

---

# Rethink AI-based Power Grid Control: Diving Into Algorithm Design

---

Xiren Zhou<sup>1</sup>, Siqi Wang<sup>2</sup>, Ruisheng Diao<sup>2</sup>, Desong Bian<sup>2</sup>, Jiajun Duan<sup>2</sup> and Di Shi<sup>2</sup>

<sup>1</sup>Columbia University

<sup>1</sup> xz2754@columbia.edu

<sup>2</sup>Global Energy Interconnection Research Institute North America(GEIRINA)

<sup>2</sup>{siqi.wang, ruisheng.diao, desong.bian, jiajun.duan, di.shi}@geirina.net

## Abstract

Recently, deep reinforcement learning (DRL)-based approach has shown promise in solving complex decision and control problems in power engineering domain. In this paper, we present an in-depth analysis of DRL-based voltage control from aspects of algorithm selection, state space representation, and reward engineering. To resolve observed issues, we propose a novel imitation learning-based approach to directly map power grid operating points to effective actions without any interim reinforcement learning process. The performance results demonstrate that the proposed approach has strong generalization ability with much less training time. The agent trained by imitation learning is effective and robust to solve voltage control problem and outperforms the former RL agents.

## 1 Introduction

Nowadays, the rapid development of artificial intelligence (AI) technologies provides new ideas and solutions for solving many challenges in the field of power grid operation and control. The application of deep reinforcement learning has been extensively explored to solve complex power engineering problems, such as grid emergency control [1], real-time autonomous energy management [2] and topology adjustment [3]. Specifically, DRL-based power grid control paradigm has become a hot spot and provides effective future development direction for both power system research and engineering community. The pioneer work in [4] presents a novel autonomous control paradigm called "Grid Mind" to derive fast and effective controls in real time with Deep Q Network (DQN) agent to eliminate voltage violations. Later on, the follow-up work in [5] expands the findings of [4] and modified the DQN algorithm to improve control performance by avoiding the choice of same actions multiple times and normalizing the observations. Similar ideas have been further expanded in [6] and [7] for controlling the voltage setpoint of generators and PV-converters in a continuous manner by employing Deep Deterministic Policy Gradient (DDPG) agent. However, the previous work is mainly focused on demonstrating that the state-of-art DRL algorithm can be plugged in to the power grid control framework. Different from the previous exploration, in this paper, we investigate the underlying nature of power grid voltage control and revisit it from the Markov Decision Process (MDP) point of view, mainly focusing on tuning the design of the problem formulation, i.e., the choice of state representation, the DRL algorithm selection, and the reward engineering. Moreover, we present the lessons learnt from training an effective DRL agent using the real world power grid data collected from the control center of SGCC Jiangsu Electric Power Company. We further propose an imitation learning-based approach to resolve the voltage control problem based upon the findings of the training process and demonstrate that the imitation learning approach is also an effective alternative without any reinforcement learning component for this type of voltage control problem.

## 2 Problem formulation

The main objective of power system voltage control is to maintain bus voltage profiles within the predefined bounds and at the same time keeping transmission line flows within limits. The state is the current operating condition of the power grid represented by the power flow snapshot at a given timestamp. We model the power grid control problem as an MDP as follows:

- $S$ : an infinite state space of continuous-valued state representation.  $S \in \mathbb{R}^{n_s}$ .  
Three types of measurement values can be adopted to construct state space: bus values (bus voltage  $V_m$  and bus angle  $V_a$ ), branch values (line flow  $S_{line}$ ), and generator values (active power  $P_g$  and reactive power  $Q_g$ ).
- $A$ : an infinite action space of continuous-valued action vector.  $A \in \mathbb{R}^{n_a}$ .  
The voltage setpoints of generators within each power plant are used to control bus voltages. Therefore  $n_a$  represents the number of power plants in the studied power grid. The plant voltage value is bounded within the range of [0.9, 1.1] in p.u.
- $P$ : a transition dynamics model that specifies  $P(s'|s, a)$ .  
In essence, the transition probability  $P(s'|s, a)$  is determined and dominated by the physics laws and all configurations of its underlying environment. In this design, we use an in-house power flow solver (simulator) to model the transitions from state  $s$  to state  $s'$  after applying action  $a$ , i.e.,  $P(s'|s, a) = 1$  where  $s'$  is obtained from the power flow solver given the input  $s$  and  $a$ . Thus,  $P$  is a sparse/one-hot matrix (of infinite dimension) in our MDP.
- $R$ : a reward function that maps a state action pair to a real number.  
Due to the sparse property of  $P$  mentioned above, the reward at time step  $t$  can also be determined by the resulting state at time step  $t + 1$ . i.e.,  $r_t = R(s_t, a_t) = R(s_{t+1})$ .
- $\gamma$ : the discount factor. It is fixed to 0.99 throughout the paper.

Moreover, we define a set of terminal states  $T$  which is composed of all the states without voltage violation for all buses and without line flow violation for all lines. A state  $s$  is "unsuccessful" iff  $s \notin T$ , meaning that there exists either voltage violations or line flow violations at the current operating condition. Likewise, we define any  $s \in T$  as a "successful" state. The MDP is episodic by defining terminal states as above and enforcing a horizon limit for each episode (e.g., 50, 100, 1000, etc.).

In summary, the ultimate goal for the DRL agent is to find a policy  $\pi$  (a mapping from state to action) that can eliminate both voltage violations and line flow violations as quickly as possible. In other words, given an unsuccessful initial state, the policy is desired to be able to make effective sequence of decisions that leads to a successful state in as few number of steps as possible.

## 3 Data preparation

We collected 10,433 power flow snapshots from the control center of SGCC Jiangsu Power Company, which represents the power grid operating conditions from January to March, 2020. The snapshots are stored in \*.dat format, each of which can be treated as an unsuccessful initial state  $s_0$ . To eliminate the seasonal impact of co-relation, the entire dataset is randomly split into 9,433 and 1,000 for training and testing, respectively, i.e.:

$$S_{\text{train}} = \{s_0^{(1)}, s_0^{(2)}, \dots, s_0^{(9433)}\}; S_{\text{test}} = \{s_0^{(9434)}, s_0^{(9435)}, \dots, s_0^{(10433)}\} \quad (1)$$

## 4 RL-based power grid voltage control

### 4.1 Algorithm selection

Both of value based RL algorithm (i.e., DQN) and policy gradient based RL algorithm (i.e., DDPG) have been studied in previous work to some extent. But policy gradient RL algorithms tend to be more appropriate and effective in solving power grid control problem for the following reasons: 1) policy gradient based algorithms are proven to converge to at least local optimum (under certain conditions) whereas DQN does not have such guarantee mathematically [8]; 2) an ordinary DQN is difficult to

deal with infinite continuous action space while both action space and state space of voltage control are continuous. In this work, we consider soft actor critic (SAC) as our policy gradient RL algorithm instead of DDPG to enable more incentives for randomness in actions, which is important in finding a good plant voltage point through exploration [9].

## 4.2 Reward design strategies

As mentioned before, the goal of training RL agents is to move from an unsuccessful state  $s_0$  to a successful state in as few steps as possible. However, in SAC formulation, the agent is optimized directly upon the discounted sum of rewards instead of the number of steps per episode. Therefore, a well-designed reward function plays a critical role in training an effective agent. The agent must have the capability to learn lessons efficiently from the variation of the rewards to shorten its number of steps within an episode.

Intuitively, the essence of the policy gradient formula is to increase the probability of taking the actions that receive large rewards and decrease the probability of taking the actions that receive small rewards. A good reward function design should be able to direct the policy network effectively through rewards and penalties (i.e. small positive rewards or even negative rewards). With this assumption, we define our reward function  $R$  by dividing it up into two separate functions  $R_-$  and  $R_+$ , according to the types of a transition step (either successful or unsuccessful). i.e.,

$$r_t = R(s_t, a_t) = \begin{cases} R_-(s_{t+1}), & \text{if } s_{t+1} \notin T \\ R_+(s_{t+1}), & \text{if } s_{t+1} \in T \end{cases} \quad (2)$$

And we propose the following two reward design strategies:

1.  $R_-(s) = f_{\text{penalty}}(V_m, S_{\text{line}})$ ;  
 $R_+(s)$  is a fixed non-negative constant.
2. CartPole-style reward:  $R \equiv -1$ .

where  $f_{\text{penalty}} = \alpha \sum_i \text{line\_overflow}[i] + \beta \sum_j \text{bus\_violation}[j]$  computes the penalty (negative reward) of a given state according to its  $S_{\text{line}}$  and  $V_m$ .

The overflow of the  $i$ th line and the voltage violation of the  $j$ th bus are defined as following:

$$\begin{aligned} \text{line\_overflow}[i] &= \max\{S_{\text{line}}[i] - \text{line\_limit}[i], 0\}^2 \\ \text{bus\_violation}[j] &= \max\{(V_m[j] - \text{bus\_lower\_limit}[j])(V_m[j] - \text{bus\_upper\_limit}[j]), 0\} \end{aligned} \quad (3)$$

$\alpha$  and  $\beta$  are used to balance the importance between line overflow and bus voltage violation, which are set to  $-0.1$  and  $-1000$  respectively throughout this paper.

Since the optimization objective is the average number of steps needed to solve an unsuccessful state, which is basically the length of the episode, we attempt to make the agent learn directly from it. Hence, strategy 2 is proposed to make the length of an episode reflected by the return value. We call it "CartPole-style" reward since this is similar to the reward style of the CartPole environment [10].

It is the goal that under either of the two strategies above, the agent will converge to a theoretically optimal policy that solves all voltage violation cases in one step. However, it is difficult to achieve in real world, because power grid is such a complex system that we are not able to theoretically ensure it is Markovian. Meanwhile, for some voltage violation cases, there might not exist any one-step solution, given that we only enable control over plant voltage value settings.

In the next section, we demonstrate the experiment results of different reward function designs, where we end up with a best reward design pattern and obtain some insights into power grid MDP.

## 4.3 Experiments & analysis

As shown in Table 1, we tried different values for  $R_+$  for reward design strategy 1. A higher  $R_+$  basically makes the training more efficient and converge to an optimal policy more quickly. The SAC algorithm with  $R_+ = 0$  fails to find an optimal policy.

Table 1: Training time needed til finding the optimal policy under different  $R_+$  for strategy 1

$R_+$	0	1	10~20	50	80~1000
training steps	fail to converge	3.4k	1.3k	1k	0.9k

For strategy 2, with  $R \equiv -1$ , the length of an episode is directly reflected by its sum of reward. However, it fails to converge to an optimal policy. We then adjusted  $R_+$  to 0 or 1, which is different from the original "CartPole-style" setting, since it is inappropriate to give the agent a penalty even on a successful step. Unfortunately, the training diverges as well. We finally set a large value (1000) for  $R_+$ , which eventually works well.

Table 2: Training time needed til finding the optimal policy under different  $R_+$  for strategy 2

$R_+$	-1	0	1	1000
training steps	fail to converge	fail to converge	fail to converge	0.9k

From these experiments, we find out that a higher positive reward makes the agent learn lessons more quickly and efficiently. As for the negative reward on any unsuccessful step, according to formula 3, the penalty increases quadratically as the line flow or bus voltage goes beyond the limit, which is a common design in this field. However, It turns out that a well-designed meaningful  $R_-$  does not play a such significant role as  $R_+$ ; even with  $R_-$  as a constant, as long as  $R_+$  is large enough, the agents can learn very well.

These observations imply that the agent learns significantly from the last successful steps. Those unsuccessful steps give little helpful information to the RL agent. It inspires us to train the agent with only successful steps and throw away all unsuccessful steps, which leads to the imitation learning based method that will be covered in the next section.

## 5 Imitation learning-based power grid control

**Algorithm 1** Imitation Learning for Training a Power Grid Control Agent

---

<p>Initialize: policy network <math>\pi</math> with random weights <math>\theta</math>  <math>D = \text{COLLECT SUCCESSFUL STEPS}(S_{\text{train}})</math>  <math>l_{\text{objective}} = 1</math> # set an objective average episode's length to terminate the runtime.  <b>while</b> True <b>do</b>              Train <math>\theta</math> with an optimizer (eg., Adam) on dataset <math>D</math> for an epoch              <math>l = \text{EVALUATE EPISODE LENGTH}(S_{\text{test}}, \pi_\theta)</math>              <b>if</b> <math>l \leq l_{\text{objective}}</math> <b>then</b>                  <b>terminate</b></p> <p><b>procedure</b> COLLECT SUCCESSFUL STEPS(<math>S</math>)              Initialize <math>D = \emptyset</math>              Set policy <math>\pi_{\text{collect}}</math> to be an arbitrary policy (eg., random policy, trained SAC policy, etc.)              <math>n =</math> number of successful steps to collect (eg., 10000)              <math>t_{\text{limit}} = 1000</math> # set a horizon limit              <b>while</b> <math> D  &lt; n</math> <b>do</b>                  Randomly sample a state <math>s_0</math> from <math>S</math>                  <b>for</b> <math>t = 0, 1, 2, \dots, t_{\text{limit}} - 1</math> <b>do</b></p>	<p><math>a_t = \pi_{\text{collect}}(s_t)</math>  <math>s_{t+1} \leftarrow \text{perform } a_t \text{ on } s_t</math>  <b>if</b> <math>s_{t+1} \in T</math> <b>then</b>              <math>D = D \cup \{(s_t, a_t)\}</math>              <b>break</b>  <b>return</b> <math>D</math></p> <p><b>procedure</b> EVALUATE EPISODE LENGTH(<math>S, \pi</math>)              Initialize <math>L = \emptyset</math>              <math>n_{\text{episodes}} = 50</math> # set number of episodes(cases) to evaluate              <math>t_{\text{limit}} = 50</math> # set a horizon limit              <b>for</b> <math>i = 1, 2, \dots, n_{\text{episodes}}</math> <b>do</b>                  Randomly sample a state <math>s_0</math> from <math>S</math>                  <b>for</b> <math>t = 0, 1, 2, \dots, t_{\text{limit}} - 1</math> <b>do</b>                      <math>a_t = \pi(s_t)</math>                      <math>s_{t+1} \leftarrow \text{perform } a_t \text{ on } s_t</math>                      <b>if</b> <math>s_{t+1} \in T</math> <b>or</b> <math>t + 1 = t_{\text{limit}}</math> <b>then</b>                          <math>L = L \cup \{t + 1\}</math>                          <b>break</b>              <b>return</b> <math>\mathbb{E}_{l \in L}[l]</math></p>
--	--

---

Inspired by the study in Section 4.3, an imitation learning method is proposed to train an agent with only successful steps. Specifically, the proposed method does not incorporate any reinforcement

learning component except a policy neural network mapping from a state vector to an action vector. The network is trained by supervised learning on a dataset  $D$ :

$$D = \{(s_t, a_t) | \text{where } s_{t+1} \in T\} \quad (4)$$

## 6 Performance & analysis

### 6.1 Random policy baseline performance

For each  $s_0 \in S_{\text{train}} \cup S_{\text{test}}$ , a random agent is employed to interact with the environment. The simulation is conducted multiple times (episodes) for each  $s_0$  and the average number of steps needed to solve each case are recorded. The maximum episode length is 1000. If a case can never be solved within 1000 steps, it is considered as an "unsolvable" case. Figure 1(a) shows the random policy performance on the entire dataset of 10,433 cases. There are totally 142/10433 (1.36%) unsolvable cases.

### 6.2 SAC agent performance

The SAC agent was trained with  $R_+ = 1000$ . The evaluation result of the agent is shown in Figure 1 and Table 3. We evaluated the agent's policy with its stochastic (i.e., the action is sampled from a normal distribution given its learned mean and standard deviation) and greedy (i.e., the action is solely determined by its mean, which means a deterministic policy) variants respectively. The greedy SAC policy is able to solve all cases in only one step; but there exist more unsolvable cases, compared with random policy. The stochastic SAC policy, on the other hand, has a higher solvable rate (98.5% on train set and 98.4% on test set) but needs more steps on average to solve a voltage violation case. One reason is that the state representation is still possibly not a Markovian representation (considering the complexity of power grid system), and adding randomness can alleviate the partial observability problem [11]. This is why a policy with distributed action or a random policy is able to solve more cases (regardless of steps cost) than a deterministic policy. Despite the advantage of stochasticity, the deterministic policy is superior in that it solves a case instantly (one step). In practice, it relies on the grid operators' engineering judgement to decide which policy is better.

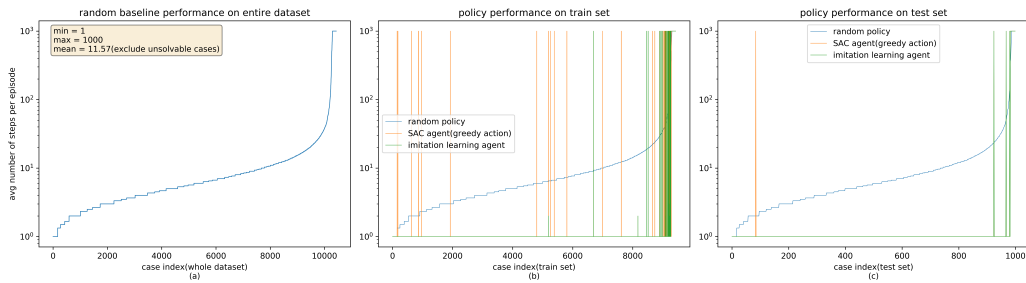


Figure 1: Policy performance comparison. (Case indices are sorted according to random policy performance. A case is considered unsolvable if its y value reaches  $10^3$ .)

### 6.3 Imitation learning agent performance

A dataset  $D$  is collected according to Algorithm 1 by running random policy, which contains  $\sim 20k$  successful transitions corresponding to  $S_{\text{train}}$ . In other words, any of the successful transition is generated within an episode starting with a certain initial state in  $S_{\text{train}}$ . Note that there might be multiple transitions in  $D$  corresponding to the same  $s_0$  in  $S_{\text{train}}$ . The network architecture is the same as the policy network of the former SAC agent (see Appendix). The training takes only three epochs until finding the optimal policy, which is much faster than the SAC algorithm. As shown in Figure 1 and Table 3, the imitation learning agent performs even better than the SAC agent with a higher solvable rate. It can solve any solvable case in one step. The strong generalization ability of the imitation learning agent can be viewed in two aspects: 1)  $D$  was collected by exploring cases in  $S_{\text{train}}$ , and the agent is able to solve most cases in  $S_{\text{test}}$  in one step; 2) consider  $(s, a) \in D$  which is the successful step of an episode starting with  $s_0^{(i)} \in S_{\text{test}}$ , since  $D$  was collected by running a random

Table 3: Policy performance details

Policy	Number of unsolvable cases	Avg steps to solve a case
random policy	<b>train set: 127/9433(1.35%)</b> <b>test set: 13/1000(1.3%)</b>	train set: 11.53 test set: 11.98
SAC agent(normal distributed action)	train set: 138/9433(1.46%) test set: 16/1000(1.6%)	train set: 3.17 test set: 2.91
SAC agent(greedy action)	train set: 226/9433(2.40%) test set: 22/1000(2.2%)	<b>train set: 1</b> <b>test set: 1</b>
imitation learning agent	train set: 196/9433(2.08%) test set: 21/1000(2.1%)	<b>train set: 1</b> <b>test set: 1</b>

policy, mostly  $s \neq s_0^{(i)}$  (the case is not solved within one step.) E.g., say  $s = s_5^{(i)}$ , training with  $(s_5^{(i)}, a_5^{(i)})$ , the agent is able to learn an  $a_0^{(i)}$  which instantly solves this case in one step.

### 6.4 PCA analysis of the state space

As demonstrated in Figure 1, and in Table 3, there are a small number of extreme cases that are unsolvable no matter what policy it is. In order to have a deep understanding about the states of the environment, the principle component analysis (PCA) is conducted to investigate the cluster distribution of unsolvable cases versus solvable cases in Euclidean space. Specifically, for  $\forall s \notin T$ , a random agent is initiated with allowable maximum steps as 1000 for multiple times. If  $s$  can never be solved within 1000 steps during multiple trials, we consider  $s$  as an "unsolvable" case (state); otherwise, it's a "solvable" case (state).

The full state has a dimension of 1,442. Figure 2(a) verifies that the first two components can explain the majority of the variance in the data and Figure 2(b) shows there is a clear boundary between unsolvable cases and solvable cases. This fact corresponds to the fact in complex power grid that there exists a small amount of extreme cases which can be proven unsolvable in the case of adjusting plant voltage setpoint only.

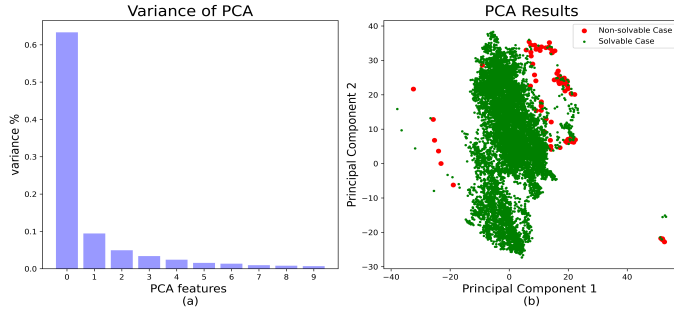


Figure 2: PCA results: (a) Variations (b) Non-solvable v.s Solvable case

## 7 Conclusion

In this work, we revisited the previous DRL-based voltage control problem of power grid. We performed an in-depth analysis on algorithm selection, state space representation, and reward engineering. Based upon the analysis result, we realize that the agent mostly learn lessons from the positive rewards of the last successful steps. Thus, we optimize the reward design which results in a sample-efficient SAC-based approach that converges to the optimal policy very quickly. Furthermore, we proposed a novel imitation learning-based approach to perform power grid voltage control. The training and testing results show that the trained imitation learning agent has strong generalization ability which even outperforms the RL agent with the same policy network architecture. Meanwhile, the imitation learning based method does not involve any complex hyper-parameter tuning or design of a reward function, and requires less training time to converge to the optimal policy.

## Broader Impact

The rapid development of renewable energy brings more and more complexity to the control of power grid. The traditional control methods based on operator's experience are difficult to cope with the complex and changeable power grid operating conditions (especially under unknown operating conditions). How to effectively realize rapid regulation of power grid is an urgent problem to be solved. With the fast development of AI technology, it is promising to provide operators with accurate and timely control plans, to improve the control efficiency during incidents. This work provides an in-depth analysis on the DRL-based methodologies for autonomous voltage control for power grid operation. Key aspects of MDP are thoroughly investigated for improving the performance of RL agents. During this process, we found that such MDP formulation aims at obtaining one-step control to fix voltage violations once detected; then we propose an imitation learning-based approach to achieve this goal, the effectiveness of which is verified via massive simulation studies conducted on actual power grid operating conditions. It is our hope that this work can help research community in better understanding the underlining principles of such control problems and promote AI-based solutions towards real-world implementation.

## References

- [1] Q. Huang, R. Huang, W. Hao, J. Tan, R. Fan, and Z. Huang, "Adaptive power system emergency control using deep reinforcement learning," *IEEE Trans. Smart Grid*, vol. 11, no. 22, pp. 1171-1182, 2020.
- [2] Y. Ye, D. Qiu, X. Wu, G. Strbac and J. Ward, "Model-Free real-time autonomous control for a residential multi-energy system using deep reinforcement learning," *IEEE Trans. Smart Grid*, vol. 11, no. 4, pp. 3068-3082, July 2020.
- [3] T. Lan, J. Duan, B. Zhang, D. Shi, Z. Wang, R. Diao, and X. Zhang, "AI-based autonomous line flow control via topology adjustment for maximizing time-series ATCs," <https://arxiv.org/abs/1911.04263>, 2019.
- [4] R. Diao, Z. Wang, D. Shi, etc., "Autonomous voltage control for grid operation using deep reinforcement learning," *IEEE PES General Meeting*, Atlanta, GA, USA, 2019.
- [5] B. L. Thayer and T. J. Overbye, "Deep reinforcement learning for electric transmission voltage control," *arXiv preprint arXiv:2006.06728*, 2020.
- [6] J. Duan et al., "Deep reinforcement learning-based autonomous voltage control for power grid operations," *IEEE Trans. Power Systems*, vol. 35, no. 1, pp. 814-817, Jan. 2020.
- [7] C. Li, C. Jin, and R. K. Sharma, "Coordination of pv smart inverters using deep reinforcement learning for grid voltage regulation," *2019 18th IEEE International Conference On Machine Learning And Applications (ICMLA)*, pp. 1930-1937, 2019
- [8] R. S. Sutton, D. Mcallester, S. Singh, and Y. Mansour. "Policy gradient methods for reinforcement learning with function approximation." In *NIPS*, pages 1057-1063. MIT Press, 2000.
- [9] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International Conference on Machine Learning (ICML)*, volume 80, pages 1856-1865, 2018.
- [10] G. Brockman et al., "OpenAI Gym", *arXiv:1606.01540*, 2016
- [11] Singh, Satinder, T. Jaakkola and Michael I. Jordan. "Learning without state-estimation in partially observable markovian decision processes." *Proc. ICML-94*, pp. 284-292.

## Appendix:

All code implementations are based on TensorFlow 2.3.0. The SAC implementation is based on TF-Agents 0.6.0. The hyperparameters used for training are shown in Table 5 and 6.

Table 4: Variables and physical quantities

Name	Description	Unit or value
$V_m$	voltage magnitude	p.u
$V_a$	voltage angle	p.u
$S_{line}$	line flow apparent power	p.u
$P_g$	active power of generator	MW
$Q_g$	reactive power of generator	MVar
$n_s$	state dimension	1442
$n_a$	action dimension(number of plants)	16

Table 5: Hyperparameters used for training SAC agents

Hyperparameter	Description	Value
critic_obs_fc	fully connected layers for observation in critic network	[512, 512]
critic_action_fc	fully connected layers for action in critic network	[256]
critic_joint_fc	fully connected layers after merging observations and actions	[256, 256]
actor_fc	fully connected layers for policy network	[512, 512]
batch_size	batch size	256
episode_max_len	horizon limit during training	50
lr	learning rate	0.0003
$\tau_{target}$	Factor for soft update of the target networks	0.005
target_update_period	Period for soft update of the target networks	1
initial $\log(\alpha)$	initial value of $\log(\alpha)$	1
collect_episodes	number of episodes to collect per training step	1

Table 6: Hyperparameters used for imitation learning

Hyperparameter	Description	Value
hidden_fc_layers	fully connected layers	[512, 512]
batch_size	batch size	32~512
lr	learning rate	0.001