

Energy Management of Smart Homes with Electric Vehicles Using Deep Reinforcement Learning

Xavier Weiss, Qianwen Xu, Lars Nordström

KTH

Teknikringen 33, 10044

Stockholm, Sweden

Telephone: +46(8)790-6830

Email: xavierw@kth.se

URL: <https://www.kth.se/>

Keywords

«Energy Management System (EMS)», «Microgrid», «Electric Vehicle», «Energy storage», «Deep learning», «Safety»

Abstract

The proliferation of electric vehicles (EVs) has resulted in new charging infrastructure at all levels, including domestically. These new domestic EVs can potentially provide vehicle to home (V2H) services where EVs are used as energy storage systems (ESSs) for the home when they are not in use. Energy management systems (EMSs) can control these EVs to minimize the electricity cost to the owner but must satisfy constraints. Uncertainty in EV availability and the microgrid environment is also a challenge and can be addressed through real-time operation. Hence this paper formulates the EV charge/discharge scheduling problem as a Markov Decision Process (MDP). A safe implementation of Proximal Policy Optimization (PPO) is proposed for real-time optimization and compared to a day-ahead Mixed Integer Linear Programming (MILP) benchmark. The resulting PPO agent is able to minimize RA and SD costs for a typical EV user 3% better than the MILP solution. It obtains a 39% higher electricity cost than MILP, but unlike MILP does not require accurate forecasting data and operates in real-time.

Introduction

The global supply of EVs is growing and is a crucial factor in decarbonizing the transport sector. The large-scale integration of EVs into the grid will greatly increase electricity consumption, and thus pose significant stress on the grid. To alleviate the stress, a smart home integrated with EV and PhotoVoltaics (PV) is a promising solution.

A smart home is composed of local loads, an ESS and generators and thus forms a natural building block for the smart grid [2]. EVs in a smart home can act as both a load and an ESS. Conventionally, an EV acts as a load with a unidirectional flow of power from home to vehicle (H2V). Increased use of bidirectional converters [3, 4] will allow the EV to provide vehicle to home (V2H) and V2G functionality. This has the potential to save costs for the home owner and provide ancillary services to the grid.

An EMS is key for the optimal operation of smart homes. Conventionally, model-based methods are used. Ref. [7] proposes an optimal control strategy for efficient utilization of available EV for storing PV power and providing grid support. In [8] a smart home with EV and PV is optimized using a dynamic programming successive algorithm to minimize the variance in household load. However, these examples do not consider uncertainty.

EMSs considering uncertainties have thus been proposed. Stochastic optimization is used to forecast uncertainty in [9] for power grids with high penetration of renewables and EVs. Similarly, [6] is able to improve the profit of a charging station and reduce departure delay by using a stochastic Lyapunov drift technique. Robust optimization is proposed in [10] for EV charging in an unregulated electricity market

and obtains real-time performance. However, these methods require the distribution of uncertainties to be known - which is difficult to obtain.

Easy access to data and advances in DRL have led to a strong interest in data-driven methods. While these cannot guarantee a global optimum, they can be more computationally efficient and versatile than model-based methods, which is beneficial for real-time applications. In [11], a DQN-based EMS is developed to optimize the cost of a residential microgrid. Ref. [5] develops a two-level actor-critic method to size and schedule components in a smart home. A PPO-based EMS is proposed for demand response in [12]. Finally, in [13] long-short term memory networks are used to build a data-driven forecasting model, which is then combined with a deep deterministic policy gradient agent to reduce charging costs for an EV owner with random arrival times by up to 70.2% - relative to an unmanaged scenario. However, these implementations are not safe. For instance, even when a DRL agent is trained and exhibiting good performance it can still violate constraints - especially in unseen scenarios. Moreover, they do not consider the EV owner's requirement to minimize range anxiety.

To address the above issues, this paper proposes a DRL-based EMS for smart homes which guarantees safety and fulfills customer requirements.

Problem Formulation

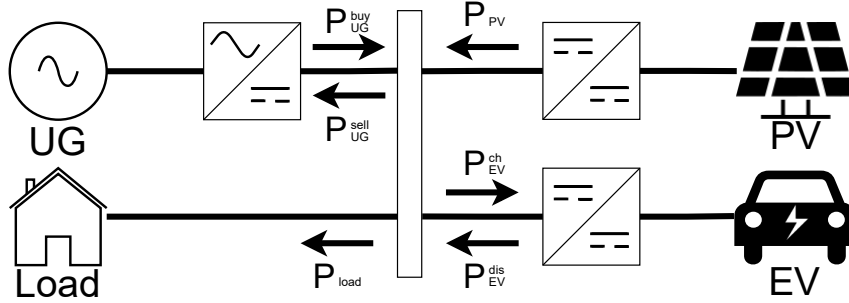


Fig. 1: Layout of residential microgrid

The goal of the EMS considered in this study is to minimize the total operating cost of the smart home shown in Fig.1. The smart home consists of PV, EV, residential load and a Utility Grid (UG) connection. The minimization of the total operating cost of this smart home can therefore be formulated as an optimization problem of the form:

$$\min \sum_{\tau=1}^T C_{UG}^{buy}(\tau) \cdot P_{UG}^{buy}(\tau) - C_{UG}^{sell}(\tau) \cdot P_{UG}^{sell}(\tau) + C_{ev}(P_{ev}^{ch}(\tau) + P_{ev}^{dis}(\tau)) \quad (1)$$

subjects to constraints (3)-(13)

where τ is a discrete time slot with a fixed interval (e.g. 15 minutes), C_{UG}^{buy} and C_{UG}^{sell} are the cost of buying and selling power from the Utility Grid (UG), P_{UG}^{buy} and P_{UG}^{sell} are the amounts of power bought and sold to the UG, P_{ev}^{ch} and P_{ev}^{dis} are the amounts of power used to charge/discharge the EV, and C_{ev} is the degradation cost per kWh of using the EV's battery.

The degradation cost of each charge/discharge cycle can be estimated linearly:

$$C_{ev}(P_{ev}(\tau)) = \rho |SoE(\tau) - SoE(\tau - \Delta\tau)| \quad (2)$$

where $SoE = SoC \cdot E_{max}$ is the State of Energy (SoE) of the EV, E_{max} is the maximum energy capacity of the battery and P_{ev} is the power transferred from/to the battery. The State of Charge (SoC) is a fraction between 0 and 1 of the EV's maximum energy storage capacity E_{max} .

Power Balance Constraints:

$$P_{load}(\tau) + \frac{1}{\eta_{BIC}} P_{UG}^{sell}(\tau) + P_{ev}^{ch}(\tau) = P_{PV}(\tau) + \eta_{BIC} P_{UG}^{buy}(\tau) + P_{ev}^{dis}(\tau) \quad \forall \tau \quad (3)$$

where η_{BIC} is the efficiency of the bi-directional converter between the smart home and the utility grid.

Energy Storage Constraints:

$$SoE(\tau) = SoE(\tau - 1) + P_{ev}^{ch}(\tau) \eta_{ch} - \frac{P_{ev}^{dis}(\tau)}{\eta_{dis}} \quad (4)$$

$$|SoE(\tau) - SoE(\tau - 1)| \leq \Delta SoE_{max} \quad (5)$$

$$E_{min}(\tau) \leq SoE(\tau) \leq E_{max}(\tau) \quad \forall \tau \quad (6)$$

where η_{ch} and η_{dis} are the charging and discharging efficiency of the EV's battery, respectively.

Power capacity constraints:

$$0 \leq P_{ev}^{ch}(\tau) \leq P_{ev}^{ch,max}(\tau) \quad (7)$$

$$0 \leq P_{ev}^{dis}(\tau) \leq P_{ev}^{dis,max}(\tau) \quad (8)$$

$$0 \leq P_{buy}^{UG}(\tau) \leq P_{buy}^{UG,max}(\tau) \quad (9)$$

$$0 \leq P_{sell}^{UG}(\tau) \leq P_{sell}^{UG,max}(\tau) \quad (10)$$

$$0 \leq P_{PV}(\tau) \leq P_{PV}^{max}(\tau) \quad (11)$$

Bi-directional Flow Constraints:

$$P_{ev}^{ch}(\tau) P_{ev}^{dis}(\tau) = 0 \quad \forall \tau \quad (12)$$

$$P_{buy}^{UG}(\tau) P_{sell}^{UG}(\tau) = 0 \quad \forall \tau \quad (13)$$

The optimization problem in Eq.1 with constraints Eqs.3-13 is a MILP problem and can be solved by commercial software like Gurobi. However, uncertainties in PV, EV and the load are not considered which leads to ineffective results. To handle uncertainties a DRL-based EMS is proposed in the next section.

Proposed Method

A safe PPO agent is proposed to manage the EV in the EMS by modelling it as a Markov Decision Process (MDP). First the EMS problem is reformulated as an MDP so it can be solved using DRL methods. A real-time PPO algorithm [14] is then suggested for its simplicity and its robustness in noisy environments. The PPO algorithm is finally equipped with a novel safety layer to guarantee it does not violate any of the constraints in Eq.3-13.

Markov Decision Process

State Space The state of the residential microgrid is given as:

$$s(\tau) = (C_{UG}, P_{load}, P_{PV}, SoE, A_{EV}) \quad (14)$$

where P_{load} and P_{PV} are the energy (in kWh) consumed/generated by the house and the PV installation and A_{EV} is a binary value indicating whether the EV is available.

Action Space The action space is continuous and is generally expressed as a 4-vector:

$$a(\tau) = (P_{buy}^{UG}, P_{sell}^{UG}, P_{ev}^{ch}, P_{ev}^{dis}) \quad (15)$$

In this implementation, $a(\tau)$ is the output of a safety layer. Through the use of a safety layer the agent

only directly controls the charge P_{ev}^{ch} and discharge P_{ev}^{dis} of the EV's battery:

$$a'(\tau) = (P_{ev}^{ch}, P_{ev}^{dis}) \quad (16)$$

Reward The reward is simply a reformulation of Eq.1:

$$r(\tau) = C_{UG}^{sell} - C_{UG}^{buy} - C_{ev} - C_{RA} - C_{SD} \quad (17)$$

The agent's objective is to maximize the reward in real-time, hence there is no discount factor to consider. All costs are specified in Table.Ib.

The Range Anxiety (RA) cost reflects the EV owner's anxiety when leaving without enough energy to complete their trip:

$$C_{RA} = \max(K_{RA} \cdot (E_{trip} - SoE_{\tau}), 0) \quad (18)$$

where E_{trip} is the energy required to complete the trip, SoE_{τ} is the SoE at time slot τ and K_{RA} is a constant.

The State Difference (SD) cost rewards the agent for ending each day with the same SoE it started with, or higher. While this is not necessary for a real-time implementation, it encourages the agent to not cause energy shortages when operating across multiple days.

$$C_{SD} = K_{SD}^{\max(T-\tau, 1)} \cdot \max(SoE_{initial} - SoE_{\tau}, 0) \quad (19)$$

where K_{SD} is a constant and $SoE_{initial}$ is the SoE at the start of the day. For $\kappa < 1$, c_{SD} increases exponentially as the agent approaches the end of the day unless the condition is met.

DRL Agent

For its performance and simplicity, the standard PPO algorithm is selected. Since the microgrid environment has several sources of noise – including the variation in PV generation due to weather, load and EV usage due to occupant behavior and electricity prices due to market forces – the PPO agent is especially appropriate since it is generally more robust against perturbations to the input s_{τ} than other DRL algorithms like Deep Deterministic Policy Gradient (DDPG). The full implementation details for PPO are provided in [14, 17]. Parameters for both the actor and the critic components of the PPO agent are given in Table.Ia.

To avoid destroying the learned policy in a single update and maintain computational efficiency, PPO adopts a clipped version of trust regions [15] directly into the surrogate loss function:

$$L^{CLIP}(\theta) = \hat{E}_{\tau}[\min(\delta_{\tau}(\theta)\hat{A}_{\tau}, \text{clip}(\delta_{\tau}(\theta), 1 - \epsilon, 1 + \epsilon)\hat{A}_{\tau})] \quad (20)$$

where ϵ puts a trust-region constraint on how much the policy can be updated and $\delta_{\tau}(\theta) = \frac{\pi_{\theta}(a(\tau)|s(\tau))}{\pi_{\theta_{old}}(a(\tau)|s(\tau))}$ is the ratio of the new policy over the old policy. The advantage \hat{A}_{τ} is given by:

$$\hat{A}_{\tau} = \delta_{\tau} + (\gamma)\delta_{\tau+1} + \dots + (\gamma)^{T-\tau+1}\delta_{T-1} \quad (21)$$

where $\delta_{\tau} = r_{\tau} + \gamma V(s_{\tau+1}) - V(s_{\tau})$, γ is the discount factor and V is the value function. The advantage therefore compares the discounted rewards r_{τ} to the baseline value estimate from the critic.

The complete loss function for PPO also adds a Mean-Square Error (MSE) and entropy term:

$$L_{\tau}^{CLIP+MSE+S}(\theta) = \hat{\mathbb{E}}_{\tau} [L_{\tau}^{CLIP}(\theta) - c_1 L_{\tau}^{MSE}(\theta) + c_2 S[\pi_{\theta}](s_{\tau})] \quad (22)$$

where L_{τ}^{CLIP} is the clipped surrogate loss function from Eq.20 at time slot τ , $L_{\tau}^{MSE}(\theta) (= V_{\theta}(s_{\tau}) - r_{\tau})^2$ is the MSE between the value function $V_{\theta}(\theta)$ and the reward r_{τ} and S is the entropy of the agent's output. In this study the constants c_1 and c_2 are set to 0.5 and $c_2 = -0.01$.

The standard deviation of the PPO agent's output layer action $a'(t)$ is also set to gradually decay according to:

$$\sigma(\text{epoch}_{no}T + \tau) = \max(\sigma_{min}, \sigma_{initial} - \sigma_{decay\ rate}(\text{epoch}_{no}T + \tau - 1)) \quad (23)$$

Safety Layer

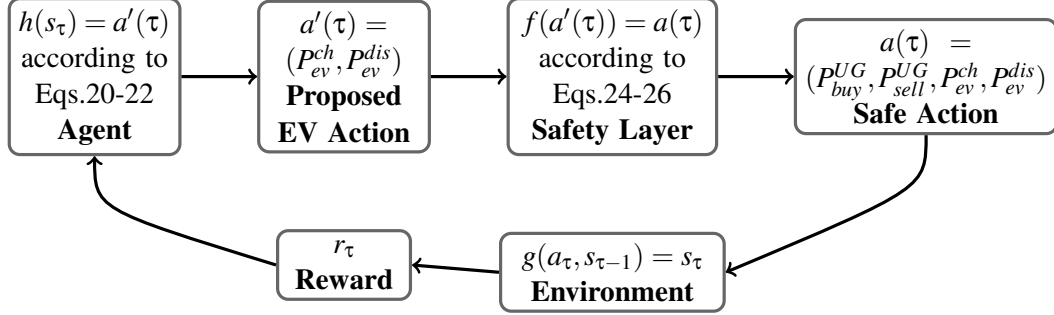


Fig. 2: Architecture of proposed safe PPO agent

In order to guarantee that the PPO agents satisfy the constraints Eqs.3-13 a safety layer is employed. The safety layer takes the proposed action from the agent $a'(t)$ (see Eq.16) and transforms it into a legal action $a(t)$ (see Eq.15) in the environment as shown in Fig.2.

Since the PPO agent's output is passed through a \tanh activation function, the output is in the $[-1, 1]$ range. This is dynamically re-scaled to match the safe limits of EV charging or discharging:

$$\begin{aligned} P_{ev}^{ch} &= lb' + \frac{(P_{ev}^{ch} - P_{ev}^{ch,lb})(ub' - lb')}{P_{ev}^{ch,ub} - P_{ev}^{ch,lb}} \\ P_{ev}^{dis} &= lb' + \frac{(P_{ev}^{dis} - P_{ev}^{dis,lb})(ub' - lb')}{P_{ev}^{dis,ub} - P_{ev}^{dis,lb}} \end{aligned} \quad (24)$$

where lb' & ub' are the lower and upper bound of the original range $[-1, 1]$ in this case). For the EV, the SoE, power capacity and availability constraints can be included by taking the minimum of the rate limit, the remaining SoE, and available SoE:

$$\begin{aligned} P_{ev}^{ch,ub} &= \min(A_{ev}P_{ev}^{ch,max}, E_{max} - SoE) \text{ and } P_{ev}^{ch,lb} = -P_{ev}^{ch,ub} \\ P_{ev}^{dis,ub} &= \min(A_{ev}P_{ev}^{dis,max}, SoE - E_{min}) \text{ and } P_{ev}^{dis,lb} = -P_{ev}^{dis,ub} \end{aligned} \quad (25)$$

where the use of a minimum function implies that the safety layer applies a non-linear transformation on the PPO's output.

As suggested by the split in the SoE constraint from Eq.4, the bidirectional-flow constraints from Eq.12-13 are guaranteed in the safety layer by splitting a single input into two. For the range $[-1, 1]$ the split is at 0.

Power balance is maintained at every time step by having the agent only output 1 value corresponding to both P_{ev}^{ch} and P_{ev}^{dis} . The other actions of buying/selling electricity are derived automatically from Eq.3:

$$B = (P_{load} + P_{ev}^{ch} - P_{ev}^{dis} - P_{PV}) \quad (26)$$

$$P_{buy}^{grid} = \frac{1}{\eta_{BIC}}(B) \text{ if } B > 0, 0 \text{ otherwise} \quad (27)$$

$$P_{sell}^{grid} = \eta_{BIC}(B) \text{ if } B < 0, 0 \text{ otherwise} \quad (28)$$

where η_{BIC} is the bidirectional converter's efficiency.

Based on the Eq.24 and Eq.26 the safety layer dynamically transforms the raw and potentially unsafe EV charge/discharge action $a'(t)$ proposed by PPO agent into safe action $a(t)$ that can be executed in the environment.

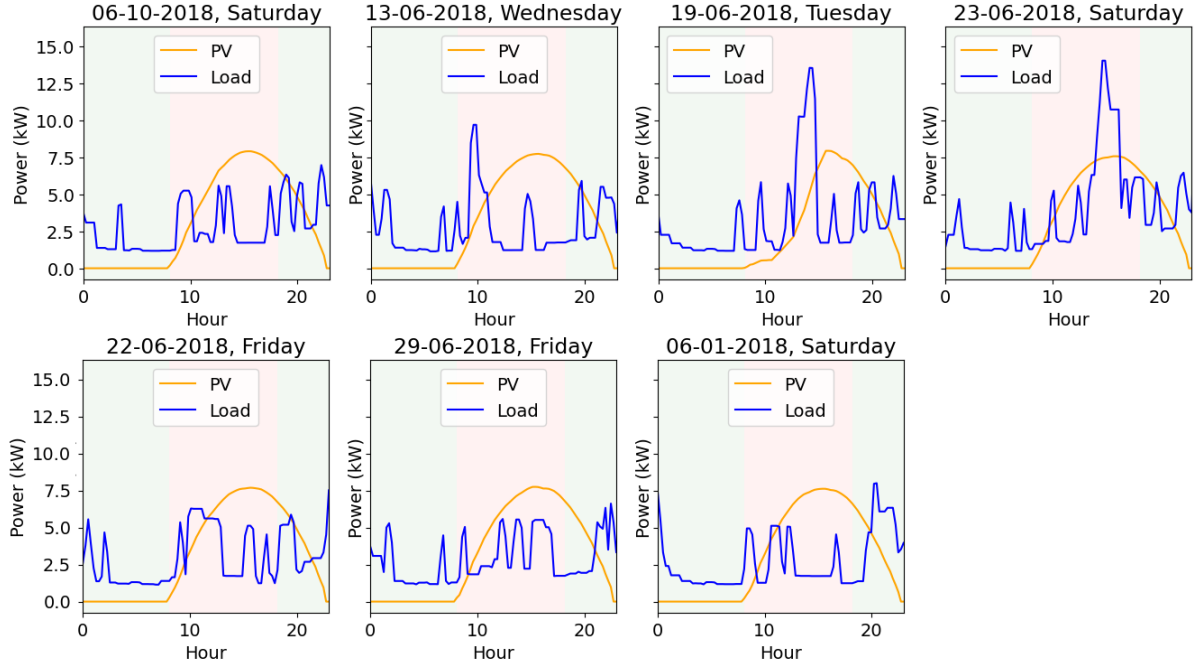


Fig. 3: Sample load, PV and EV profiles in training set

Case Study

Environment		PPO	
Time Resolution	15 minutes	Learning rate _{actor}	1e-5
Train Size	15 days	Learning rate _{critic}	0.9
Validation Size	6 Days	γ	0.99999
Test Size	8 Days	ϵ	1.0
Patience	20 epochs	$\sigma_{initial}$	1.0
Grace Period	100 epochs	σ_{min}	0.1
Max Epochs	800 epochs	$\sigma_{decay\ rate}$	0.01
Random Seed	1812	Batch size	5

(a) Hyperparameters

Costs		EV Battery		Microgrid Components	
Parameter	Value	Parameter	Value	Parameter	Value
κ_{SD}	0.5	$SoE_{initial}$	4.8 kWh	eff_{bat}	0.95
κ_{RA}	0.1 per kWh	E_{max}	24 kWh	eff_{conv}	0.95
c_{ev}	0.01 per kWh	E_{min}	0 kWh	$\Delta Conv_{max}$	200 kWh
c_{buy}	See Fig. 5	ΔSoE_{max}	1.75 kWh	$\Delta U_{G_{max}}$	1000 kWh
c_{sell}	0.5*c _{buy}	E_{rip}	10 kWh		

(b) Fixed parameters

Table I: Hyperparameters and fixed parameters for the experimental setup

In this case study we consider the DC microgrid setup shown in Fig.1 using data from the Open Energy Data Initiative (OEDI) accessible at [1] and shown in Fig.3. The PPO agent is trained on 30 days of simulated time series data. The data is at a 15 minute time resolution and is from a fictional house in California. The house has a 7.4kW PV installation and a mean load of 2.6 kW. The variation in PV and load profiles will allow an assessment of the PPO agent's ability to plan.

The EV profile was selected to reflect a weekday commuter lifestyle, where the EV departs at 8am and returns at 6pm. The capacity of the battery was set to 24 kWh based on [16], with a maximum charge/discharge rate of 7 kWh every 15 minutes. The battery is assumed to be 95% efficient, and starts with 4.8 kWh of charge, up to a maximum of 24 kWh and a minimum of 0kWh.

The converter is assumed to be 95% efficient and is able to convert up to 200 kWh per time step. Electricity can be bought from/to the utility grid according to the price profile shown in Fig.5. Electricity can also be sold using the same price profile but at 50% of the buying price.

Method	Test Score (price only)	Test Score (w. RA and SD)
PPO	-7.7053	-8.55218
MILP	-4.7332	-8.8287

Table II: Comparison of performances on unseen test data for 3 different methods

Results

The scores of the PPO and MILP algorithms on unseen data show that the DRL agent results in less cost savings for the EV owner than the conventional method. However, the MILP solution was solved by assuming all information for the 24 hour window is perfectly forecast while the PPO solution operates in real-time and hence does not use a forecast. Once the soft constraints of RA and SA are considered, the PPO algorithm even marginally outperforms the MILP algorithm - though these are not explicitly minimized in the MILP. That a similar performance is obtained using data alone and with real-time operation therefore shows the merit of the data-driven technique.

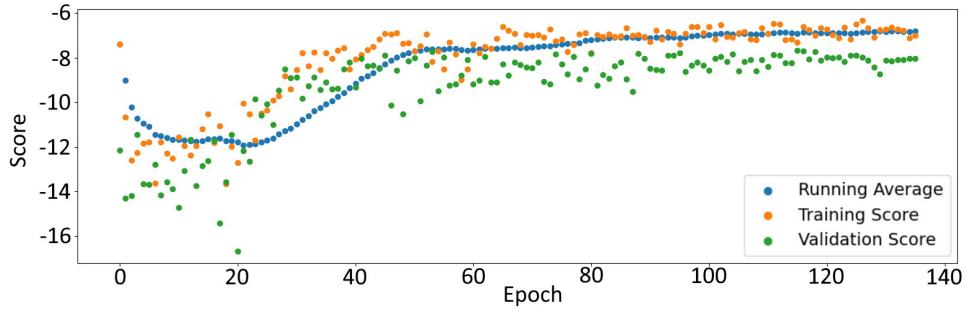


Fig. 4: PPO algorithm training curve

The training performance curve of the PPO agent in Fig.4 shows there is room for further refinement of the agent. The high sample efficiency of the PPO algorithm is clearly visible since it converges in approximately 50 epochs. There is a small gap between the validation set scores and the training set scores, hence the model does not appear to be overfitting. This poor performance may be due to the reward being calculated on the re-scaled safe action rather than the raw output of the agents, hence there is not a clear mapping between the reward and the agent's proposed action. Hence a large penalty could be given for proposed actions that would violate a constraint, but nevertheless are rendered safe through the safety layer.

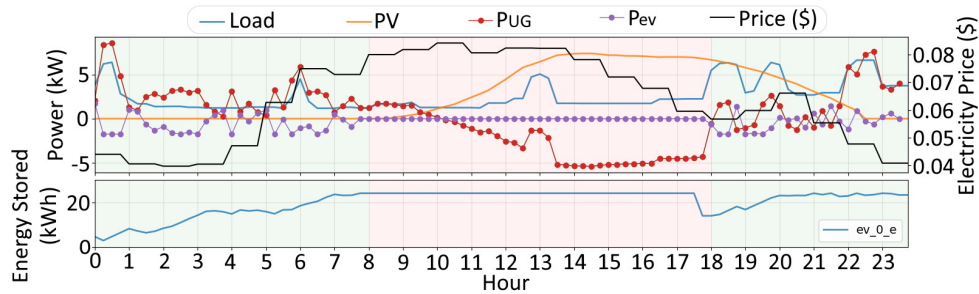


Fig. 5: PPO agent behavior

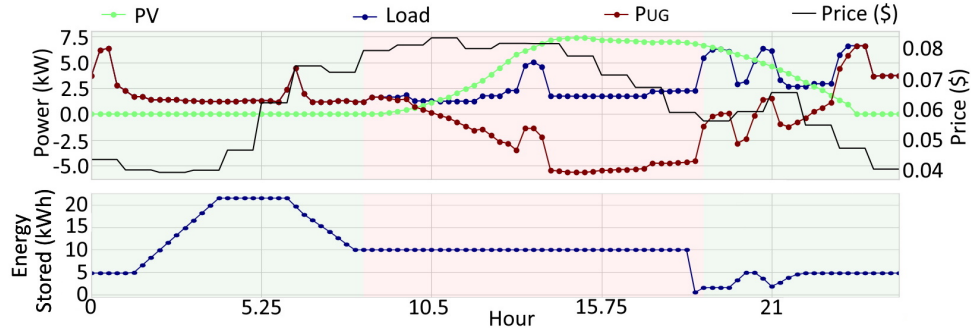


Fig. 6: Behaviour of MILP solution

Based on the optimal behaviour shown by the MILP solution, the data-driven method does make optimal use of the variable electricity price. As shown in Fig.5 the PPO agent appears to charge preferentially in the early hours (midnight to 4am) as electricity is cheaper then. However, more cost savings would be made if the agent discharged the battery when electricity prices were high. Hence the agent appears to have converged to a local minimum and will require further exploration to find the global optimum.

Conclusion

The black-box behaviour and unconstrained output of standard deep reinforcement learning agents can lead to safety violations during operation. To address this concern in the energy management of a smart home with an electric vehicle, a safety layer is appended to the output of a PPO agent to dynamically re-scale any outputs into a safe range, thus guaranteeing the safe operation of the system. The proposed method is able to operate reliably in real-time and on unseen data. It is compared to a day-ahead MILP algorithm, which has perfect information about the next 24 hours of solar production and household demand. The PPO algorithm has a 39% higher electricity cost than the MILP algorithm, but obtains a comparable performance (within 3%) in minimizing the range anxiety of the EV owner and ensuring the battery charge at the end of the day is at least as high as the start of the day. Hence despite the worse savings the shielded PPO is a more practical option than the MILP algorithm, since the latter does not run in real-time.

References

- [1] "OEDI: Commercial and Residential Hourly Load Profiles for all TMY3 Locations in the United States", <https://data.openei.org/submissions/4520>.
- [2] Q. Xu, T. Zhao, Y. Xu, Z. Xu, P. Wang and F. Blaabjerg, "A Distributed and Robust Energy Management System for Networked Hybrid AC/DC Microgrids", *IEEE Transactions on Smart Grid*, vol. 11, no. 4, pp.3496-3508, 2020.
- [3] Y. Shen, H. Wang, A. Al-Durra, Z. Qin and F. Blaabjerg, "A Bidirectional Resonant DC–DC Converter Suitable for Wide Voltage Gain Range," in *IEEE Transactions on Power Electronics*, vol. 33, no. 4, pp. 2957-2975, April 2018, doi: 10.1109/TPEL.2017.2710162.
- [4] S. Liu, X. Xie and L. Yang, "Analysis, Modeling and Implementation of a Switching Bi-Directional Buck-Boost Converter Based on Electric Vehicle Hybrid Energy Storage for V2G System," in *IEEE Access*, vol. 8, pp. 65868-65879, 2020, doi: 10.1109/ACCESS.2020.2985772.
- [5] S. Lee and D. Choi, "Energy Management of Smart Home with Home Appliances, Energy Storage System and Electric Vehicle: A Hierarchical Deep Reinforcement Learning Approach" in *Sensors 2020*, Vol. 20, Page 2157.
- [6] E. Bagherzadeh, A. Ghiasian and A. Rabiee, "Long-term profit for electric vehicle charging stations: A stochastic optimization approach", *Sustainable Energy, Grids and Network*, Vol.24, 2020
- [7] M. J. E. Alam, K. M. Muttaqi, and D. Sutanto, "Effective utilization of available PEV battery capacity for mitigation of solar PV impact and grid support with integrated V2G functionality," *IEEE Trans. Smart Grid*, vol. 7, no. 3, pp. 1562–1571, 2016.
- [8] F. Hafiz, P. Fajri and I. Husain, "Load regulation of a smart household with PV-storage and electric vehicle by dynamic programming successive algorithm technique," 2016 IEEE Power and Energy Society General Meeting (PESGM), 2016, pp. 1-5

- [9] B. Wang, P. Dehghanian and D. Zhao, "Chance-Constrained Energy Management System for Power Grids With High Proliferation of Renewables and Electric Vehicles," in *IEEE Transactions on Smart Grid*, vol. 11, no. 3, pp. 2324-2336, May 2020
- [10] N. Korolko and Z. Sahinoglu, "Robust Optimization of EV Charging Schedules in Unregulated Electricity Markets," in *IEEE Transactions on Smart Grid*, vol. 8, no. 1, pp. 149-157, Jan. 2017.
- [11] Y. Liu, D. Zhang and H. B. Gooi, "Optimization strategy based on deep reinforcement learning for home energy management," in *CSEE Journal of Power and Energy Systems*, vol. 6, no. 3, pp. 572-582, Sept. 2020
- [12] H. Li, Z. Wan and H. He, "A Deep Reinforcement Learning Based Approach for Home Energy Management System," 2020 IEEE Power Energy Society Innovative Smart Grid Technologies Conference (ISGT), 2020, pp. 1-5.
- [13] S. Li et al., "Electric Vehicle Charging Management Based on Deep Reinforcement Learning," in *Journal of Modern Power Systems and Clean Energy*.
- [14] Schulman, J., Wolski, F., Dhariwal, P., Radford, A. and Klimov, O., 2017. Proximal policy optimization algorithms. *arXiv*
- [15] J. Schulman, S. Levine, P. Moritz, M. Jordan and P. Abbeel, "Trust Region Policy Optimization", *arXiv*, 2017
- [16] R. Lian, J. Peng, Y. Wu, H. Tan, H. Zhang, "Rule-interposing deep reinforcement learning based energy management strategy for power-split hybrid electric vehicle" *Energy*, Vol. 197, 2020.
- [17] N.Barhate, "Minimal PyTorch Implementation of Proximal Policy Optimization", *GitHub*, <https://github.com/nikhilbarhate99/PPO-PyTorch>, 2021.