# A Reinforcement Learning Framework for Optimisation of Power Grid Operations and Maintenance

R. Rocchetta [1,*], M. Compare [1,2], E. Patelli [1], and E. Zio [2,3,4,5]

[1] *Institute for Risk and Uncertainty, Liverpool University, United Kingdom*
[2] *Aramis s.r.l., Milano, Italy*
[3] *Energy Department, Politecnico di Milano,Italy*
[4] *Chair on Systems Science and the Energetic challenge,Foundation Electricite' de France at CentraleSupelec, France*
[5] *Eminent Scholar, Department of Nuclear Engineering, College of Engineering, Kyung Hee University, Republic of Korea*
[*] Corresponding author: roberto.rocchetta@liverpool.ac.uk

**Abstract**

In this work, we investigate Reinforcement Learning (RL) for managing Operation and Maintenance (O&M) of power grids equipped with Prognostic and Health Management (PHM) capabilities, which allow tracking the health state of the grid components. RL exploits this information to select optimal O&M actions on the grid components giving rise to state-action-reward trajectories maximising the expected profit. A scaled-down case study is solved for a power grid, and strengths and weaknesses of the framework are discussed.

Keywords: Reinforcement Learning, Prognostic and Health Management, Operation and Maintenance, Degradation, Power Grid, Uncertainty

## 1 Introduction

Modern power grids are complex systems, including many highly interconnected components. Maximising the grid productivity while ensuring a safe and reliable delivery of power is of uttermost importance for grid operators. This requires developing robust decision-making frameworks, which give account to both the complexity of the asset and the uncertainties on its operational conditions, component degradation, failure behaviours, external environment, etc.

Nowadays, the grid management issue is further challenged by the possibility of equipping grid elements with Prognostics and Health Management (PHM) capabilities, which allow tracking the health state evolution. This information can be exploited by grid operators to further increase the profitability of their assets [1–6].

Reinforcement Learning (RL) [7, 8] has been used in the last decades to solve a variety of realistic control and decision-making issues in the presence of uncertainty, including power grid management. In the RL paradigm, a controller (i.e. the decision maker) learns from the interaction with the environment (e.g. the grid) by observing states, collecting rewards and selecting actions to maximise the future revenues, considering the aleatory uncertainties in the environment behavior. The state-action-reward trajectories [9] can be gathered from direct interaction with the real system (e.g. [10]), or from its realistic simulation [7]. This makes RL suitable to power grid management optimization, as it can cope with both the complexity of the asset and the unavoidable uncertainties related to its operation.

In [6], an RL framework based on Q-learning is proposed to solve constrained load flow and reactive power control problems in power grids. Kuznetsova et al. [5] develop an optimisation scheme for consumers actions management in the microgrid contest and accounting for renewable volatility and environmental uncertainty. In [9], a comparison between RL and a predictive control model is presented for a power grid damping problem. In [4], the authors review recent advancements in intelligent control of micro grids including few attempts using RL methods. However, none of the revised works employs RL to find optimal combined Operation and Maintenance (O&M) policies for power grids with degrading elements.

We present an RL framework to support O&M decisions for power grids equipped with PHM systems, which seeks for the settings of the generator power outputs and the scheduling of preventive maintenance actions that maximize the grid load balance and expected profit over an infinite time horizon, while considering the uncertainty of power production from Renewable Energy Sources (RES), power loads

and component failure behaviors.

The rest of this paper is organized as follows: Section 2 presents the RL framework for optimal decision making under uncertainty. A scaled-down power grid application is proposed in Section 3, whereas the results and limitations of RL are discussed in Sections 4 and 5, respectively. Section 6 closes the paper.

## 2 Modelling framework for optimal decision making under uncertainty

As anticipated above, developing a RL framework for power grid O&M management requires defining the environment, the actions that the agent can take in every state of the environment, the state transitions the actions lead to and, finally, the rewards associated to each state-action-transition step.

### 2.1 State space

Consider a power grid made up of elements $C = \{1, ..., N\}$, physically and/or functionally interconnected, according to the given grid structure. Similarly to [11], the features of the grid elements defining the environment are the $N^d$ degradation mechanisms affecting the degrading components $d \in D \subseteq C$ and the $N^p$ setting variables of power sources $p \in P \subseteq C$. For simplicity, we assume $D = \{1, ..., |D|\}$, $P = \{|D| + 1, ..., |D| + |P|\}$ and $|D| + |P| \leq N$.

The degradation processes evolve independently on each other according to a Markov process defining the transition probability from state $s_i^d(t)$ at time $t$ to the next state $s_i^d(t+1)$, where $s_i^d(t) \in \{1, ..., S_i^d\}$ $\forall t$, $d \in D, i = 1, ..., N^d$. Similarly, for the power sources production, a Markov process defines the probabilistic dynamic of power setting variables from $s_j^p(t)$ at time $t$ to the next state $s_j^p(t+1)$, where $s_j^p(t) \in \{1, ..., S_j^p\}$ $\forall t, p \in P, j = 1, ..., N^p$. Then, system state vector $\mathbf{S} \in \mathcal{S}$ at time $t$ reads:

$$\mathbf{S}(t) = \left[ \ s_1^1(t), s_2^1(t), \ldots, s_{N^{|P|+|D|}}^{|P|+|D|}(t) \ \right] \in \mathcal{S} \tag{1}$$

### 2.2 Actions

Actions can be performed on the grid components $g \in G \subseteq C$ at each $t$. The system action vector $\mathbf{a} \in \mathcal{A}$ at time $t$ is:

$$\mathbf{a}(t) = \left[ \ a_{g_1}(t), \ldots, a_{g_\varrho}(t), \ldots, a_{|g|_{|G|}}(t) \ \right] \in \mathcal{A} \tag{2}$$

were action $a_{g_\varrho}(t)$ is selected for component $g_\varrho \in G$ among a set of mutually exclusive actions $a_{g_\varrho} \in \mathbf{A}_g$. The action set $\mathbf{A}_{g_\varrho}$ can include operational actions (e.g. closure of a valve, generator power ramp up, etc.) and maintenance actions (e.g. preventive and corrective). Constraints can be defined for reducing $\mathbf{A}_{g_\varrho}$ to a subset $\hat{\mathbf{A}}_{g_\varrho} \subseteq \mathbf{A}_{g_\varrho}$. For example, Corrective Maintenance (CM), cannot be taken on As-Good-As-New (AGAN) components and, similarly, it is mandatory action for failed components. In an optimistic view [11], both Preventive Maintenance (PM) and CM actions are assumed to restore the AGAN state for each component. An example of Markov process for a 4 degradation state component is presented in Fig.1, where circle markers indicate maintenance actions and squared markers indicate other actions, i.e. operational actions.

### 2.3 Transition probabilities

Transition probability matrices are associated to each feature $f$ of each component $c \in P \cup D$ and to each action $\mathbf{a} \in \mathcal{A}$, where $f \in \{1, .., N^d\}$ if $c \in D$ and $f \in \{1, .., N^p\}$ otherwise, as follows:

$$\mathcal{P}_{c,f}^{\mathbf{a}} = \begin{bmatrix} p_{1,1} & p_{1,2} & \cdots & p_{1,S_f^c} \\ p_{2,1} & p_{2,2} & \cdots & p_{2,S_f^c} \\ \vdots & \vdots & \ddots & \vdots \\ p_{S_f^c,1} & p_{S_f^c,2} & \cdots & p_{S_f^c,S_f^c} \end{bmatrix}_{c,f}^{\mathbf{a}} \tag{3}$$

where $p_{i,j}$ represents the probability of transition from state $i$ to state $j$ of feature $f$ of component $c$ and conditional to the action $\mathbf{a}$ in a time varying setting, i.e. $\mathcal{P}_{c,f}^{\mathbf{a}}(s_j|\mathbf{a}, s_i)$. The normalization propriety holds, i.e. $\sum_{j=1}^n p_{i,j} = 1$. In practice, element $p_{i,j}$ of the transition probability matrix $\mathcal{P}_{c,f}^{\mathbf{a}}$ can be estimated as the relative frequency of the measured component state to fall into the $j^{th}$ state at time $t+1$ provided that it was at the $i^{th}$ state in the previous time step when the action $\mathbf{a}$ was taken.
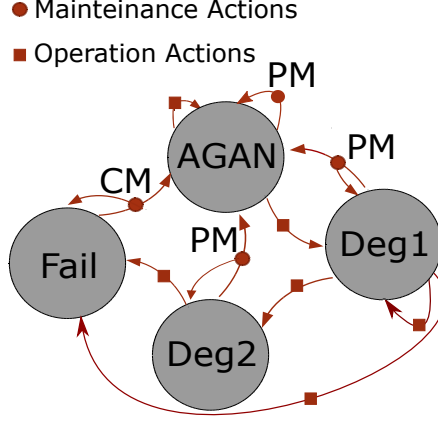
Figure 1: The Markov Decision Process associated to the health state of a degrading component.

## 2.4 Rewards

Numerical rewards are case-specific and obtained by solving a physic-economic model of the system, which evaluates how good is the transition from one state to another given that $\mathbf{a}$ is taken:

$$R(t) = \mathfrak{F}\left(\mathbf{S}(t+1)\ ,\ \mathbf{a}(t)\ ,\ \mathbf{S}(t)\right)\ \in \mathbb{R}$$

## 2.5 Reinforcement Learning and SARSA($\lambda$) method

Generally speaking, the goal of RL methods for optimal control is to find the optimal action-value function $Q_{\pi^*}(\mathbf{S}, \mathbf{a})$, which provides an estimation of future revenues when an action $\mathbf{a}$ is taken in state $\mathbf{S}$, following the optimal policy $\pi^*$:

$$Q_{\pi^*}(\mathbf{S}, \mathbf{a}) = \mathbb{E}_{\pi^*}\left[\sum_{t=0}^{\infty} R(t) | \mathbf{S}(t), \mathbf{a}(t)\right] \tag{4}$$

Among the wide range of RL algorithms, we adopt SARSA($\lambda$), which is a temporal difference learning methods (i.e. it changes an earlier estimate of $Q$ based on how it differs from a later estimate) employing eligibility traces to carry out backups over n-steps and not just over one step [7]. Details on SARSA($\lambda$) are provided in Algorithm 1 in the Appendix.

## 3 Case study

A scaled-down power grid case study is used to test the RL decision making framework. The grid includes: 2 controllable generators; 5 cables for the power transmission; 2 non-controllable RES which are connected to 2 loads and provide them electric power depending on random weather conditions (Fig. 2). Then, $|C|$=11. Two traditional generators (Gen1 and Gen2) are installed as displayed in Fig. 2 and controlled to minimize power unbalances on the grid. We assume that the 2 controllable generators and links 3 and 4 are affected by degradation and, thus, are equipped with PHM capabilities to inform the decision-maker on their degradation states, then $D = \{1, 2, 3, 4\}$. The two loads and the two renewable generators define the grid power setting, $P = \{5, 6, 7, 8\}$

### 3.1 States and Actions

In the case study, we consider $N^d = 1$ degradation features, $d = 1, .., 4$ and $N^p = 1$ power features $p = 1, .., 4$. We consider 4 degradation states for the generators, $s_1^d = \{1, .., S_1^d = 4\}$ for $d = 1, 2$, whereas three states are associated to the power lines $s_1^d = \{1, .., S_1^d = 3\}$, $d = 3, 4$. State 1 refers to the AGAN conditions, state $S_1^d$ to the failure state and states $1 < s_1^d < S_1^d$ to degraded states in ascending order. For each load, we consider 3 states of increasing power demand $s_1^p = \{1, .., S_1^p = 3\}$ for $p = 5, 6$ and three states of increasing power production are associated to renewable sources, $s_1^p = \{1, .., S_1^p = 3\}$ for
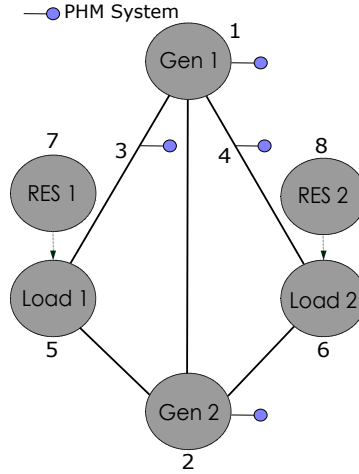
Figure 2: The power grid structure and the position of the 4 PHM capabilities, 2 renewable sources, 2 loads and 2 controllable generators.

$p = 7, 8$. Then, the total number of state vectors combinations is 11664 and the grid state vector at time $t$ is defined as follows:

$$\mathbf{S}(t) = \{s_1^1, s_1^2, s_1^3, s_1^4, s_1^5, s_1^6, s_1^7, s_1^8\}$$

The agent can operate both generators with the aim to maximise the system revenue by minimizing unbalance between demand and production, while preserving the structural and functional integrity of the system, $g \in G = \{1, 2\}$. Other actions can be performed by other agents on other components (e.g. transmission lines), but being outside from the control domain of the first agent those are assumed included in the environment. Then, the action vector reads $\mathbf{a} = [a_1, a_2]$. Five O&M actions can be performed on each controllable generator, for a total of 25 combinations, thus giving rise to a 291600 state-action pairs. The action set for each generator is the following:

$$\mathbf{A}_g = \{1, .., 5\} \quad g \in \{1, 2\}$$

where the first 3 (operational) actions affect the power output of the generator, changing it to one of the 3 allowed power levels. The last 2 actions are preventive and corrective maintenance actions, respectively. It is assumed that CM is mandatory for failed generators. Furthermore, highly degraded generators (i.e. $S_g^d = 3, \ d = 1, 2$) are assumed degraded in their operational performance and only the lower power output can be obtained (only $a_g = 1$ action is allowed). Tables 1-3 display the costs for each action and the corresponding power output of the generator, the line electric parameters and the relation between state indices $s_1^p$ and the power variable settings, respectively.

Table 1: The power output of the 2 generators in [MW] associated to the 5 available actions and action costs in monetary unit [m.u.].

| Action: | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| $P_{g=1}$ [MW] | 40 | 50 | 100 | 0 | 0 |
| $P_{g=2}$ [MW] | 50 | 60 | 120 | 0 | 0 |
| $C_{a,g}$ [m.u.] | 0 | 0 | 0 | 10 | 500 |

Table 2: The transmission lines proprieties.

| From | To | Am [A] | X |
|---|---|---|---|
| Gen 1 | Load 1 | 125 | 0.0845 |
| Gen 1 | Load 2 | 135 | 0.0719 |
| Gen 1 | Gen 2 | 135 | 0.0507 |
| Load 1 | Gen 2 | 115 | 0.2260 |
| Load 2 | Gen 2 | 115 | 0.2260 |

Table 3: The physical values of the power settings in [MW] associated to each state $S_1^p$ of component $p \in P$.

| | State index $s_1^p$ | 1 | 2 | 3 |
|---|---|---|---|---|
| $p = 5$ | Demanded [MW] | 60 | 100 | 140 |
| $p = 6$ | Demanded [MW] | 20 | 50 | 110 |
| $p = 7$ | Produced [MW] | 0 | 20 | 30 |
| $p = 8$ | Produced [MW] | 0 | 20 | 60 |

## 3.2 Probabilistic Model

State transitions may occur from time $t$ to the next time step $t+1$ and are specifically defined for each feature of each component. The 2 loads have identical transition probability matrices and also the degradation of the transmission cables and generators are described by the same Markov process. Thus, for ease of notation, the component subscripts have been dropped. Each action $\mathbf{a} \in \mathcal{A}$ is associated to a specific transition probability matrix $\mathcal{P}_g^{\mathbf{a}}$ describing the evolution of the generator health state conditioned by its operative state or maintenance action. It can be noticed that probabilities associated to operational actions, namely $a_g = 1, 2, 3$, affect differently the degradation of the component. For those actions, the bottom row corresponding to the failed state has only zero entries. This is to indicate that operational actions cannot be taken for failed generators, but only CM is allowed. The transition matrices for the considered features are defined as follows:

$$\mathcal{P}_d^{a_d=1} = \begin{pmatrix} 0.98 & 0.02 & 0 & 0 \\ 0 & 0.95 & 0.05 & 0 \\ 0 & 0 & 0.9 & 0.1 \\ 0 & 0 & 0 & 0 \end{pmatrix} d = 1,2 \qquad \mathcal{P}_d^{a_d=2} = \begin{pmatrix} 0.97 & 0.03 & 0 & 0 \\ 0 & 0.95 & 0.05 & 0 \\ 0 & 0 & 0.9 & 0.1 \\ 0 & 0 & 0 & 0 \end{pmatrix} d = 1,2$$

$$\mathcal{P}_d^{a_d=3} = \begin{pmatrix} 0.95 & 0.04 & 0.01 & 0 \\ 0 & 0.95 & 0.04 & 0.01 \\ 0 & 0 & 0.97 & 0.03 \\ 0 & 0 & 0 & 0 \end{pmatrix} d = 1,2 \qquad \mathcal{P}_d^{a_d=4} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0.5 & 0 & 0.5 & 0 \\ 0.5 & 0 & 0 & 0.5 \\ 0 & 0 & 0 & 0 \end{pmatrix} d = 1,2$$

$$\mathcal{P}_d^{a_d=5} = \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0.15 & 0 & 0 & 0.85 \end{pmatrix} d = 1,2$$

$$\mathcal{P}_d^{\mathbf{a}} = \begin{pmatrix} 0.9 & 0.08 & 0.02 \\ 0 & 0.97 & 0.03 \\ 0.1 & 0 & 0.9 \end{pmatrix} \forall \, \mathbf{a}, d = 3,4 \qquad \mathcal{P}_p^{\mathbf{a}} = \begin{pmatrix} 0.4 & 0.3 & 0.3 \\ 0.3 & 0.3 & 0.4 \\ 0.2 & 0.4 & 0.4 \end{pmatrix} \forall \, \mathbf{a}, p = 5,6$$

$$\mathcal{P}_7^{\mathbf{a}} = \begin{pmatrix} 0.5 & 0.1 & 0.4 \\ 0.3 & 0.3 & 0.4 \\ 0.1 & 0.4 & 0.5 \end{pmatrix} \forall \, \mathbf{a} \qquad \mathcal{P}_8^{\mathbf{a}} = \begin{pmatrix} 0.5 & 0.2 & 0.3 \\ 0.4 & 0.4 & 0.2 \\ 0 & 0.5 & 0.5 \end{pmatrix} \forall \, \mathbf{a}$$

## 3.3 Reward Model

When the agent performs an action at time $t$; the environment provides a reward and leads the system to its state at time $t+1$. The reward is calculated as the sum of 4 different terms: (1) the revenue from selling electric power, (2) the cost of producing electric power by traditional generators, (3) the cost associated to the performed actions and (4) the cost of not serving energy to the customers. Mathematically, the reward reads:

$$R(t) = \sum_{p=5}^{6} (L_p(t) - ENS_p(t)/\Delta_t) \cdot C_{el} - \sum_{g=1}^{2} P_g \cdot C_g - \sum_{g=1}^{2} C_{a,g} - \sum_{p=5}^{6} ENS_p(T) \cdot C_{ENS} \qquad (5)$$

where $L_p$ is the power demanded by component $p$, $C_{el}$ is the price paid by the loads for per-unit of electric power, $P_g$ is the power produced by the generators, $C_g$ is the cost of producing the unit of power, $C_{a,g}$ is the cost of the action a on the generator $g$, $\Delta_t$ is the time difference between the present and the next system state and it is assumed to be 1 h, $ENS_p$ is the energy not supplied to the load $p$ and is a function of the grid state vector and lines and generators electrical proprieties and availability, i.e. $\mathbf{ENS}(t) = \mathcal{G}(\mathbf{S}, \mathbf{Am}, \mathbf{X})$ where $\mathcal{G}$ defines the constrained DC power flow solver [12]. $C_{ENS}$ is the cost of the energy not supplied. The costs $C_{ENS}$, $C_g$ and $C_{el}$ are set to 5, 4 and 0.145 monetary unit (m.u.) per-unit of energy or power, respectively.

## 4 Results and Discussions

The SARSA($\lambda$) algorithm (Algorithm 1 in the Appendix) has been used to provide an approximate solution to the decision problem. The stochastic grid model is used to sample control trajectories only, i.e. it provides a reward and a new state when an action and the old state is provided as input. The SARSA method has been run changing parameters setting and accumulating eligibility traces. The initial state

$s = \mathbf{S}(t = 0)$ has been selected for the episodic loop randomly, using a degradation-weighted probability mass function $f_S(s) \propto \sum_{d=1}^{|D|} s_1^d$. This sampling scheme is used to better estimate action-value functions in rarely visited sates (i.e. low-probability states with many failed/highly degraded components), which speeds up the convergence of the SARSA method. For validation, Bellman's optimality [13]-[14] has been solved to provide a reference optimal action-value function. The Bellman's results are in good agreement with the SARSA results, as it can be seen from Fig. 3.
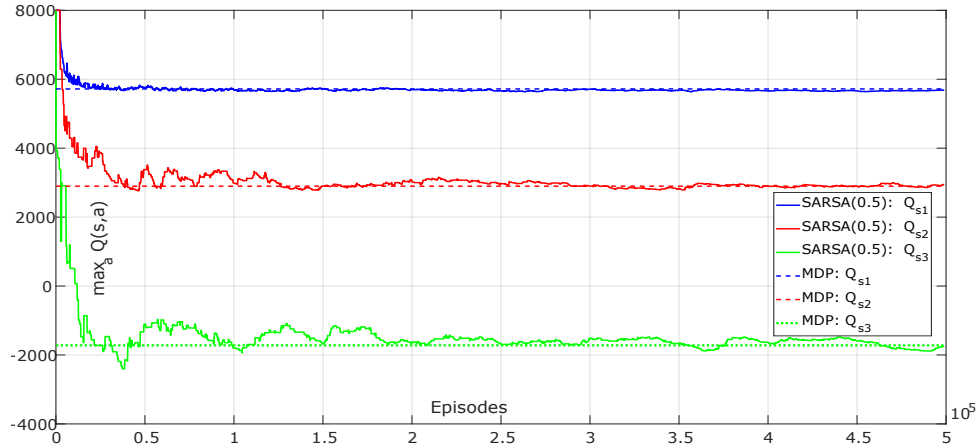


Figure 3: The plot shows a comparison of the maximum $Q_{\pi^*}(\mathbf{S}, \mathbf{a})$ for 3 states indicative of the different state-action value levels, obtained by SARSA(0.5) algorithm and $T = 50$ (solid lines) and the reference Bellman's solution of the underlying Markov Decision Process (dashed lines).

The SARSA($\lambda$) results are summarised in Fig. 4, where the curves provide a compact visualization of the distribution of $Q_{\pi^*}(\mathbf{S}, \mathbf{a})$ over the states for the available 25 combinations of actions. Three clusters can be identified: on the far left, we find the set of states from which CM on both generators is performed; being CM a costly action, this leads to a negative expectation of the discounted reward. The second cluster ($C\ 2$) corresponds to the 8 combination of one CM and any other action on the operating generator. The final cluster ($C\ 1$) of 16 combinations of actions includes only PM and operational actions. If corrective maintenance is not performed, higher rewards are expected.
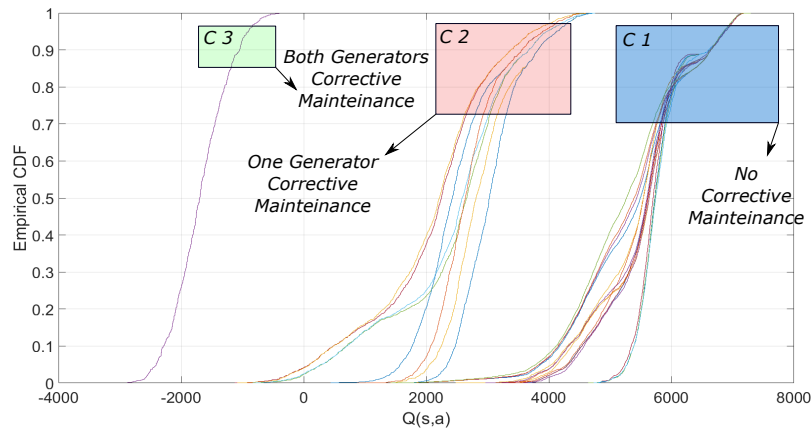


Figure 4: The $Q(s, a)$ values displayed using ECDFs and the 3 clusters.

In Fig. 5, each sub-plot shows the the highest expected discounted power grid return, $Q_{\pi^*}(\mathbf{S}, \mathbf{a})$, adopting the optimal policy, conditional to a specific degradation states of the generators and for increasing electric load demand. It can be noticed that if the generators are both healthy or slightly degraded (i.e. $\sum_{d=1}^{2} s_1^d = 2, 3, 4$) an increment in the overall load demand leads to an increment in the expected reward, due to the larger revenues from selling more electric energy to the customers. On the other hand,

if the generators are highly degraded or failed (i.e. $\sum_{d=1}^{2} s_1^d = 7, 8$), an increment in the load demand leads to a drop in the expected revenue. This is due to the increasing risk of load curtailments and associated cost (i.e. cost of energy not supplied), and to the impacting PM and CM actions costs.
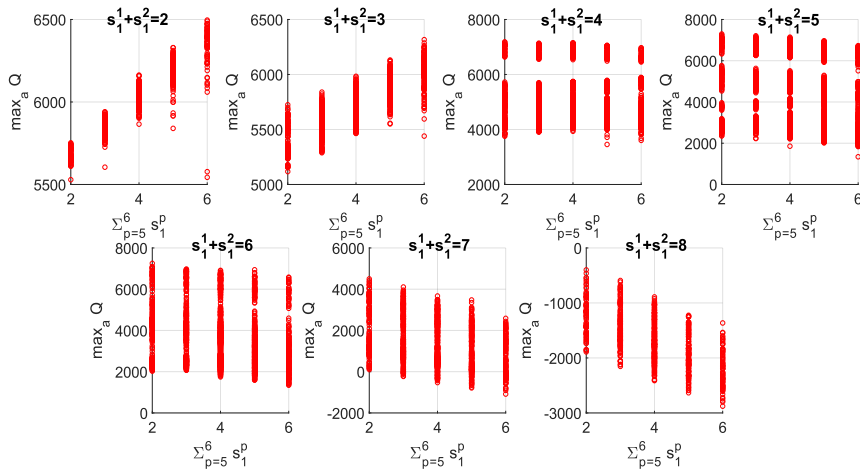


Figure 5: The maximum $Q_{\pi^*}(\mathbf{S}, \mathbf{a})$ (i.e. maximum expected discounted cumulative reward) for increasing total load and different degrading condition of the generators.

## 4.1 Policies comparison

We have empirically found that SARSA(0.5) policies outperform SARSA($\lambda$), $\lambda = 0$ and $\lambda = 1$. Thus, two SARSA(0.5) have been further investigated, by setting the truncation windows $T$ to 50 and 250 time steps for each episode, respectively. Table 4 shows the results of the SARSA($\lambda$) algorithms (columns 3 and 4, respectively) and compares them with the MDP (Bellman's optimality) solution (column 2) and 2 artificial suboptimal policies: $Q_{50rnd}$ (column 5), which is artificially obtained randomizing the action to be selected in 50 % of the states and selecting the MDP optimal action for the remaining states and $Q_{100rnd}$ (column 6), where all states have a random action associated with. Three representative system states $\mathbf{S}_1 = [1, 1, 1, 1, 1, 1, 1, 1]$, $\mathbf{S}_2 = [4, 1, 1, 1, 1, 1, 1, 1]$ and $\mathbf{S}_3 = [4, 4, 3, 3, 3, 3, 3, 3]$ are used to compare the expected discounted return $Q$. The 3 states are associated with substantially different rewards as they have been selected from the 3 clusters $C\ 1$, $C\ 2$ and $C\ 3$, respectively (see Fig. 4): $\mathbf{S}_1$ has both generators in the AGAN state, $\mathbf{S}_2$ has on generator out of service whilst $\mathbf{S}_3$ has both generators failed. $Act$ is defined as the portion of actions taken from the SARSA($\lambda$) policies that are equal to those taken using the reference MDP optimal policy in the corresponding states; $\mathbb{E}[R(t)]$ is the expected averaged non-discounted return, i.e. $\mathbb{E}\left[\frac{\sum_{t=1}^{T} R(t)}{T}\right]$, independent from the initial state of the system. It is interesting to notice that SARSA(0.5) provides better policies (i.e. higher expected discounted and non-discounted returns) compared to $Q_{50rnd}$ and $Q_{100rnd}$. This is true even if $Q_{50rnd}$ has higher $Act$ compared to the SARSA policies, i.e. more than 60 % of the $Q_{50rnd}$ actions are equal to the MDP actions whilst less than 50 % for the SARSA. This points out that the optimal policy is very sensitive to some of the state-action combinations and less to others. In other words, taking the wrong action in some states can lead to a catastrophic drop in the expected return, whilst in other cases a sub-optimal action affects less the expected revenue (e.g. making generator 1 produce power rather than generator 2 or vice versa).

Fig. 6 presents in details 2 control trajectories obtained selecting greedily actions with the MDP Bellman's policy (top plot) and the SARSA($\lambda$) policy (bottom plot), rewards are displayed on the y-axis and actions and states (see Table 5) are associated to each time step. It is interesting to observe that by following an optimal policy, PM actions are sometimes recommended even if the generators are As-Good-As-New. This might seem counter intuitive, but it can be explained considering the degradation model settings. A PM action taken in an AGAN degradation state will assure a transition to the AGAN state. In this sense, the MDP policy is ready to accept a slightly lower revenue (due to PM costs), but with the advantage of suspending the degradation process, especially when the power produced by RES can be used to minimise unbalances between power production and the 2 loads are small.

Table 4: The MDP Bellman's optimality and the RL results compared with suboptimal policies.

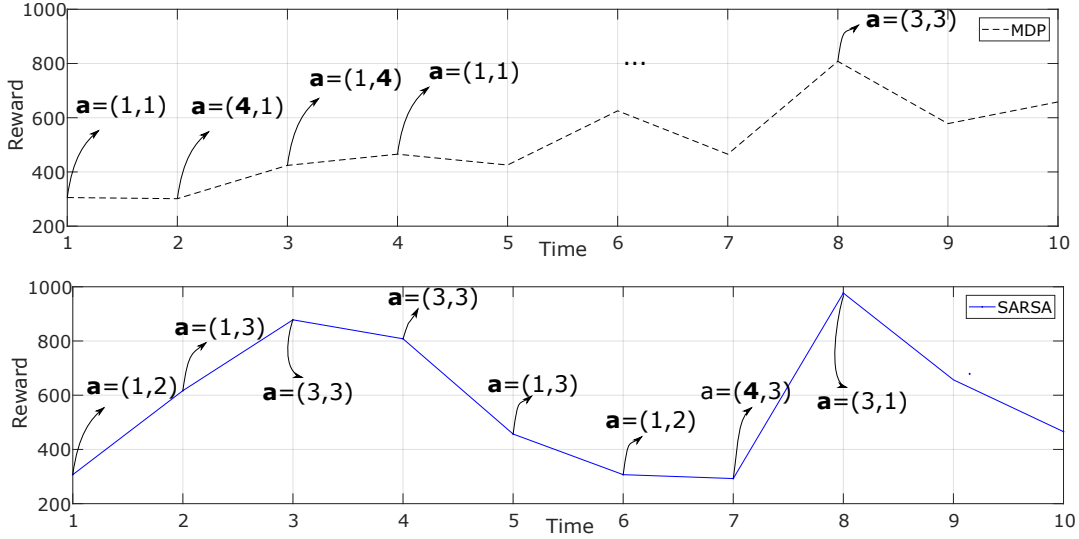| | MDP | SARSA(0.5) | | $Q_{50rnd}$ | $Q_{100rnd}$ |
|---|---|---|---|---|---|
| $Q_{\mathbf{S}_1}$ | 5719 | 5511 | 5555 | 4191 | 2028 |
| $Q_{\mathbf{S}_2}$ | 2898 | 2577 | 2664 | 1297 | -1229 |
| $Q_{\mathbf{S}_3}$ | -1721 | -1816 | -1813 | -2956 | -4288 |
| $Act$ top1 | 100 % | 48.8 % | 49.1 % | 62.1% | 24.8% |
| $Act$ top3 | 100 % | 66.5 % | 66.5 % | 71.4% | 43.1% |
| $\mathbb{E}[R(t)]$ | 529.8 | 478.8 | 488.1 | 370.3 | 190.4 |
| $N_e$ | - | 5e5 | 5e5 | - | - |
| T | - | 50 | 250 | - | - |



Figure 6: Actions taken in 2 separate control trajectories using MDP and SARSA policies. Initial state $s_1$ and next states are randomly generated by the underlying probabilistic model (see Table 5).

Table 5: The state vectors for the MDP and SARSA control trajectories in Figure 6.

| MDP states trajectory | | | | | | | |
|---|---|---|---|---|---|---|---|
| Gens | | Loads | | RES | | Lines | |
| $s_1^1$ | $s_1^2$ | $s_1^5$ | $s_1^6$ | $s_1^7$ | $s_1^8$ | $s_1^3$ | $s_1^4$ |
| 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| 1 | 1 | 1 | 1 | 3 | 1 | 1 | 1 |
| 1 | 1 | 1 | 2 | 2 | 3 | 1 | 1 |
| 1 | 1 | 2 | 1 | 3 | 2 | 1 | 1 |
| 1 | 1 | 1 | 2 | 3 | 1 | 1 | 1 |
| 1 | 1 | 3 | 1 | 3 | 3 | 3 | 1 |
| 2 | 1 | 2 | 1 | 3 | 2 | 1 | 1 |
| 2 | 2 | 2 | 3 | 1 | 2 | 1 | 1 |
| 2 | 2 | 2 | 2 | 1 | 1 | 1 | 1 |
| 2 | 2 | 1 | 3 | 1 | 2 | 1 | 1 |

| SARSA states trajectory | | | | | | | |
|---|---|---|---|---|---|---|---|
| Gens | | Loads | | RES | | Lines | |
| $s_1^1$ | $s_1^2$ | $s_1^5$ | $s_1^6$ | $s_1^7$ | $s_1^8$ | $s_1^3$ | $s_1^4$ |
| 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| 1 | 1 | 3 | 1 | 3 | 1 | 1 | 1 |
| 1 | 1 | 3 | 3 | 2 | 1 | 1 | 1 |
| 1 | 1 | 2 | 3 | 2 | 1 | 1 | 1 |
| 1 | 1 | 2 | 1 | 1 | 1 | 3 | 3 |
| 1 | 1 | 1 | 1 | 1 | 2 | 1 | 3 |
| 1 | 1 | 1 | 1 | 1 | 2 | 2 | 3 |
| 1 | 1 | 3 | 3 | 3 | 3 | 2 | 3 |
| 1 | 1 | 1 | 3 | 2 | 2 | 2 | 1 |
| 1 | 1 | 2 | 1 | 3 | 2 | 2 | 1 |

## 5 Discussion on Limitation

While RL, like stochastic dynamic programming (DP), has in principle a very broad scope of application, it has to face computational issues when the state-action spaces of the control problem are very large. In such a case, RL has to be combined with regression techniques capable of interpolating over the state-action space the data obtained from (relatively) few control trajectories [9]. Most of the research in this context has focused on parametric function approximators, representing either some (state-action) value functions or parameterized policies, together with some stochastic gradient descent algorithms (see e.g. [8] or [15]).

Further research work will focus on the development of enhanced RL algorithms, capable of dealing with

imprecise rewards (e.g. due to unavailable/unreliable models), partial observability and issues related to scarcity of samples due to low-probability of specific state-action pairs.

## 6 Conclusion

A framework based on Reinforcement Learning for optimal decision making of power grid systems affected by uncertain operations and degradation mechanisms has been investigated. Power grid models can include PHM devices, which are used to inform the agent about the health state of the system components. This information helps to select optimal O&M actions on the system components.

The SARSA($\lambda$) method was used to solve a control problem for a scale down power grid with renewable and PHM capabilities. The RL results have been compared to the reference Bellman's optimality solution and are in good agreement, although inevitable approximation errors have been observed.

The framework proved to be flexible and effective in tackling a small but representative case study and future works will test its applicability to more realistic (larger) state-action spaces. To this aim, artificial neural networks will be used in future research work for state-action space regression to scale up to larger grids. This necessary verification for a possible future applicability of the method.

## Acknowledgments

## References

[1] M. Papageorgiou, C. Diakaki, V. Dinopoulou, A. Kotsialos, and Y. Wang, Review of road traffic control strategies, *Proceedings of the IEEE* **91**, 2043 (2003), ISSN 0018-9219.

[2] J. Jin and X. Ma, Hierarchical multi-agent control of traffic lights based on collective learning, *Engineering Applications of Artificial Intelligence* **68**, 236 (2018), ISSN 0952-1976.

[3] R. Yousefian, R. Bhattarai, and S. Kamalasadan, Transient stability enhancement of power grid with integrated wide area control of wind farms and synchronous generators, *IEEE Transactions on Power Systems* **32**, 4818 (2017), ISSN 0885-8950.

[4] M. S. Mahmoud, N. M. Alyazidi, and M. I. Abouheaf, Adaptive intelligent techniques for microgrid control systems: A survey, *International Journal of Electrical Power & Energy Systems* **90**, 292 (2017), ISSN 0142-0615.

[5] E. Kuznetsova, Y.-F. Li, C. Ruiz, E. Zio, G. Ault, and K. Bell, Reinforcement learning for microgrid energy management, *Energy* **59**, 133 (2013), ISSN 0360-5442.

[6] J. G. Vlachogiannis and N. D. Hatziargyriou, Reinforcement learning for reactive power control, *IEEE Transactions on Power Systems* **19**, 1317 (2004), ISSN 0885-8950.

[7] R. S. Sutton, D. Precup, and S. Singh, Between mdps and semi-mdps: A framework for temporal abstraction in reinforcement learning, *Artificial Intelligence* **112**, 181 (1999), ISSN 0004-3702.

[8] C. Szepesvari, *Algorithms for Reinforcement Learning* (Morgan and Claypool Publishers, 2010), ISBN 1608454924, 9781608454921.

[9] D. Ernst, M. Glavic, F. Capitanescu, and L. Wehenkel, Reinforcement learning versus model predictive control: A comparison on a power system problem, *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)* **39**, 517 (2009), ISSN 1083-4419.

[10] D. P. Bertsekas and J. N. Tsitsiklis, *Neuro-Dynamic Programming* (Athena Scientific, 1996), 1st ed., ISBN 1886529108.

[11] M. Compare, P. Marelli, P. Baraldi, and E. Zio, A markov decision process framework for optimal operation of monitored multi-state systems, *Proceedings of the Institution of Mechanical Engineers Part O Journal of Risk and Reliability* (2018).

[12] R. Rocchetta and E. Patelli, Assessment of power grid vulnerabilities accounting for stochastic loads and model imprecision, *International Journal of Electrical Power & Energy Systems* **98**, 219 (2018), ISSN 0142-0615.

[13] R. Bellman, A markovian decision process, *Journal of Mathematics and Mechanics* **6**, 679 (1957).

[14] E. Gross, On the bellmans principle of optimality, *Physica A: Statistical Mechanics and its Applications* **462**, 217 (2016), ISSN 0378-4371.

[15] H. Li, T. Wei, A. Ren, Q. Zhu, and Y. Wang, in *2017 IEEE/ACM International Conference on Computer-Aided Design (ICCAD)* (2017), 847–854.

[16] S. Singh, T. Jaakkola, M. L. Littman, and C. Szepesvári, Convergence results for single-step on-policy reinforcement-learning algorithms, *Machine Learning* **38**, 287 (2000), ISSN 1573-0565.

## Appendix

The SARSA($\lambda$) algorithm starts initializing the action-value function $Q$ and eligibility traces $Z$ tables. Then, the values for the learning rate $\alpha$, the discount factor $\gamma$, the decay rate of the traces $\lambda \in [0, 1]$ and the greediness factor $\epsilon$ (or a policy $\pi$ to be evaluated) are selected. After this initialization, the episodic loop starts with a random sample (or selection) of an initial state $s_t$, then, an action $a_t$ is selected based on the adopted policy, e.g. $\epsilon$-greedy or $\pi(\cdot|s_t)$. A $\epsilon$-greedy policy consists of random actions, taken with probability $\epsilon$, or greedy actions taken with probability 1-$\epsilon$ (i.e. actions for which $Q$ is maximised). Once the initial state-action pair is obtained, the episode $e$ is evaluated (i.e. a sequence of action-rewards-state-actions). Temporal difference errors $\delta_t$ at the time step $t$ are calculated, traces replaced or accumulated and $Q$ updated.The episode terminates when a predefined truncation horizon $T$ is reached (i.e. maximum time length of the episode). The procedure is iterated until a predefined number of events $N_E$ is obtained. The SARSA(0) is guaranteed to convergence to an optimal action-value function for a Robbins-Monro sequence of step-sizes $\alpha_t$, for further details regarding stopping criteria and convergence the reader is referred to [16]. RL approaches can tackle control problems with infinite optimisation horizon by approximating the solution with a T-stage approach. In this sense, windows of $T$ time steps are used to truncating the time horizon, thus reducing the computational burdens [9]. The SARSA($\lambda$) algorithm works as follows [7]:

---

**Data:** Set $e = 1$, $N_E$, $\epsilon$ (or a policy $\pi$ to be evaluated), $\alpha$, $\gamma$, $\lambda$;
Initialize $Q(s, a)$, for all $s \in S$ and $a \in A$, arbitrarily (e.g. $Q = 0$);
Initialize traces $Z(s, a) = 0$, for all $s \in S$ and $a \in A$;
**while** $e < N_E$ *(Episodic Loop)* **do**
    Set $t = 1$;
    Initialize starting state $s_t$ e.g. randomly;
    Select action $a_t \in A(s_t)$ using policy derived from $Q$ (e.g. $\epsilon$-greedy) or $\pi(\cdot|s_t)$;
    **while** $t < T$ *(run an episode)* **do**
        Take action $a_t$, observe $s_{t+1}$ and reward $R_t$;
        Select action $a_{t+1} \in A(s_{t+1})$ using policy derived from $Q$ (e.g. $\epsilon$-greedy) or $\pi(\cdot|s_{t+1})$;
        Compute temporal difference $\delta_t$ and update traces: $\delta_t = R_t + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)$;
        $Z(s_t, a_t) = Z(s_t, a_t) + 1$ (accumulate traces) or;
        $Z(s_t, a_t) = 1$ (replace traces);
        Update $Q$ and $Z$ for each $s$ and $a$: $Q(s, a) = Q(s, a) + \alpha \delta_t Z(s, a)$;
        $Z(s, a) = \gamma \lambda Z(s, a)$;
        Set $t = t + 1$;
    **end**
    go to next episode $e = e + 1$;
**end**

**Algorithm 1:** The SARSA($\lambda$) algorithm adopting replacing or accumulating eligibility traces settings.