Energy Storage Arbitrage in Real-Time Markets via Reinforcement Learning

Hao Wang, Baosen Zhang

Department of Electrical Engineering, University of Washington, Seattle, WA 98195 Email: {hwang16,zhangbao}@uw.edu

Abstract

In this paper, we derive a temporal arbitrage policy for storage via reinforcement learning. Real-time price arbitrage is an important source of revenue for storage units, but designing good strategies have proven to be difficult because of the highly uncertain nature of the prices. Instead of current model predictive or dynamic programming approaches, we use reinforcement learning to design an optimal arbitrage policy. This policy is learned through repeated charge and discharge actions performed by the storage unit through updating a value matrix. We design a reward function that does not only reflect the instant profit of charge/discharge decisions but also incorporate the history information. Simulation results demonstrate that our designed reward function leads to significant performance improvement compared with existing algorithms.

I. INTRODUCTION

Energy storage can provide various services (e.g., load shifting, energy management, frequency regulation, and grid stabilization) [1] to the power grid and its economic viability is receiving increasing attention. One of the most discussed revenue sources for energy storage is *real-time temporal arbitrage* (i.e., charging at low prices and discharging at higher prices), where storage units take advantage of the price spreads in real-time electricity market prices to make profits over time [2]. This application has received significant attention from the research community, especially since the growing penetration of intermittent renewable generations are resulting in more volatile real-time electricity market prices [3]. However, even with this increase in price spread, it remains nontrivial to design arbitrage policies that make significant (or even positive) profit [4]. The difficulties come from the fact that future prices are unknown, difficult to forecast and may even be non-stationary [5, 6]. In this paper, we aim to develop an arbitrage policy for energy storage in a data-driven framework by using reinforcement learning [7].

For example, arbitrage using energy storage has been studied in [2, 8–11] (and see the references within). The authors in [8] studied using sodium-sulfur batteries and flywheels for arbitrage in NYISO found the batteries can be potentially profitable using data from 2001 to 2004. The authors in [2] analyzed a generic storage system in the PJM real-time market and discovered that the arbitrage value was nearly doubled from 2002 to 2007 due to higher price variations. The authors in [9] formulated a linear optimization problem to compare the arbitrage profits of 14 energy storage technologies in several major U.S. real-time electric markets. Similar studies have also been carried out in different markets, e.g., Australian national electricity market [10] and European electricity markets [11].

Crucially, all of these studies assumed perfect knowledge of electricity prices and therefore cannot be implemented as real-time arbitrage strategies.

Some recent works [12–15] have started to explicitly take the electricity price uncertainty into account when designing arbitrage strategies. The authors in [12] proposed a stochastic dynamic program to optimally operate an energy storage system using available forecast. The authors in [13] formulated a stochastic optimization problem for a storage owner to maximize the arbitrage profit under uncertainty of market prices. Both studies need to forecast electricity prices and their performances heavily rely on the quality of the forecast. However, real-time market prices are highly stochastic and notoriously difficult to forecast well [16]. To overcome the reliance on price predictions, the authors in [14] employed approximate dynamic programming to derive biding strategy for energy storage in NYISO real-time market without requiring prior knowledge of the price distribution. However, this strategy is often highly computationally expensive. The authors in [15] proposed an online modified greedy algorithm for arbitrage which is computationally straightforward to implement and does not require the full knowledge of price distributions. But it needs to estimate the bounds of prices and assume that storages are "big enough", which is not always true in practice.

The aforementioned challenges motivate us to develop an easily implementable arbitrage policy using *reinforce-ment learning (RL)*. This policy is both price-distribution-free and outperforms existing ones. Without explicitly assuming a distribution, our policy is able to operate under constantly changing prices that may be non-stationary. Over time, by repeatedly performing charge and discharge actions under different real-time prices, the proposed RL-based policy learns the best strategy that maximizes the cumulative reward. The key technical challenge turns out to be the design of a good reward function that will guide the storage to make the correct decisions. Specifically, we make the following two contributions in this paper:

- We formulate the energy storage operation as a Markov decision process (MDP) and derive a Q-learning policy [17] to optimally control the charge/discharge of the energy storage for temporal arbitrage in the real-time market.
- 2) We design a reward function that does not only reflect the instant profit of charge/discharge decisions but also incorporate historical information. Simulation results demonstrate that the designed reward function leads to significant performance improvements compared to the natural instant reward function. In addition, using real historical data, we show the proposed algorithm also leads to much higher profits than existing algorithms.

The remainder of the paper is ordered as follows. In Section II, we present the optimization problem for energy storage arbitrage. In Section III, we provide a reinforcement learning approach to obtain the arbitrage policy. Numerical simulations using real data are discussed in Section IV. Section V concludes this paper.

II. ARBITRAGE MODEL AND OPTIMIZATION PROBLEM

We consider an energy storage (e.g., a battery) operating in a real-time electricity market over a finite operational horizon $\mathcal{T} = \{1, ..., T\}$. The objective of the energy storage is to maximize its arbitrage profit by charging at low prices and discharging when prices are high. We assume the energy storage is a price taker, and its operation will not affect the market prices. We denote d_t as the discharged power from the storage at time t and c_t as the charged

power into the storage at time t. Let the real-time prices be denoted as p_t . We formulate the Arbitrage Maximization Problem (AMP) as follows:

$$\max \sum_{t=1}^{T} p_t \left(\eta_d d_t - \frac{1}{\eta_c} c_t \right)$$
 (AMP)

subject to
$$E_t = E_{t-1} + c_t - d_t, \ \forall t \in \mathcal{T}$$
 (1)

$$E_{\min} \le E_t \le E_{\max}, \ \forall t \in \mathcal{T}$$
 (2)

$$0 \le c_t \le C_{\text{max}}, \ t \in \mathcal{T} \tag{3}$$

$$0 \le d_t \le D_{\max}, \ t \in \mathcal{T} \tag{4}$$

variables: c_t , d_t , $\forall t \in \mathcal{T}$,

where $\eta_c \in (0,1)$ and $\eta_d \in (0,1)$ denote the charge/discharge efficiencies. The constraint in (1) specifies the dynamics of energy level E_t over time, (2) constraints the amount of energy in the storage to be between E_{\min} and E_{\max} , (3) and (4) bounds the maximum charge and discharge rates (denoted by C_{\max} and D_{\max} , respectively) of the storage.

The optimization problem in AMP is a linear program, and we characterize its optimal solution in the next lemma.

Lemma 1. The optimal charge and discharge profiles $\{c_t^{\star}, d_t^{\star}, \forall t \in \mathcal{T}\}$ satisfy

1) At least one of c_t^* or d_t^* is 0 at any time t;

2)
$$c_t^{\star} = \{0, \min\{C_{\max}, E_{\max} - E_{t-1}\}\},\$$

 $d_t^{\star} = \{0, \min\{D_{\max}, E_{t-1} - E_{\min}\}\}.$

Lemma 1 states that the energy storage will not charge and discharge at the same time. Also, the optimal charge and discharge power will hit the boundary per the operational constraints (1)-(4). Specifically, when the storage decides to charge, it will charge either at the maximum charge rate $C_{\rm max}$ or reaching the maximum energy level $E_{\rm max}$. Similarly, the discharge power will be either the maximum discharge rate $D_{\rm max}$ or the amount to reach the minimum energy level $E_{\rm min}$. This binary charging/discharging structure will be important when we design the reinforcement learning algorithm in the next section.

If the future prices are known, the optimization problem in AMP can be easily solved to provide an offline optimal strategy for the charge/discharge decisions. However, the offline solution is only practical if a good price forecast is available. In reality, future prices are not known in advance and the energy storage needs to make decisions based on only the current and historical data. In other words, the charge/discharge decisions $\{\hat{c}_t, \ \hat{d}_t\}$ are functions of price information up to the current time slot t, denoted by $\{p_1, ..., p_t\}$:

$$\{\hat{c}_t, \ \hat{d}_t\} = \pi(p_1, ..., p_t),$$
 (5)

where $\pi(\cdot)$ is the arbitrage policy for maximizing the profit. Therefore, AMP is a constrained sequential decision problem and can be solved by dynamic programming [18]. But the potentially high dimensionality the state space

makes dynamic programming computationally expensive, and potentially unsuitable for applications like real-time price arbitrage. Moreover, price forecast in real-time markets is extremely challenging, as the mismatch between power supply and demand can be attributed to many different causes.

III. REINFORCEMENT LEARNING ALGORITHM

To solve the *online version* of AMP, we use reinforcement learning (RL). Reinforcement learning is a general framework to solve problems in which [7]: (i) actions taken depend on the system states; (ii) a cumulative reward is optimized; (iii) only the current state and past actions are known; (iv) the system might be non-stationary. The energy storage arbitrage problem has all of the four properties: (i) different electricity prices lead to different actions (e.g., charge/discharge), and the future energy storage level depends on past actions; (ii) the energy storage aims at maximizing the total arbitrage profit; (iii) the energy storage does not have a priori knowledge of the prices, while it knows the past history; (iv) the actual price profiles are non-stationary. In the following, we describe the RL setup for AMP in more detail.

A. State Space

We define the state space of the energy arbitrage problem as a finite number of states. To be specific, the system's state can be fully described by the current price p_t and previous energy level E_{t-1} . We discretize the price into M intervals, and energy level into N intervals.

$$S = \{1, ..., M\} \times \{1, ..., N\},\$$

where $\{1,...,M\}$ represents M even price intervals from the lowest to the highest, and $\{1,...,N\}$ denotes N energy level intervals ranging from E_{\min} to E_{\max} .

B. Action Space

Per Lemma 1, the energy storage will not charge and discharge at the same time. Moreover, the optimal charge and discharge power always reach their maximum allowable rates. We denote the maximum allowable charge/discharge rates as $\tilde{D}_{\max} = \min\{D_{\max}, E_{t-1} - E_{\min}\}$ and $\tilde{C}_{\max} = \min\{C_{\max}, E_{\max} - E_{t-1}\}$. Therefore, the action space of the energy storage consists of three actions: charge at full rate, hold on, and discharge at full rate:

$$\mathcal{A} = \{ -\tilde{D}_{\text{max}}, 0, \tilde{C}_{\text{max}} \},$$

where action $a = -\tilde{D}_{\max}$ denotes discharge either at maximum rate D_{\max} or until the storage hits the minimum level E_{\min} . Action $a = \tilde{C}_{\max}$ denotes charge at maximum rate C_{\max} or until the storage reaches the maximum level E_{\max} .

C. Reward

At time t, after taking an action $a \in \mathcal{A}$ at state $s \in \mathcal{S}$, the energy storage will receive a *reward*, such that the energy storage knows how good its action is. According to the objective function of AMP, the energy storage aims

to maximize the arbitrage profit by charging at low prices and discharge at high prices. Therefore, we can define the reward as

$$r_t^1 = \begin{cases} -p_t \tilde{C}_{\max} & \text{if charge} \\ 0 & \text{if hold on} \\ p_t \tilde{D}_{\max} & \text{if discharge} \end{cases} \tag{Reward 1}$$

which is the instant reward of charge or discharge. If the energy storage charges at the rate of \tilde{C}_{\max} at time t, it will pay at the spot price and reward is negative, i.e., $-p_t\tilde{C}_{\max}$. In contrast, the energy storage discharges at the rate of \tilde{D}_{\max} and will earn a revenue of $p_t\tilde{D}_{\max}$.

Reward 1 is a straightforward and natural design, but is actually not very effective. The reason is that the negative reward for charge makes the energy storage perform conservatively in the learning process and thus the arbitrage opportunity is under explored. This motivates us to develop a more effective reward. To avoid conservative actions, we introduce an average price in the reward. The idea comes from the basic principle of arbitrage: to charge at low prices and discharge at high prices. The average price works as a simple indicator to determine whether the current price is low or high *compared to the historical values*. Specifically, the new reward is defined as

$$r_t^2 = \begin{cases} (\overline{p}_t - p_t) \tilde{C}_{\text{max}} & \text{if charge} \\ 0 & \text{if hold on} \\ (p_t - \overline{p}_t) \tilde{D}_{\text{max}} & \text{if discharge} \end{cases}$$
 (Reward 2)

where the average price \overline{p}_t is calculated by

$$\overline{p}_t = (1 - \eta)\overline{p}_{t-1} + \eta p_t, \tag{6}$$

in which η is the smoothing parameter. Note that \overline{p}_t is not a simple average that weighs all past prices equally. Instead, we use moving average in (6), such that we not only leverage the past price information but also adapt to the current price change.

We see from Reward 2 that when the energy storage charges at a price lower than the average price (i.e., $p_t < \overline{p}_t$), it will get a positive reward $(\overline{p}_t - p_t)\tilde{C}_{\max} > 0$, otherwise it will receive a loss if the spot price is greater. Similarly, Reward 2 encourages the energy storage to discharge at high price by giving a positive reward, i.e., $(p_t - \overline{p}_t)\tilde{D}_{\max} > 0$. Reward 2 outperforms Reward 1 in exploring more arbitrage opportunities and achieving higher profits. It also mitigates the non-stationarity of prices, since it weights the current price much heavier than prices in the more distant past. We will show the numerical comparisons in Section IV.

D. Q-Learning Algorithm

With the state, action and reward defined, we obtain the real-time charge and discharge policy using Q-learning (a popular subclass of RL algorithms [7]). Here the energy storage maintains a state-action value matrix Q, where each entry Q(s,a) is defined for each pair of state s and action a. When the energy storage takes a charge/discharge action under a spot price, the value matrix is updated as follows:

$$Q(s,a)_{t} = (1-\alpha)Q(s,a)_{t-1} + \alpha[r_{t} + \gamma \max_{a'} Q(s',a')], \tag{7}$$

where the parameter $\alpha \in (0,1]$ is the learning rate weighting the past value and new reward. $\gamma \in [0,1]$ is the discount rate determining the importance of future rewards. After taking an action a, the state transits from s to s', and the energy storage updates the value matrix incorporating the instant reward r_t (e.g., Reward 1 or 2) and the future value $\max_{a'} Q(s', a')$ in state s'. Over time, the energy storage can learn the value each action in all states. When Q(s, a) converges to the optimal state-action values, we obtain the optimal arbitrage policy. Specifically, the Q-learning algorithm can derive an arbitrage policy for (5) as

$$a^* = \pi(s) = \arg\max_{a} Q(s, a), \tag{8}$$

which is the optimal arbitrage policy guranteed for finite MDP [17]. For any state s, the energy storage always chooses the best action a^* which maximizes the value matrix Q(s, a).

Algorithm 1 Q-learning for energy storage arbitrage

- 1: **Initialization**: In each time slot $t \in \{1, ..., T\}$, set the iteration count k = 1, $\alpha = 0.5$, $\alpha = 0.9$, and $\epsilon = 0.9$. Initialize the Q-matrix, i.e., Q = 0.
- 2: repeat
- 3: **Step1:** Observe state s based on price and energy level;
- 4: **Step2:** Decide the best action a (using ϵ -greedy method) based on Q(s, a);
- 5: **Step3:** Calculate the reward (using Reward 1 or 2);
- 6: **Step4:** Update Q(s, a) according to (7) and energy level in (1);
- 7: $s \leftarrow s' \text{ and } k \leftarrow k+1;$
- 8: **until** end of operation, i.e., t = T.
- 9: **end**

The step-by-step Q-learning algorithm for energy arbitrage is presented in Algorithm 1. To avoid the learning algorithm getting stuck at sub-optimal solutions, we employ ϵ -greedy [17]. The algorithm not only exploits the best action following (8) but also explores other actions, which could be potentially better. Specifically, using ϵ -greedy, the algorithm will randomly choose actions with probability $\epsilon \in [0,1]$, and choose the best action in (8) with probability $1 - \epsilon$.

IV. NUMERICAL RESULTS

In this section, we evaluate two reward functions and also compare our algorithm to a baseline in [15] under both synthetic prices and realistic prices. For synthetic prices, we generate i.i.d. (independent and identically distributed) prices, and for the realistic price, we use hourly prices from ISO New England real-time market [19] from January 1, 2016 to December 31, 2017. The realistic price is depicted in Figure 1. We see that the averaged price is flat but the instantaneous prices fluctuate significantly with periodic spikes.

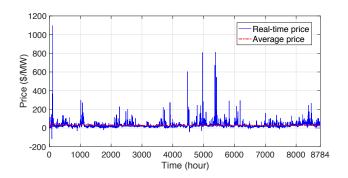


Fig. 1: PJM Real-time price.

A. Synthetic Price

We first evaluate the two reward functions under synthetic price, which is uniformly distributed in [0,1] over 1500 hours. We set $C_{\text{max}} = D_{\text{max}} = 1$, $E_{\text{min}} = 0$ and $E_{\text{max}} = 1$. The cumulative profits for both rewards are depicted in Figure 2. Both profits stay flat over the first 300 hours, as the algorithm is exploring the environment with different prices. Afterwards, the algorithm using Reward 2 starts to make profit and achieves 166% more than Reward 1 in the end.

To further understand the how Reward 1 and Reward 2 affect the storage operation, we plot the evolution of energy level over a 48 hour horizon in Fig. 3. We see that algorithm using Reward 1 performs conservatively while Reward 2 makes the algorithm actively charge and discharge to take advantage of price spread. Therefore, Reward 2 leads to a more profitable arbitrage strategy.

B. Real Historical Price

We evaluate the two reward functions using realistic prices from ISO New England real-time market in 2016. We plot the cumulative profits of two rewards during training in Figure 4. We see that Reward 1 fails to make profit while using Reward 2 produces a high profit. This demonstrates the effectiveness of our designed reward: it is able to adapt to price changes and makes profit continuously.

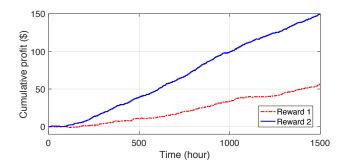
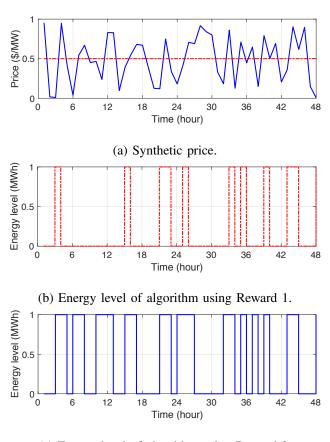


Fig. 2: Cumulative profits under synthetic prices.



(c) Energy level of algorithm using Reward 2.

Fig. 3: Price and energy levels over a 48 hour period using reward 1 and reward 2 under synthetic prices.

We also plot the evolution of energy levels over a 48-hours operational horizon in Figure 5. We see that algorithm using Reward 1 cannot capture the price differences but makes charge/discharge when the real-time price is flat. In contrast, our algorithm using Reward 2 is able to charge at low prices at hours 2 and 29, hold the energy when prices are low, and discharge at hours 12 and 44, respectively, when the price reaches a relatively high point.

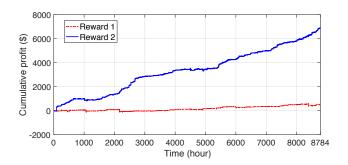


Fig. 4: Cumulative profits under real-time prices.

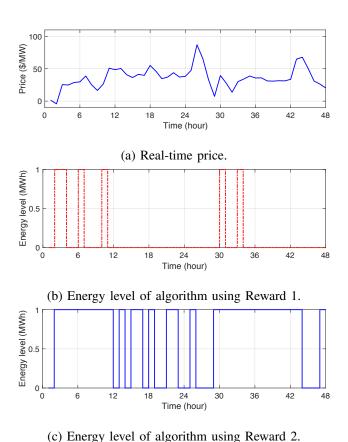


Fig. 5: Price and energy levels over a 48 hour horizon for reward 1 and reward 2 under historical data.

C. Comparison with baseline algorithm

Above discussion demonstrates that Reward 2 performs much better than Reward 1, and thus we stick to Reward 2 and compare our algorithm with a baseline algorithm called online modified greedy algorithm in [15]. This algorithm uses a thresholding strategy to control charge and discharge in an online fashion. We configure the parameters for the baseline according to [15]. The arbitrage profits of two algorithms are simulated on an 8–MWh battery, with a charge/discharge rate of 1MW as depicted in Figure 6. The baseline algorithm can only get \$5,845, while our algorithm earns \$28,027 that is 4.8 times of the baseline profit. The profit of the baseline decreases when the charge/discharge rate increases to 2MW. But our algorithm achieves even a higher profit, i.e., \$39,690, which is 8.6 times of the baseline profit \$4,603. The reason is that the baseline algorithm relies on the off-line estimate of the price information and lacks adaptability to the real-time prices. Our algorithm updates the average price to adapt to the price changes and thus performs better.

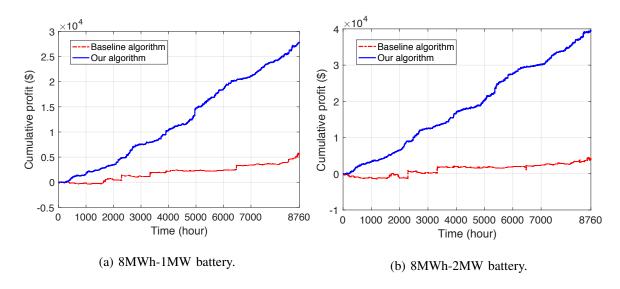


Fig. 6: Cumulative profits of the baseline algorithm in [15] and our algorithm.

V. CONCLUSION

In this paper, we derive an arbitrage policy for energy storage operation in real-time markets via reinforcement learning. Specifically, we model the energy storage arbitrage problem as an MDP and derive a Q-learning policy to control the charge/discharge of the energy storage. We design a reward function that does not only reflect the instant profit of charge/discharge decisions but also incorporate the history information. Simulation results demonstrate our designed reward function leads to significant performance improvement and our algorithm achieves much more profit compared with existing baseline method. We will consider self-discharge and degradation of battery in our future work.

ACKNOWLEDGMENT

This work was partially supported by the University of Washington Clean Energy Institute.

REFERENCES

- [1] J. Eyer and G. Corey, "Energy storage for the electricity grid: Benefits and market potential assessment guide," *Sandia National Laboratories*, vol. 20, no. 10, p. 5, 2010.
- [2] R. Sioshansi, P. Denholm, T. Jenkin, and J. Weiss, "Estimating the value of electricity storage in pjm: Arbitrage and some welfare effects," *Energy economics*, vol. 31, no. 2, pp. 269–277, 2009.
- [3] C.-K. Woo, I. Horowitz, J. Moore, and A. Pacheco, "The impact of wind generation on the electricity spot-market price level and variance: The texas experience," *Energy Policy*, vol. 39, no. 7, pp. 3939–3944, 2011.
- [4] R. H. Byrne and C. A. Silva-Monroy, "Estimating the maximum potential revenue for grid connected electricity storage: Arbitrage and regulation," *Sandia National Laboratories*, 2012.
- [5] T. T. Kim and H. V. Poor, "Scheduling power consumption with price uncertainty," *IEEE Transactions on Smart Grid*, vol. 2, no. 3, pp. 519–527, 2011.
- [6] S. Borenstein, "The long-run efficiency of real-time electricity pricing," The Energy Journal, pp. 93-116, 2005.
- [7] R. S. Sutton and A. G. Barto, "Reinforcement learning: An introduction," 2011.

- [8] R. Walawalkar, J. Apt, and R. Mancini, "Economics of electric energy storage for energy arbitrage and regulation in new york," *Energy Policy*, vol. 35, no. 4, pp. 2558–2568, 2007.
- [9] K. Bradbury, L. Pratson, and D. Patiño-Echeverri, "Economic viability of energy storage systems based on price arbitrage potential in real-time us electricity markets," *Applied Energy*, vol. 114, pp. 512–519, 2014.
- [10] D. McConnell, T. Forcey, and M. Sandiford, "Estimating the value of electricity storage in an energy-only wholesale market," *Applied Energy*, vol. 159, pp. 422–432, 2015.
- [11] D. Zafirakis, K. J. Chalvatzis, G. Baiocchi, and G. Daskalakis, "The value of arbitrage for energy storage: Evidence from european electricity markets," *Applied Energy*, vol. 184, pp. 971–986, 2016.
- [12] K. Abdulla, J. De Hoog, V. Muenzel, F. Suits, K. Steer, A. Wirth, and S. Halgamuge, "Optimal operation of energy storage systems considering forecasts and battery degradation," *IEEE Transactions on Smart Grid*, 2016.
- [13] D. Krishnamurthy, C. Uckun, Z. Zhou, P. Thimmapuram, and A. Botterud, "Energy storage arbitrage under day-ahead and real-time price uncertainty," *IEEE Transactions on Power Systems*, 2017.
- [14] D. R. Jiang and W. B. Powell, "Optimal hour-ahead bidding in the real-time electricity market with battery storage using approximate dynamic programming," *INFORMS Journal on Computing*, vol. 27, no. 3, pp. 525–543, 2015.
- [15] J. Qin, Y. Chow, J. Yang, and R. Rajagopal, "Online modified greedy algorithm for storage control under uncertainty," *IEEE Transactions on Power Systems*, vol. 31, no. 3, pp. 1729–1743, 2016.
- [16] R. Weron, "Electricity price forecasting: A review of the state-of-the-art with a look into the future," *International Journal of Forecasting*, vol. 30, no. 4, pp. 1030 1081, 2014.
- [17] C. J. Watkins and P. Dayan, "Q-learning," Machine learning, vol. 8, no. 3-4, pp. 279-292, 1992.
- [18] R. A. Howard, Dynamic programming and markov processes. Oxford, England: John Wiley, 1960.
- [19] "Hourly real-time lmp." [Online]. Available: https://www.iso-ne.com