R, t = 1, \dots, T, i = 1, \dots, N.$

Assumption 2: The graph topology about the units is undirected and connected. At each time slot t , each agent i can only access the local power demand $D_{i,t}$, adjust the local power output $p_{i,t}$, and obtain the local generation cost $F_i(p_{i,t})$.

III. DISTRIBUTED REINFORCEMENT LEARNING OPTIMIZATION ALGORITHM

In order to solve the DED problem with unknown cost functions, we apply reinforcement learning ideas. Suppose each agent corresponding to each unit was assigned a unique identifier ID, e.g., its IP address. By using the graph discovery algorithm proposed in [15], each agent can get the total number of agents. A distributed reinforcement learning optimization algorithm is proposed based on seven steps.

1) *Discover the total demand at time slot t :* Define $\bar{D}_t[0] = (D_{1,t}, D_{2,t}, \dots, D_{N,t})^T$. Apply the average-consensus protocol (18) for each agent i as follows:

$$\bar{D}_t[k+1] = \bar{D}_t[k] - \epsilon L \bar{D}_t[k] \quad (2)$$

where L is the Laplacian matrix of graph \mathcal{G} , $\epsilon \in (0, \frac{1}{\max_i l_{ii}})$.

From Lemma 1 in Appendix B, we get $\lim_{k \rightarrow \infty} \bar{D}_t[k] = (\frac{1}{N} D_t) \mathbf{1}_N$ where $\mathbf{1}_N$ is a N -dimensional column vector with each entry being one. Hence, the local estimation of the average power demand converges to the actual average power demand at time slot t . As result, the total demand at time slot t can be obtained as D_t .

2) *Find an FPOC at time slot t :* Choose $p_{i,t} \in (\max\{\underline{p}_i, p_{i,t-1} - p_i^R\}, \min\{\bar{p}_i, p_{i,t-1} + p_i^R\})$. Define the mismatch of demand generations $m_t[0] = (D_{1,t} - p_{1,t}, \dots, D_{N,t} - p_{N,t})^T$, and apply Lemma 1 in Appendix B again as follows:

$$m_t[k+1] = m_t[k] - \epsilon L m_t[k]. \quad (3)$$

It holds that $\lim_{k \rightarrow \infty} m_t[k] = \frac{1}{N} \sum_{i=1}^N (D_{i,t} - p_{i,t}) \mathbf{1}_N = \alpha \mathbf{1}_N$.

Adjust $p_{i,t}$ according to the following policy:

$$p_{i,t} \leftarrow \begin{cases} p_{i,t} + \text{sign}(\alpha) \min\{\min\{\bar{p}_i, p_{i,t-1} + p_i^R\} \\ - p_{i,t}, \alpha\}, \alpha \geq 0 \\ p_{i,t} + \text{sign}(\alpha) \min\{-\max\{\underline{p}_i, p_{i,t-1} - p_i^R\} \\ + p_{i,t}, |\alpha|\}, \alpha < 0 \end{cases} \quad (4)$$

where $\text{sign}(\cdot)$ is symbolic function. Repeat (3) and (4) till $\alpha = 0$.

When $\alpha = 0$, $P_t = (p_{1,t}, p_{2,t}, \dots, p_{N,t})^T$ is an FPOC at time slot t .

3) *Measure the total generation cost at time slot t :* Define $\bar{c}_t[0] = (c_{1,t}, \dots, c_{N,t})^T$ and \bar{c}_t as the local estimation of the average generation cost at time slot t , where $c_{i,t} = F_i(p_{i,t})$ for each agent i . Apply the average-consensus protocol

$$\bar{c}_t[k+1] = \bar{c}_t[k] - \epsilon L \bar{c}_t[k]. \quad (5)$$

As a result of Lemma 1 in Appendix B, we can obtain $\lim_{k \rightarrow \infty} \bar{c}_t[k] = \bar{c}_t \mathbf{1}_N$, i.e., the local estimation of the average

generation cost converges to the actual average generation cost at time slot t . Hence, the total generation cost is $N\bar{c}_t$.

4) *Update the parameters of approximate function at time slot t :* Define $J_t(D_t, P_t, \theta^t) = \phi(P_t)^T \theta^t$ to be the approximate state-action-value function, where $\phi(P_t)$ is a feature vector. The update of the parameters θ^t is

$$\begin{cases} \theta^t \leftarrow \theta^t + \beta [N\bar{c}_t + \min_{P_{t+1}} J_{t+1}(D_{t+1}, P_{t+1}, \theta^{t+1}) \\ - J_t(D_t, P_t, \theta^t)] \phi(P_t) \end{cases} \quad (6)$$

The feature vector may be constructed from P_t in many different ways. For easier analysis, it is smart to design $\phi(P_t)$ such that the approximate state-action-value function is a convex function. For example, let $\phi(P_t) = (p_{1,t}, \dots, p_{N,t}, p_{1,t}^2, \dots, p_{N,t}^2)^T$, $\theta^t = (\theta_1^t, \dots, \theta_{2N}^t)^T$, and $f_i(p_{i,t}) = \theta_i^t p_{i,t} + \theta_{i+N}^t p_{i,t}^2$. Then, $J_t(D_t, P_t, \theta^t) = \phi(P_t)^T \theta^t = \sum_{i=1}^N f_i(p_{i,t})$, (6) becomes

$$\begin{cases} \theta^t \leftarrow \theta^t + \beta [N\bar{c}_t + \min_{P_{t+1}} J_{t+1}(D_{t+1}, P_{t+1}, \theta^{t+1}) \\ - J_t(D_t, P_t, \theta^t)] p_{i,t} \\ \theta_{i+N}^t \leftarrow \theta_{i+N}^t + \beta [N\bar{c}_t + \min_{P_{t+1}} J_{t+1}(D_{t+1}, P_{t+1}, \theta^{t+1}) \\ - J_t(D_t, P_t, \theta^t)] p_{i,t}^2 \end{cases} \quad (7)$$

Remark 2: $\min_{P_{t+1}} J_{t+1}(D_{t+1}, P_{t+1}, \theta^{t+1})$ in (7) can be obtained through step 5. Taking into account the particularity of the finite horizon in (1), we use time-varying parameters θ^t for each time slot t . This is done in order to guarantee that the approximate state-action-value function is a convex function (necessary for the analysis in Section IV). Equation (7) can be seen as a semigradient method applied to the state-action-value function [25].

5) *Obtain $\min_{P_t} J_t(D_t, P_t, \theta^t)$ in a distributed way:* Solve the following problem about approximate state-action-value function

$$\begin{aligned} \min \quad & \sum_{i=1}^N f_i(p_{i,t}) \\ \text{s.t.} \quad & \sum_{i=1}^N p_{i,t} = D_t \\ & p_{i,t} \in \mathcal{P}_i, i = 1, 2, \dots, N \\ & |p_{i,t} - p_{i,t-1}^*| \leq p_i^R, i = 1, 2, \dots, N \end{aligned} \quad (8)$$

where $p_{i,0}^* = p_{i,0}$ for each i . Before moving on, let $\mathcal{P}_{i,t}^{\text{new}} = \mathcal{P}_i \cap [p_{i,t-1}^* - p_i^R, p_{i,t-1}^* + p_i^R]$. Problem (8) can be solved under the following standard assumption.

Assumption 3: There exists a finite optimal solution P_t^{a*} to problem (8). The Slater's constraint condition is satisfied for (8), that is, there exist $\hat{p}_{i,t} \in \text{int}(\mathcal{P}_{i,t}^{\text{new}}) \forall i$, such that $\sum_{i=1}^N \hat{p}_{i,t} = D_t$.

Here is the procedure to solve (8). The duality of (8) with $\lambda \in \mathbb{R}$ is

$$\max_{\lambda \in \mathbb{R}} \sum_{i=1}^N q_i(\lambda) = \sum_{i=1}^N \inf_{p_{i,t} \in \mathcal{P}_{i,t}^{\text{new}}} \{f_i(p_{i,t}) - \lambda p_{i,t} + \lambda \frac{1}{N} D_t\}.$$

We formulate a constrained optimization problem with Laplacian matrix L and $\Lambda = (\lambda_1, \lambda_2, \dots, \lambda_N)^T \in \mathbb{R}^N$ as

$$\begin{aligned} \max_{\Lambda} \quad & \sum_{i=1}^N q_i(\lambda_i) \\ \text{s.t.} \quad & L\Lambda = 0_N. \end{aligned} \quad (9)$$

The augmented Lagrangian duality of (9) with multipliers $Z = (z_1, z_2, \dots, z_N)^T \in \mathbb{R}^N$ is

$$\min_Z \max_{\Lambda} \sum_{i=1}^N q_i(\lambda_i) - Z^T L\Lambda - \frac{1}{2} \Lambda L \Lambda.$$

The distributed algorithm for agent i is given as follows:

$$\begin{cases} \dot{p}_{i,t} = \mathbb{P}_{\mathcal{P}_{i,t}^{\text{new}}} (p_{i,t} - \nabla f_i(p_{i,t}) + \lambda_i) - p_{i,t} \\ \dot{\lambda}_i = (\frac{1}{N} D_t - p_{i,t}) - \sum_{j \in \mathcal{I}_i} (z_i - z_j) - \sum_{j \in \mathcal{I}_i} (\lambda_i - \lambda_j) \\ \dot{z}_i = \sum_{j \in \mathcal{I}_i} (\lambda_i - \lambda_j) \end{cases}. \quad (10)$$

From the Karush-Kuhn-Tucker (KKT) condition, the equilibrium point of (10) is the optimal solution to (8) (cf. analysis in Section IV). Denote one of such equilibrium points $\text{col}(P_t^{a*}, \Lambda^{a*}, Z^{a*})$ as the column vector stacked with vectors P_t^{a*} , Λ^{a*} , and Z^{a*} . Then, the value of $\sum_{i=1}^N (F_i(p_{i,t}^{a*}))$ can be obtained by Lemma 1 in Appendix B.

6) *Renew the local operating policy*: Renew the local operating policy according to the following algorithm.

Denote $W_{a*} = \sum_{t=1}^T \sum_{i=1}^N (F_i(p_{i,t}^{a*}))$, $W_p = \sum_{t=1}^T \sum_{i=1}^N (F_i(p_{i,t}))$, $W_\pi = \sum_{t=1}^T \sum_{i=1}^N (F_i(\pi_i(D_t)))$, then, the local operating policy can be renewed by

$$\pi_i(D_t) \leftarrow \begin{cases} p_{i,t}^{a*}, & \text{if } W_{a*} = \min\{W_{a*}, W_p, W_\pi\} \\ p_{i,t}, & \text{if } W_p = \min\{W_{a*}, W_p, W_\pi\} \\ \pi_i(D_t), & \text{otherwise} \end{cases} \quad (11)$$

where $P_t^{a*} = (p_{1,t}^{a*}, \dots, p_{N,t}^{a*})^T = \arg \min_{P_t} J_t(D_t, P_t, \theta^t)$. In particular, $\pi(D_t)$ is a determined policy in DED problem.

7) *Balance exploration and exploitation*: In order to balance exploration and exploitation, we use the ε -greedy policy, i.e., selecting the action $(\pi_1(D_t), \dots, \pi_N(D_t))^T$ with probability $1 - \varepsilon$, and other FPOC with probability ε .

The distributed reinforcement learning optimization algorithm for the DED problem is summarized in Algorithm 1.

Remark 3: In the process of developing the distributed algorithm, the key difficulties are: 1) How to determine the total power demand at each time slot by agents in a distributed way in the absence of a centralized decision-making agent with global information? 2) How to find an FPOC in a distributed way? 3) How to renew the local operating policy in a distributed manner? For issue 1, the total power demand D_t can be obtained by the average-consensus protocol (2). The aim of (3) and (4) is to solve issue 2 by finding an FPOC in a distributed way. Issue 3 is addressed through (11).

IV. THEORETICAL ANALYSIS

In this section, the main theoretical results of the proposed distributed reinforcement learning optimization algorithm are provided and proven via convex analysis and projection.

Algorithm 1: DED With Distributed Reinforcement Learning Optimization.

- 1: Initialize $t = 0, k = 0$;
 - 2: Initialize ε with ε -greedy policy;
 - 3: **Repeat**
 - 4: $t \leftarrow t + 1$;
 - 5: Obtain the total power demand D_t at time t via (2);
 - 6: Initialize the parameters θ^t of the approximate state action-value function;
 - 7: Set J_t with $\theta^t = \mathbf{0}$;
 - 8: **Until** $t = T$
 - 9: Define $J_{T+1} = 0$.
 - 10: **Repeat**
 - 11: $k \leftarrow k + 1$;
 - 12: $\tilde{r} = \text{rand}(1)$;
 - 13: Reset $t = 1, W_p = 0, W_{a*} = 0$;
 - 14: **Repeat**
 - 15: **If** $k \geq 2$ and $\tilde{r} \geq \varepsilon$ **then**
 - 16: **Repeat**
 - 17: Choose power output as $\pi(D_t)$;
 - 18: Obtain immediate generation cost of $\pi(D_t)$ via (5);
 - 19: Update the parameter θ^t through (7);
 - 20: $W_p \leftarrow W_p + N\bar{c}_t$;
 - 21: Find the P_t^{a*} of (8) by (10);
 - 22: Obtain immediate generation cost of P_t^{a*} via (5);
 - 23: $W_{a*} \leftarrow W_{a*} + N\bar{c}_t^{a*}$;
 - 24: $t \leftarrow t + 1$;
 - 25: **Until** $t = T + 1$
 - 26: **Else**
 - 27: **Repeat**
 - 28: Propose a power output $p_{i,t}$ of unit i ;
 - 29: **Repeat**
 - 30: Predict the average demand-generation mismatch α based on (3);
 - 31: Adjust $p_{i,t}$ according to (4);
 - 32: **Until** $\alpha \rightarrow 0$
 - 33: **If** $k = 1$ **then**
 - 34: Denote the local operation policy $\pi(D_t)$ as P_t ;
 - 35: $W_\pi \leftarrow W_\pi + N\bar{c}_t$;
 - 36: **Else**
 - 37: Choose power output as P_t ;
 - 38: Obtain immediate generation cost via (5);
 - 39: Update the parameter θ^t through (7);
 - 40: $W_p \leftarrow W_p + N\bar{c}_t$;
 - 41: Find the P_t^{a*} of (8) by (10);
 - 42: Obtain immediate generation cost of P_t^{a*} via (5);
 - 43: $W_{a*} \leftarrow W_{a*} + N\bar{c}_t^{a*}$;
 - 44: **End if**
 - 45: $t \leftarrow t + 1$;
 - 46: **Until** $t = T + 1$
 - 47: **End if**
 - 48: **Until** $t = T + 1$
 - 49: Update the local operation policy by (11);
 - 50: $W_\pi = \min\{W_{a*}, W_p, W_\pi\}$;
 - 51: **Until** $k = K$
 - 52: /* K is the maximum number of trials */
-

First of all, the equilibrium point of (10) with P_t^{a*} is analyzed to be the optimal solution of (8), and the convergence of (10) to the exact optimal solution of (8) is also proved. Denote

$$\mathcal{P}_t^{\text{new}} = \mathcal{P}_{1,t}^{\text{new}} \times \mathcal{P}_{2,t}^{\text{new}} \times \cdots \times \mathcal{P}_{N,t}^{\text{new}}$$

$$P_t = (p_{1,t}, p_{2,t}, \dots, p_{N,t})^T$$

$$\Lambda = (\lambda_1, \lambda_2, \dots, \lambda_N)^T$$

$$Z = (z_1, z_2, \dots, z_N)^T$$

$$\nabla f(P_t) = (\nabla f_1(p_{1,t}), \nabla f_2(p_{2,t}), \dots, \nabla f_N(p_{N,t}))^T.$$

Then, the compact form of (10) is

$$\begin{cases} \dot{P}_t = \mathbb{P}_{\mathcal{P}_t^{\text{new}}}(P_t - \nabla f(P_t) + \Lambda) - P_t \\ \dot{\Lambda} = -L\Lambda - LZ + \frac{1}{N}D_t \mathbf{1}_N - P_t \\ \dot{Z} = L\Lambda \end{cases} \quad (12)$$

The following theorem is given for the equilibrium point of (12), which indicates that P_t^{a*} in the equilibrium point $(P_t^{a*}, \Lambda^{a*}, Z^{a*})$ of (12) is corresponding to the optimal solution of (8).

Theorem 1: Suppose that Assumptions 1–3 hold and with the equilibrium point of distributed algorithm (12) with $(P_t^{a*}, \Lambda^{a*}, Z^{a*})$, then, P_t^{a*} is the optimal solution of (8).

Proof: By the property of the equilibrium point $(P_t^{a*}, \Lambda^{a*}, Z^{a*})$ of (12), we get the following equations.

1) $L\Lambda^{a*} = 0$, i.e., $\Lambda^{a*} = \lambda^{a*} \mathbf{1}_N$, $\lambda^{a*} \in \mathbb{R}$, because the undirected graph \mathcal{G} is connected.

2) $-L\Lambda^{a*} - LZ^{a*} + \frac{1}{N}D_t \mathbf{1}_N - P_t^{a*} = 0$, which implies

$$\text{that } D_t = \mathbf{1}_N^T P_t^{a*}, \text{ i.e., } \sum_{i=1}^N p_{i,t}^{a*} = D_t.$$

3) $\mathbb{P}_{\mathcal{P}_t^{\text{new}}}(P_t^{a*} - \nabla f(P_t^{a*}) + \Lambda^{a*}) - P_t^{a*} = 0$, which implies that $-\nabla f(P_t^{a*}) + \Lambda^{a*} \in N_{\mathcal{P}_t^{\text{new}}}(P_t^{a*})$.

Therefore, the equilibrium point $(P_t^{a*}, \Lambda^{a*}, Z^{a*})$ of (12) satisfies the KKT condition for (8)

$$\begin{cases} 0 \in \nabla f_i(p_{i,t}^{a*}) - \lambda^{a*} + N_{\mathcal{P}_t^{\text{new}}}(p_{i,t}^{a*}) \\ \sum_{i=1}^N p_{i,t}^{a*} = D_t \end{cases} \quad (13)$$

Hence, P_t^{a*} in the equilibrium point $(P_t^{a*}, \Lambda^{a*}, Z^{a*})$ of (12) is the optimal solution of (8).

Based on the above result, our next task is to prove that the trajectories of (12) with P_t will converge to the optimal solution P_t^{a*} .

Theorem 2: Under Assumptions 1–3, given the initial points $p_{i,t} \in \mathcal{P}_{i,t}^{\text{new}}$, $i \in 1, 2, \dots, N$, the trajectories of the algorithm of (12) are bounded and the power output $p_{i,t}$ of agent i converges to $p_{i,t}^{a*}$.

Proof: Denote $\bar{\mathcal{P}}_t^{\text{new}} = \mathcal{P}_t^{\text{new}} \times \mathbb{R}^N \times \mathbb{R}^N$. We define a new vector $M = \text{col}(P_t, \Lambda, Z)$ and the function $F(M) : \mathbb{R}^{3N} \rightarrow \mathbb{R}^{3N}$ as

$$F(M) = \begin{pmatrix} \nabla f(P_t) - \Lambda \\ L\Lambda + LZ - (\frac{1}{N}D_t \mathbf{1}_N - P_t) \\ -L\Lambda \end{pmatrix}. \quad (14)$$

Then, (12) can be written as $\dot{M} = \mathbb{P}_{\bar{\mathcal{P}}_t^{\text{new}}}(M - F(M)) - M$.

Define $H(M) = \mathbb{P}_{\bar{\mathcal{P}}_t^{\text{new}}}(M - F(M))$ and the dynamics become $\dot{M} = H(M) - M$. Consider the candidate Lyapunov function

$$\begin{aligned} V = & -\langle F(M), H(M) - M \rangle - \frac{1}{2} \|H(M) - M\|^2 \\ & + \frac{1}{2} \|M - M^{a*}\|^2 \end{aligned}$$

where $M^{a*} = \text{col}(P_t^{a*}, \Lambda^{a*}, Z^{a*})$ is the equilibrium point of (12). Via convex analysis and projection, we obtain

$$\begin{aligned} V = & -\langle F(M), H(M) - M \rangle - \frac{1}{2} \|H(M) - M\|^2 \\ & + \frac{1}{2} \|M - M^{a*}\|^2 \\ = & \frac{1}{2} [\|M - F(M) - M\|^2 - \|H(M) - (M - F(M))\|^2] \\ & + \frac{1}{2} \|M - M^{a*}\|^2 \\ \geq & \frac{1}{2} \|M - H(M)\|^2 + \frac{1}{2} \|M - M^{a*}\|^2. \end{aligned}$$

Hence, $V = 0$ if and only if $M = M^{a*}$. The derivative of V along (12) is

$$\begin{aligned} \dot{V} = & (F(M) - [J_F(M) - I](H(M) - M))^T (H(M) - M) \\ & + (M - M^{a*})^T (H(M) - M) \end{aligned} \quad (15)$$

where $J_F(M)$ is the Jacobian matrix of $F(M)$

$$J_F(M) = \begin{pmatrix} \nabla^2 f(P_t) & -I & 0 \\ I & L & L \\ 0 & -L & 0 \end{pmatrix} \quad (16)$$

which is positive semidefinite.

With the property of projection, it is obvious that $\langle M - F(M) - H(M), H(M) - M^{a*} \rangle \geq 0$, which implies $\langle M - H(M) - F(M), H(M) - M + M - M^{a*} \rangle \geq 0$.

Hence, $\langle H(M) - M, M - M^{a*} + F(M) \rangle \leq -\|H(M) - M\|^2 - \langle F(M), M - M^{a*} \rangle$. We may further get that

$$\begin{aligned} \dot{V} = & \langle M - M^{a*} + F(M), H(M) - M \rangle + \|H(M) - M\|^2 \\ & - (H(M) - M)^T J_F(M) (H(M) - M) \\ \leq & - (H(M) - M)^T J_F(M) (H(M) - M) \\ & - \langle F(M), M - M^{a*} \rangle \\ \leq & - \langle F(M), M - M^{a*} \rangle \\ \leq & - \langle F(M) - F(M^{a*}), M - M^{a*} \rangle - \langle F(M^{a*}), M - M^{a*} \rangle \\ \leq & 0. \end{aligned}$$

The last inequality holds because the Laplacian matrix is positive semidefinite, $f(P_t)$ is convex, and because of the variational inequality of the optimal solution M^{a*} . Therefore, there exists

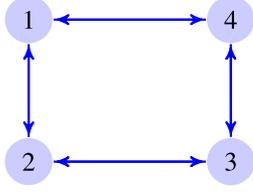


Fig. 1. Communication graph in Example 1.

TABLE I
PARAMETERS OF GENERATION UNITS

Unit number	\underline{p}_i	\bar{p}_i	a_i	b_i	c_i	e_i	f_i	p_i^R
1	200	600	0.0020	10	500	300	0.03	50
2	100	400	0.0025	8	300	200	0.04	50
3	100	300	0.0050	6	100	150	0.05	50
4	50	200	0.0060	5	90	130	0.06	50

a forward compact invariance set given as

$$IS = \{M \mid \frac{1}{2} \|M - M^{a*}\|^2 \leq V(M(0))\}.$$

From the KKT condition, there exist $p^{a*} \in N_{\mathcal{P}_t^{\text{new}}}(P_t^{a*})$ such that $p^{a*} = -\nabla f(P_t^{a*}) + \Lambda^{a*}$. Furthermore, we can obtain

$$\begin{aligned} \dot{V} &\leq -\langle F(M), M - M^{a*} \rangle \\ &= -\langle P_t - P_t^{a*}, \nabla f(P_t) - \Lambda - \nabla f(P_t^{a*}) \rangle \\ &\quad - \langle \Lambda - \Lambda^{a*}, L\Lambda + LZ - (\frac{1}{N}D_t \mathbf{1}_N - P_t^{a*}) \rangle \\ &\quad - \langle Z - Z^{a*}, -L\Lambda \rangle - \langle P_t - P_t^{a*}, \Lambda^{a*} - p^{a*} \rangle \\ &\leq -\langle P_t - P_t^{a*}, \nabla f(P_t) - \nabla f(P_t^{a*}) \rangle \\ &\quad + \langle P_t - P_t^{a*}, p^{a*} \rangle - \langle \Lambda - \Lambda^{a*}, L(\Lambda - \Lambda) \rangle \\ &\leq -\langle P_t - P_t^{a*}, \nabla f(P_t) - \nabla f(P_t^{a*}) \rangle \\ &\quad - \langle \Lambda - \Lambda^{a*}, L(\Lambda - \Lambda) \rangle. \end{aligned}$$

Denote the set $\mathcal{M} = \{M \mid \dot{V} = 0\}$. Because of the positive definite Hessian matrix $\nabla^2 f(P_t)$ and the null space for the Laplacian matrix L , we can obtain $\mathcal{M} = \{P_t = P_t^{a*}, \Lambda \in \text{span}\{\mathbf{1}_N\}\}$.

Next, we claim that the maximal invariance set within the set \mathcal{M} is the equilibrium point of (8). Because of $\Lambda \in \text{span}\{\mathbf{1}_N\}$, $Z = Z^{a*}$. According to (13), it is obvious that $\dot{\Lambda} = LZ^{a*} - (\frac{1}{N}D_t \mathbf{1}_N - P_t^{a*})$. We claim that $LZ^{a*} - (\frac{1}{N}D_t \mathbf{1}_N - P_t^{a*}) = 0$. Assume that $LZ^{a*} - (\frac{1}{N}D_t \mathbf{1}_N - P_t^{a*}) \neq 0$, then Λ will go to infinity, which contradicts that \mathcal{M} is a compact set within IS . Hence, $\dot{\Lambda} = 0$ and $\Lambda = \Lambda^{a*}$. By the LaSalle invariance principle, the power output $p_{i,t}$ of agent i converges to $p_{i,t}^{a*}$.

V. SIMULATION

In this section, the proposed distributed reinforcement learning optimization algorithm is tested through several examples.

Example 1: Consider four units connected via the undirected graph shown in Fig. 1. The cost function for each unit i is taken as $F_i(p_i) = a_i p_i^2 + b_i p_i + c_i + |e_i \cdot \sin(f_i \cdot (p_i - p_i))|$, with coefficients shown in Table I (known only to the purpose of

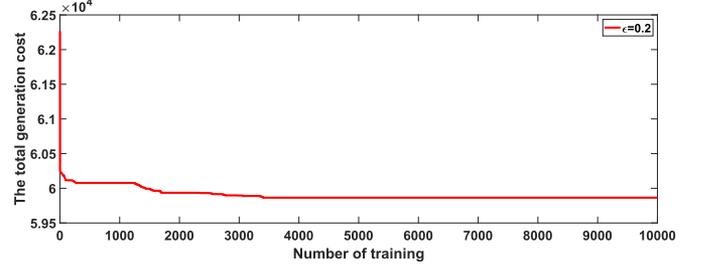
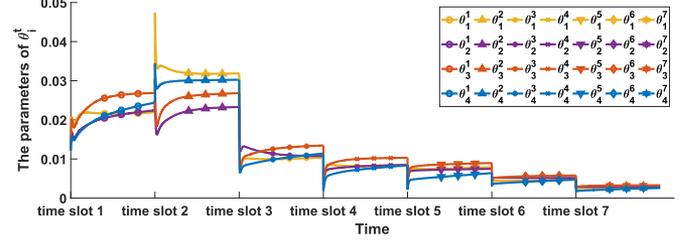
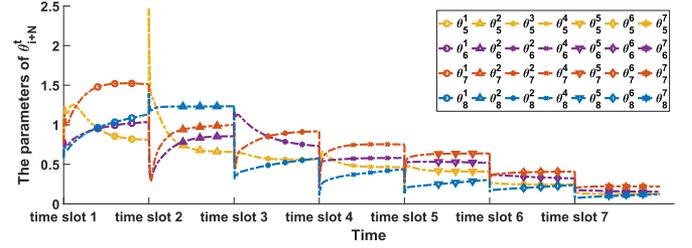


Fig. 2. Total generation cost of policy produced by the distributed reinforcement learning optimization algorithm in Example 1.

Fig. 3. Time-varying parameters θ_i^t in approximate state-action-value function.Fig. 4. Time-varying parameters θ_{i+N}^t in approximate state-action-value function.

simulation). The admissible power outputs of each unit are set as follows: $\mathcal{P}_1 = [200, 600]$, $\mathcal{P}_2 = [100, 400]$, $\mathcal{P}_3 = [100, 300]$, and $\mathcal{P}_4 = [50, 200]$ (MW). The total power demand D_t is 800, 850, 880, 900, 860, 930, and 950 (MW) for time periods $[0,2)$, $[2,6)$, $[6,8)$, $[8,18)$, $[18,22)$, and $[22,24)$, respectively.

We take for simplicity ε in the ε -greedy policy to be constant and equal to 0.2. As shown in Fig. 2, the total generation cost of updated policy is getting better and better during the training process. Figs. 3 and 4 show the time-varying parameters θ^t in approximate state-action-value functions for all time slots. In this example, the approximate state-action-value functions take the form $J_t(D_t, P_t, \theta^t) = \sum_{i=1}^N (\theta_i^t p_{i,t} + \frac{1}{4} \theta_{i+N}^t P_{i,t}^2)$. The optimal solutions P_t^{a*} of the approximate function for all time slots after training are shown in Fig. 5.

Remark 4: As the approximate state-action-value function $J_t(D_t, P_t, \theta^t)$ is the sum of total generation cost from time slot t to time slot T in the DED problem considered in this article. It can be seen from Figs. 3 and 4 that θ_i^t and θ_{i+N}^t are almost decreasing from time slot 1 to time slot T . Note that θ_i^t for time slot 2 is larger than θ_i^1 over time slot 1 which does not satisfy

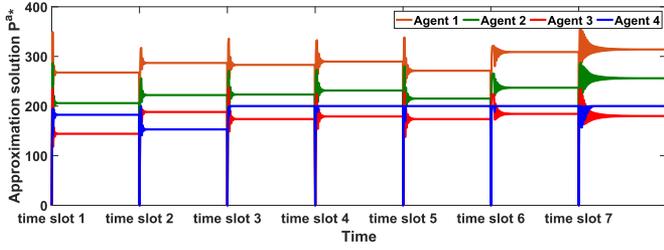


Fig. 5. P^{a*} of approximate state-action-value function after training.

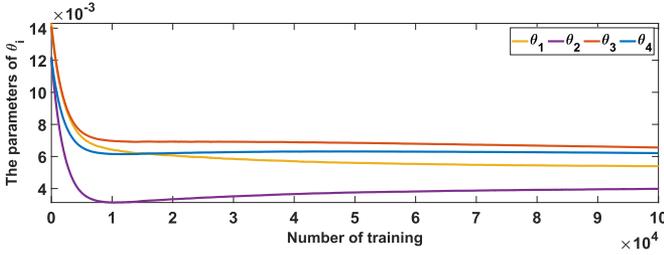


Fig. 6. Time-invariant parameters θ_i in approximate state-action-value function.

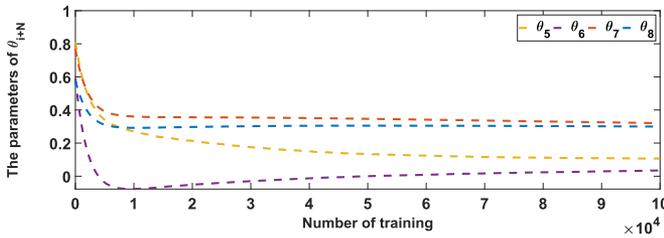


Fig. 7. Time-invariant parameters θ_{i+N} in approximate state-action-value function.

the property of decreasing; however, it has no effect according to the form of approximate state-action-value function.

In order to show the advantage of using time-varying parameters θ^t in the function approximation, a time-invariant parameter θ will be considered for all time slots. In other words, the approximate function takes the form $J(D_t, P_t, \theta) = \sum_{i=1}^N (\theta_i p_{i,t} + \frac{1}{4} \theta_{i+N} p_{i,t}^2)$. The parameters θ_i and θ_{i+N} are updated according to

$$\begin{cases} \theta_i \leftarrow \theta_i + \beta [N \bar{c}_t + \min_{P_{t+1}} J(D_{t+1}, P_{t+1}, \theta) \\ - J(D_t, P_t, \theta)] p_{i,t} \\ \theta_{i+N} \leftarrow \theta_{i+N} + \frac{\beta}{4} [N \bar{c}_t + \min_{P_{t+1}} J(D_{t+1}, P_{t+1}, \theta) \\ - J(D_t, P_t, \theta)] p_{i,t}^2 \end{cases} \quad (17)$$

Figs. 6 and 7 show the updating process. As shown in Fig. 7, θ_6 goes below zero, which contradicts the assumption of convexity of approximate state-action-value function. In this case, the step 5 cannot be performed as the necessary assumptions are violated.

Remark 5: By the definition of the state-action-value function, one gets that using time-invariant parameters θ for each time slot will cause severe fluctuations for the update of θ .

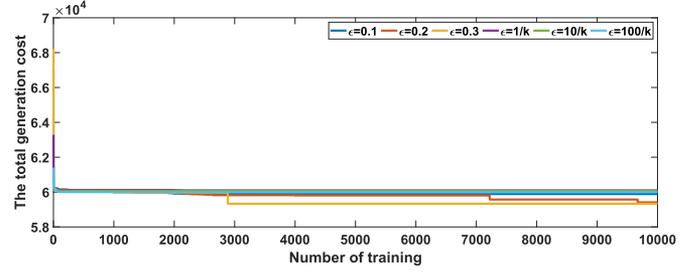


Fig. 8. Evolution of the total generation cost of updated policies in difference ε .

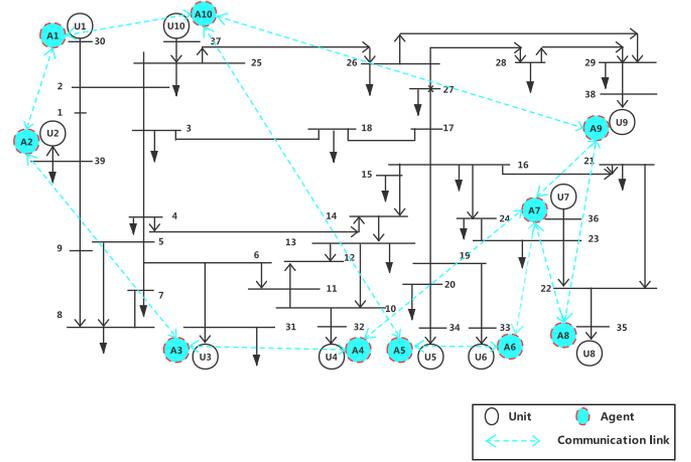


Fig. 9. IEEE 39-bus system.

The reinforcement learning optimization algorithm associated with time-varying parameters θ^t for each time slot t can reduce the concussion in the process of update of θ^t . It is also worth pointing out that using the time-varying parameter is also an efficient way when there exist the same FPOC in different time slots.

For the purpose of considering the effect of different ε in the ε -greedy policy, we take fixed $\varepsilon = 0.1$, $\varepsilon = 0.2$, and $\varepsilon = 0.3$, and also take $\varepsilon = \frac{1}{k}$, $\varepsilon = \frac{10}{k}$, and $\varepsilon = \frac{100}{k}$ which decreases gradually such that the operating policy is greedy limit with infinite exploration (GLIE) in ε -greedy. Fig. 8 shows the evolution of the total generation cost of each updated policy through 10,000 times training. As shown in Fig. 8, distributed reinforcement learning optimization algorithm yields a favorable policy when taking $\varepsilon = 0.3$.

Remark 6: It can be seen from the results given in Example 1 that the exploration in the distributed reinforcement learning optimization algorithm is very important as the number of FPOC is infinite in each time slot.

Example 2. We consider the IEEE-39 bus system with ten units. The communication network of these agents, which is described by the blue lines in Fig. 9, is undirected and connected. The cost function of each unit i is determined as $F_i(p_i) = a_i p_i^2 + b_i p_i + c_i$, where the coefficients are shown in Table II together with the minimum and maximum power

TABLE II
PARAMETERS OF UNITS

Unit number	a_i	b_i	c_i	p_i	\bar{p}_i	p_i^R
1	0.0072	5.56	30	60	339.69	50
2	0.0168	4.32	25	25	479.10	50
3	0.0216	6.60	25	28	290.4	50
4	0.0141	7.90	16	40	306.34	50
5	0.0273	7.54	6	35	593.80	50
6	0.0054	3.28	54	29	137.19	50
7	0.0159	7.31	23	45	595.40	50
8	0.0189	2.45	15	56	162.17	50
9	0.0084	7.63	20	12	165.1	50
10	0.0138	4.76	12	30	443.41	50

TABLE III
EXACT OPTIMAL SOLUTION FOR ALL TIME

Agent	P_1^*	P_2^*	P_3^*	P_4^*	P_5^*
1	332.01	339.68	339.69	339.69	339.69
2	179.47	195.27	212.33	229.38	246.42
3	87.01	99.12	112.38	125.64	138.90
4	87.58	105.82	126.13	146.45	166.77
5	51.78	61.27	71.79	82.31	92.82
6	137.18	137.18	137.19	137.19	137.19
7	96.35	112.44	130.48	148.53	166.58
8	162.16	162.16	162.17	162.17	162.17
9	163.43	165.09	165.10	165.10	165.10
10	202.98	221.90	242.71	263.52	284.32

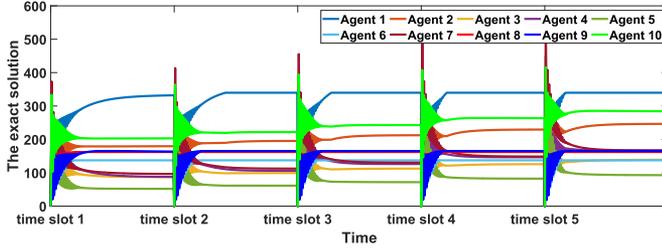


Fig. 10. Exact optimal solution.

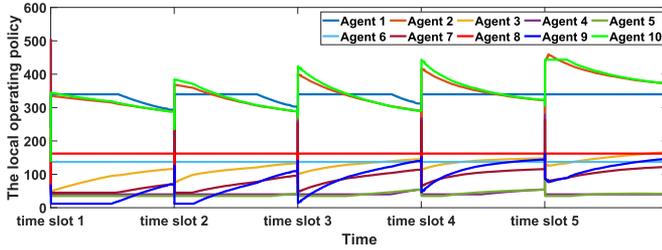


Fig. 11. Evolution of policy after 1087 times training.

generation of each unit. In this example, we consider the DED problem in five time slots. The power demand D_t is assumed to be 1500, 1600, 1700, 1800, and 1900 (MW) for time slot 1, 2, 3, 4, and 5, respectively. At first, consider the object function to be a quadratic convex function and the feasible set to be also a convex set. We use the distributed optimization algorithm based on the multiplier splitting method to find the exact optimal solution at time slot 1, 2, 3, 4, and 5 in Fig. 10. However, we do not know the form of the cost functions and the exact parameters in cost functions of units actually. Under this premise, we use the distributed reinforcement optimization algorithm to find the optimal policy. Fig. 11 shows that the evolution of the operating policy produced by the distributed reinforcement learning optimization algorithm after 1087 times training in this example. The exact optimal solution and the operating policy after 1087 times training are, respectively, shown in Tables III and IV. The error between exact optimal cost and the operating policy cost is less than 4% of exact optimal cost. In contrast to the ED problem studied in [34], the DED problem under consideration is more difficult due to the ramp-rate limit in each time slot.

TABLE IV
OPERATING POLICY AFTER 1087 TIMES TRAINING

Agent	$\pi(D_1)$	$\pi(D_2)$	$\pi(D_3)$	$\pi(D_4)$	$\pi(D_5)$
1	293.62	302.95	312.25	339.69	339.69
2	288.83	289.10	291.46	322.48	372.48
3	115.78	132.80	145.35	147.90	165.72
4	40	40	54.01	54.25	40
5	35	41.41	53.54	53.73	41.41
6	137.19	137.19	137.19	137.19	137.19
7	69.84	95.43	113.71	115.26	121.75
8	162.17	162.17	162.17	162.17	162.17
9	69.91	109.37	139.93	143.96	146.25
10	287.64	289.55	290.36	323.31	373.31

VI. CONCLUSION

In this article, we formulated a DED problem with little prior information of the generation cost functions in smart grid. To solve the DED problem, we combined the state-action-value function approximation and the distributed optimization algorithm based on multiplier splitting to get a distributed reinforcement learning optimization algorithm. Each step in the proposed algorithm was fully distributed. Theoretical analysis as well as case studies were presented to demonstrate the effectiveness of these proposed algorithms.

With respect to future works, the case that the total power demand D_{t+1} is decided by the feasible power output P_t at time slot t should be considered. Some constraints such as energy storage can also be considered in the future.

APPENDIX

A. Preliminaries on Convex Analysis

The following definitions and properties about convex set, convex function, and projection can be found in [32].

A set $\Omega \subset \mathbb{R}^n$ is called a convex set, if $\alpha x + (1 - \alpha)y \in \Omega \forall x, y \in \Omega \forall \alpha \in [0, 1]$. A function $f(\cdot) : \Omega \rightarrow \mathbb{R}$ is called to be a convex function, if $f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y) \forall x, y \in \Omega \forall \alpha \in [0, 1]$. If $f(\cdot) : \Omega \rightarrow \mathbb{R}$ is differentiable at $x \in \Omega$, its gradient, denoted by $\nabla f(x)$. $f(\cdot) : \Omega \rightarrow \mathbb{R}$, is called differentiable on Ω , if $f(x)$ is differentiable at any point $x \in \Omega$. Denote $N_\Omega(x)$ as the normal cone of Ω at x , that is, $N_\Omega(x) = \{y : \langle y, x' - x \rangle \leq 0 \forall x' \in \Omega\}$.

For a closed set Ω , define the projection of x onto Ω is $\mathbb{P}_\Omega(x) = \operatorname{argmin}_{y \in \Omega} \|x - y\|$. The common properties of

projection are as follows:

$$\langle x - \mathbb{P}_\Omega(x), \mathbb{P}_\Omega(x) - x' \rangle \geq 0 \forall x' \in \Omega \forall x \in \mathbb{R}^n$$

$$\|x - \mathbb{P}_\Omega(x)\|^2 + \|\mathbb{P}_\Omega(x) - x'\|^2 \leq \|x - x'\|^2 \forall x' \in \Omega \forall x \in \mathbb{R}^n.$$

Further, the normal cone $N_\Omega(x)$ can also be defined as $N_\Omega(x) = \{y : \mathbb{P}_\Omega(x + y) = x\}$.

B. Algebraic Graph Theory

The interaction topology of a system consisting of N units can be described by a graph. Let $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ be a graph with the set of nodes (i.e., units) $\mathcal{V} = \{1, 2, \dots, N\}$, the set of edges $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$. A directed edge $e_{ij} \in \mathcal{E}$ represents that node i can get the information from node j ; the graph \mathcal{G} is said to be undirected when $e_{ij} \in \mathcal{E}$ if and only if $e_{ji} \in \mathcal{E}$. The in-degree neighbors \mathcal{I}_i of node i is the set of nodes who can send their information to node i , i.e., $\mathcal{I}_i = \{j | e_{ij} \in \mathcal{E}\}$. A path is a sequence of distinct nodes in \mathcal{V} such that any consecutive nodes in the sequence correspond to an edge of graph. The undirected graph is connected, if there exists at least one path between any two nodes. The adjacency matrix A has the entries $a_{ij} = 1$ if $e_{ij} \in \mathcal{E}$, and $a_{ij} = 0$, otherwise. The Laplacian matrix $L = [l_{ij}]_{N \times N}$ of $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ is defined as

$$l_{ij} = \begin{cases} -a_{ij}, & i \neq j \\ \sum_{k=1, k \neq i}^N a_{ik}, & i = j \end{cases}.$$

Lemma 1: [7] Assume that the undirected graph \mathcal{G} is connected, the first-order discrete-time protocol

$$x[k+1] = x[k] - \epsilon Lx[k] \quad (18)$$

where $\epsilon \in (0, \frac{1}{\max_i l_{ii}})$, achieves asymptotic average consensus, i.e., $\lim_{k \rightarrow \infty} x_i[k] = \frac{1}{N} \sum_{i=1}^N x_i[0] \forall i \in \{1, 2, \dots, N\}$, where $x_i[k]$ is the i th element of $x[k]$.

C. Reinforcement Learning

Reinforcement learning is a framework of the problem of learning from interaction to achieve a goal. The learner is called the agent, which interacts with the environment by getting some immediate reward as a consequence of taking an action. Reinforcement learning with discrete states and actions is usually formulated as a Markov decision process (MDP). The MDP is defined as a tuple $\{\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \gamma\}$, where \mathcal{S} is the set of states, \mathcal{A} is the set of actions. $\mathcal{T} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ is the state transition function, $\mathcal{R} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ represents the reward function, and $\gamma \in [0, 1]$ is a discount factor. A policy $\pi : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ is a probability distribution over actions for each state. The state-action-value function $q_\pi(s, a)$ under policy π is defined as the expected discount of the long-term reward to the agent at the initial state s , taking action a , and then, following policy π . The aim of reinforcement learning is to find the optimal policy π^* . The policy π^* to maximize (minimize) cumulative reward is called to be the optimal policy, if $q_{\pi^*}(s, a) \geq q_\pi(s, a)$ (or $q_{\pi^*}(s, a) \leq q_\pi(s, a)$) $\forall s \in \mathcal{S}, a \in \mathcal{A} \forall \pi$. In standard reinforcement learning

problem, the environment is unknown, i.e., the transition function \mathcal{T} and reward function \mathcal{R} are unknown but static.

For large state and action spaces, function approximation in reinforcement learning is usually employed. Let $J(s, a, \theta)$ be an approximate function of the state-action-value function. We assume that $J(s, a, \theta)$ is a differential function of parameter vector θ for all $s \in \mathcal{S}, a \in \mathcal{A}$. The update of θ is as follows:

$$\theta \leftarrow \theta + \kappa \delta \nabla_\theta J(s, a, \theta)$$

where $\kappa \in (0, 1)$ and δ is the one-step temporal difference (TD) error given by

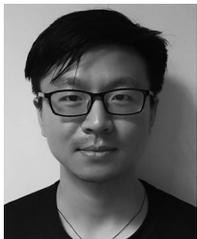
$$\delta = r + \gamma J(s', a', \theta) - J(s, a, \theta)$$

where r is immediate reward after taking action a on state s , γ is the discount factor, and (s', a') is state-action pair immediately after (s, a) .

REFERENCES

- [1] X. Fang, S. Misra, G. Xue, and D. Yang, "Smart grid—The new and improved power grid: A survey," *IEEE Commun. Surv. Tut.*, vol. 14, no. 4, pp. 944–980, 4Q 2012.
- [2] H. Farhangi, "The path of the smart grid," *IEEE Power Energy Mag.*, vol. 8, no. 1, pp. 18–28, Jan./Feb. 2010.
- [3] P. Siano, "Demand response and smart grids—A survey," *Renewable Sustain. Energy Rev.*, vol. 30, pp. 461–478, 2014.
- [4] V. C. Güngör *et al.*, "Smart grid technologies: Communication technologies and standards," *IEEE Trans. Ind. Inform.*, vol. 7, no. 4, pp. 529–539, Nov. 2011.
- [5] M. Pipattanasomporn, H. Feroze, and S. Rahman, "Multi-agent systems in a distributed smart grid: Design and implementation," in *Proc. IEEE/PES Power Syst. Conf. Expo.*, Seattle, WA, USA, 2009, pp. 1–8.
- [6] P. Gaj, J. Jasperneite, and M. Felser, "Computer communication within industrial distributed environment—A survey," *IEEE Trans. Ind. Inform.*, vol. 9, no. 1, pp. 182–189, Feb. 2013.
- [7] R. Olfati-Saber and R. M. Murray, "Consensus problems in networks of agents with switching topology and time-delays," *IEEE Trans. Autom. Control*, vol. 49, no. 9, pp. 1520–1533, Sep. 2004.
- [8] A. Nedic and A. Ozdaglar, "Distributed subgradient methods for multi-agent optimization," *IEEE Trans. Autom. Control*, vol. 54, no. 1, pp. 48–61, Jan. 2009.
- [9] A. Nedic, A. Ozdaglar, and A. P. Parrilo, "Constrained consensus and optimization in multi-agent networks," *IEEE Trans. Autom. Control*, vol. 55, no. 4, pp. 922–938, Apr. 2010.
- [10] S. Yang, S. Tan, and J.-X. Xu, "Consensus based approach for economic dispatch problem in a smart grid," *IEEE Trans. Power Syst.*, vol. 28, no. 4, pp. 4416–4426, Nov. 2013.
- [11] Z. Yang, J. Xiang, and Y. Li, "Distributed consensus based supply–demand balance algorithm for economic dispatch problem in a smart grid with switching graph," *IEEE Trans. Ind. Electron.*, vol. 64, no. 2, pp. 1600–1610, Feb. 2017.
- [12] T. Yang *et al.*, "A distributed algorithm for economic dispatch over time-varying directed networks with delays," *IEEE Trans. Ind. Electron.*, vol. 64, no. 6, pp. 5095–5106, Jun. 2017.
- [13] A. Cherukuri and J. Cortés, "Distributed generator coordination for initialization and anytime optimization in economic dispatch," *IEEE Trans. Control Netw. Syst.*, vol. 2, no. 3, pp. 226–237, Sep. 2015.
- [14] P. Yi, Y. Hong, and F. Liu, "Initialization-free distributed algorithms for optimal resource allocation with feasibility constraints and application to economic dispatch of power systems," *Automatica*, vol. 74, no. 1, pp. 259–269, Dec. 2016.
- [15] F. Guo, C. Wen, J. Mao, and Y.-D. Song, "Distributed economic dispatch for smart grids with random wind power," *IEEE Trans. Smart Grid*, vol. 7, no. 3, pp. 1572–1583, May 2016.
- [16] G. Binetti, A. Davoudi, D. Naso, B. Turchiano, and F. L. Lewis, "A distributed auction-based algorithm for the nonconvex economic dispatch problem," *IEEE Trans. Ind. Inform.*, vol. 10, no. 2, pp. 1124–1132, May 2014.

- [17] P. Yi, Y. Hong, and F. Liu, "Distributed gradient algorithm for constrained optimization with application to load sharing in power system," *Syst. Control Lett.*, vol. 83, no. 9, pp. 45–52, 2015.
- [18] G. Wen, X. Yu, Z. Liu, and W. Yu, "Adaptive consensus-based robust strategy for economic dispatch of smart grids subject to communication uncertainties," *IEEE Trans. Ind. Inform.*, vol. 14, no. 6, pp. 2484–2496, Jun. 2018.
- [19] W. Yu, C. Li, X. Yu, G. Wen, and J. Lü, "Economic power dispatch in smart grids: A framework for distributed optimization and consensus dynamics," *Sci. China Inf. Sci.*, vol. 61, no. 1, pp. 1–16, 2018.
- [20] C. Li, X. Yu, W. Yu, T. Huang, and Z.-W. Liu, "Distributed event-triggered scheme for economic dispatch in smart grids," *IEEE Trans. Ind. Inform.*, vol. 12, no. 5, pp. 1775–1785, Oct. 2016.
- [21] X. He, J. Yu, T. Huang, and C. Li, "Distributed power management for dynamic economic dispatch in the multimicrogrids environment," *IEEE Trans. Control Syst. Technol.*, vol. 27, no. 4, pp. 1651–1658, Jul. 2019.
- [22] W. Zhao, M. Liu, J. Zhu, and L. Li, "Fully decentralised multi-area dynamic economic dispatch for large-scale power systems via cutting plane consensus," *IET Gener., Transmiss. Distrib.*, vol. 10, no. 10, pp. 2486–2495, 2016.
- [23] G. Chen, C. Li, and Z. Dong, "Parallel and distributed computation for dynamical economic dispatch," *IEEE Trans. Smart Grid*, vol. 8, no. 2, pp. 1026–1027, Mar. 2017.
- [24] J. Qin, Y. Wan, X. Yu, F. Li, and C. Li, "Consensus-based distributed coordination between economic dispatch and demand response," *IEEE Trans. Smart Grid*, vol. 10, no. 4, pp. 3709–3719, Jul. 2019.
- [25] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*. Cambridge, MA, USA: MIT Press, 1998.
- [26] H. Zhang, Q. Wei, and Y. Luo, "A novel infinite-time optimal tracking control scheme for a class of discrete-time nonlinear systems via the greedy HDP iteration algorithm," *IEEE Trans. Syst., Man, Cybern. B, Cybern.*, vol. 38, no. 4, pp. 937–942, Aug. 2008.
- [27] H. Zhang, H. Jiang, Y. Luo, and G. Xiao, "Data-driven optimal consensus control for discrete-time multi-agent systems with unknown dynamics using reinforcement learning method," *IEEE Trans. Ind. Electron.*, vol. 64, no. 5, pp. 4091–4100, May 2017.
- [28] B. Wang, D. Zhao, C. Alippi, and D. Liu, "Dual heuristic dynamic programming for nonlinear discrete-time uncertain systems with state delay," *Neurocomputing*, vol. 134, pp. 222–229, 2014.
- [29] K. G. Vamvoudakis and F. L. Lewis, "Online solution of nonlinear two-player zero-sum games using synchronous policy iteration," *Int. J. Robust Nonlinear Control*, vol. 22, no. 13, pp. 1460–1483, 2012.
- [30] V. Mnih *et al.*, "Playing Atari with deep reinforcement learning," 2013, *arXiv:1312.5602*.
- [31] V. Mnih *et al.*, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [32] R. T. Rockafellar, *Convex Analysis*. Princeton, NJ, USA: Princeton Univ. Press, 1970.
- [33] A. J. Wood and B. F. Wollenberg, *Power Generation, Operation, and Control*. New York, NY, USA: Wiley, 2012.
- [34] F. Li, J. Qin, and Y. Kang, "Multi-agent system based distributed pattern search algorithm for non-convex economic load dispatch in smart grid," *IEEE Trans. Power Syst.*, vol. 34, no. 3, pp. 2093–2102, May 2019.



Pengcheng Dai (S'19) received the B.S. degree in statistics from Yancheng Normal University, Yancheng, China, in 2016, and the M.S. degree in applied mathematics in 2019 from Southeast University, Nanjing, China, where he is currently working toward the Ph.D. degree in applied mathematics.

His current research interests include distributed optimization and reinforcement learning.



Wenwu Yu (S'07–M'12–SM'15) received the B.Sc. degree in information and computing science and the M.Sc. degree in applied mathematics from the Department of Mathematics, Southeast University, Nanjing, China, in 2004 and 2007, respectively, and the Ph.D. degree in electronic engineering from the City University of Hong Kong, Hong Kong, in 2010.

He is currently the Founding Director of the Laboratory of Cooperative Control of Complex Systems and the Deputy Associate Director of the Jiangsu Provincial Key Laboratory of Networked Collective Intelligence, an Associate Director with the Research Center for Complex Systems and Network Sciences, an Associate Dean with the School of Mathematics, and a Full Professor with the Young Endowed Chair Honor in Southeast University, China. He has held several visiting positions in Australia, China, Germany, Italy, the Netherlands, and the USA. He has authored or coauthored about 100 science citation index SCI journal papers with more than 10 000 citations. His research interests include multi-agent systems, complex networks and systems, disturbance control, distributed optimization, neural networks, game theory, cyberspace security, smart grids, intelligent transportation systems, big-data analysis, etc.

Dr. Yu was the recipient of a National Natural Science Fund for Excellent Young Scholars in 2013, the National Ten Thousand Talent Program for Young Top-notch Talents in 2014, and the Cheung Kong Scholars Programme of China for Young Scholars in 2016. He was also the recipient of the Second Prize of State Natural Science Award of China in 2016. He is an Editorial Board Member of several flag journals, including the IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS II, IEEE TRANSACTIONS ON INDUSTRIAL INFORMATICS, IEEE TRANSACTIONS ON SYSTEMS, MAN, AND CYBERNETICS: SYSTEMS, *Science China Information Sciences*, *Science China Technological Sciences*, etc. He was listed as Highly Cited Researchers in Engineering by Clarivate Analytics/Thomson Reuters in 2014–2018.



Guanghui Wen (S'11–M'13–SM'17) received the Ph.D. degree in mechanical systems and control from Peking University, Beijing, China, in 2012.

He is currently a Professor with the Department of Systems Science, School of Mathematics, Southeast University, Nanjing, China. His current research interests include cooperative control of multiagent systems, analysis and synthesis of complex networks, cyber-physical systems, and resilient control.

Dr. Wen was the recipient of the Best Student Paper Award at the Sixth Chinese Conference on Complex Networks in 2010 and a National Natural Science Fund for Excellent Young Scholars in 2017. As coadvisor and coauthor, he has been a finalist for the IEEE International Symposium on Circuits and Systems (ISCAS) 2014 Best Student Paper Award. He was named a Highly Cited Researcher by Clarivate Analytics in 2018. He is a Reviewer for *American Mathematical Review* and is an active Reviewer for many journals. He is currently an Editorial Board Member and an Associate Editor for the IEEE TRANSACTIONS ON SYSTEMS, MAN, AND CYBERNETICS: SYSTEMS and the *Asian Journal of Control*.



Simone Baldi (M'14) received the B.Sc. degree in electrical engineering, and the M.Sc. and Ph.D. degrees in automatic control systems engineering from the University of Florence, Florence, Italy, in 2005, 2007, and 2011, respectively.

He is currently a Professor with the School of Mathematics, Southeast University, Nanjing, China, with a Guest Position with the Delft Center for Systems and Control, Delft University of Technology, Delft, The Netherlands, where he was an Assistant Professor. Previously, he was a Postdoctoral Researcher with the University of Cyprus, and the Information Technologies Institute, Centre for Research and Technology Hellas. His research interests include adaptive and learning systems with applications in networked control systems, smart energy, and intelligent vehicle systems.

Prof. Baldi was the recipient of the Outstanding Reviewer Award of *Applied Energy* (2016), *Automatica* (2017), and *IET Control Theory and Applications* (2018). Since March 2019, he has been the Subject Editor for the *International Journal of Adaptive Control and Signal Processing*.