

Optimal Storage Scheduling Using Markov Decision Processes

Samuele Grillo, *Member, IEEE*, Antonio Pievatolo, and Enrico Tironi

Abstract—The paper presents a method based on Markov decision processes to optimally schedule energy storage devices in power distribution networks with renewable generation. The time series of renewable generation is modeled as a Markov chain which allows for the implementation of a stochastic dynamic programming algorithm. The output of this algorithm is an optimal scheduling policy for the storage device achieving the minimization of an objective function including cost of energy and network losses. Besides this, other properties, such as energy storage placement and size, can be assessed and compared in optimized systems with different layouts.

Index Terms—Energy storage, Markov processes, renewable generation, stochastic dynamic programming.

NOMENCLATURE

t	index of an epoch
γ_t	load
π	a control policy
a_t	energy drawn from, or supplied to, the storage device
b_t	state of charge
c_t	clearness
d_t	decision function from the state space to the action space
E	nominal capacity of the storage device
n_{clr}	number of clearness levels
n_{SOC}	number of states of charge
r_t	reward function
s_t	system state
$u_t^\pi(s_t)$	value-to-go at epoch t for policy π and current system state s_t
w_t	system energy losses
z_t	price of energy

I. INTRODUCTION

THE increasing presence of electrical storage systems in power networks requires to optimize their management [1], [2]. Network losses and the price of energy are two key

Manuscript received April 15, 2015; revised August 09, 2015 and October 07, 2015; accepted October 30, 2015. Date of publication November 26, 2015; date of current version March 18, 2016. Paper no. TSTE-00294-2015.

S. Grillo and E. Tironi are with the Dipartimento di Elettronica, Informazione e Bioingegneria, Politecnico di Milano, I-20133 Milan, Italy (e-mail: samuele.grillo@polimi.it; enrico.tironi@polimi.it).

A. Pievatolo is with CNR-IMATI, I-20133 Milan, Italy (e-mail: antonio.pievatolo@mi.imati.cnr.it).

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

factors that must be taken into account when solving this decision problem at a minimum cost, with the additional complication of randomness of both the power supply (especially if provided by renewable sources) and the load.

There are several contributions to the coordinated management of energy storage devices and renewable sources. [3]–[6] consider the renewable generation profile as known with no uncertainty. As for the algorithmic approach to storage management optimization, [3], [4] use dynamic programming, whereas in [5], [6] the optimization is performed simultaneously for all decision epochs by means of mixed integer linear programming and particle swarm optimization, respectively. However, in order to include the intermittent nature of renewable energy sources and the sequential character of the decision problem in the analysis, a stochastic dynamic programming approach is to be preferred. Therefore, in this work, the optimal storage management policy (the optimal policy from now on) is based on Markov decision processes (MDP, see [7]), assuming a random renewable energy source connected to the grid and a deterministic load profile. We solved a 24-hour finite-horizon optimal policy problem, in which the solar energy produced every 15-minute interval is obtained through a Markov chain model of the clearness, as done by [8].

The Markov chain framework is attractive for storage management policy because the state of charge (SOC) at discrete time steps can be regarded as a Markov chain. In fact, the SOC at any time, given all the past SOC's, can be written as a function of the previous SOC and some random variables which do not depend on the past. The surrounding processes, which cause the SOC to change, can be non-Markov, however a Markov chain could approximate the true process satisfactorily enough. This viewpoint is shared by other authors. For example in [9] an MDP approach to the management of the batteries of a solar-powered sensor network which maximizes the ability of the sensor node to detect and transmit an event of interest, while aiming at preserving the battery energy level, has been proposed. In [10] an enhanced stochastic dynamic programming for the minimization of electric vehicles' charging costs in a smart grid framework. The enhancement consists in the possibility of using a continuous space of decision variables. In [11] authors addressed the problem of organizing energy storage purchases (for households and data centres) to minimize long-term energy costs under variable demands and prices in an MDP framework. They found that the optimal policy has a two-threshold structure, that is, the battery must be discharged when the SOC is above a certain threshold and must be charged when the SOC is below another threshold. They have also briefly suggested how to incorporate a renewable energy source and

battery replacement costs within the optimization problem, but have not implemented the suggested solution. In [12] a simplified solution to the same problem has been obtained, by considering the price the only (Markovian) stochastic quantity and by relaxing ramping constraints on storage charging and discharging; they found an optimal policy based on a threshold for the price of energy. The distinctive feature of the MDP methodology is that it guarantees that no policy can obtain a better performance under non-restrictive assumptions, as shown for example by [7, Section IV.4]. Coming to the subject matter of this work, [13] have proven that a dynamic offset policy for storage scheduling based on a MDP outperforms heuristic methods, such as those aiming at maintaining a constant level of the storage.

In [14], the authors addressed the day-ahead wind commitment problem with storage, in a stochastic dynamic programming framework. At each decision epoch, the wind farm manager must pledge the electricity that the wind farm must produce the next 24 hours, knowing the 24 hourly market prices of electricity. It is assumed that the transition of the state variable, which includes the wind as its only stochastic part, is Markovian, reducing their problem to the MDP framework, with a continuous state space and in high dimension. The proposed solution is based on a discretization of the state-space obtained from a finite sample of sequences of states over the decision horizon and on a convex approximation of the value-to-go function, so that it can be easily optimized for each state of the system.

Our approach belongs to the same framework of [14], however it solves the storage scheduling optimization problem exactly, because the state and action spaces are both discrete and low dimensional: after executing the policy optimization, one is left with a decision table which, for every decision epoch, matches each state of the system (made up of both the actual SOC of the storage device and the observed radiation) with an optimal action, i.e., how much energy should be exchanged between the grid and the storage device. The optimal policy is associated with an expected revenue. Other differences can be found in the emphasis given to the various facets of the problem. The renewable source is the solar energy, which, more than wind, is correlated with the load, because most of the energy demand occurs during daylight hours. The effect of this relationship will appear in the optimal storage scheduling solution, in which the charge and discharge cycles tend to be anticyclic with respect to the price of energy, in the test case we have examined. Finally, we have also focused on the number of full storage cycles, having found that the optimization with respect to cost of energy and network losses has the side effect of reducing this number, in certain system configurations [15], while the differences between the optimal cost and the cost associated with other solutions are not large, in absolute terms. Other authors have also considered network losses as a performance measure. For example, deterministic dynamic programming has been used by [16] to define an optimal control strategy of distributed energy storages controlled by a distribution system operator. The optimality criterion is to minimize network losses, which are then used to compare different distributed energy storage placements and ratings.

The method illustrated in this article can be employed in at least two ways. The first way as an analysis tool for use by a utility company, for example for assessing what one can gain by including a storage in the system, by comparing the optimal revenue with the storage against that without it. In this case, the Markov model of the clearness for representative periods in the year is estimated from historical data. In the second way, this method provides an actual decision tool for one-day-ahead storage system management planning, where the Markov model for the clearness would rather be obtained from weather forecasts for the next day. Here we have analysed the first usage, evaluating the just mentioned optimal performance measures (revenues and number of cycles) for different storage system sizes.

The plan of the paper is as follows. In Section II we describe the MDP framework applied to our problem, we define our objective function and the action space, and we illustrate the dynamic programming algorithm, which has been obtained by joining a deterministic (during night hours) and a stochastic part (during solar production hours). In Section III we explain our case study, introducing the test network and the parameter used for the solar radiation model based on the clearness. In Section IV we describe our results. Section V contains concluding remarks, emphasising uses to which the algorithm could be put to, other than simple storage scheduling.

II. THE SETUP FOR OPTIMAL POLICY FINDING

In this section we illustrate the main variables and the method for finding an optimal policy.

As done by [8], the solar power at epoch t is expressed in terms of the clearness c_t , which is the square root of the ratio between the observed radiation and the expected radiation. The advantage of this approach is that daily and seasonal patterns in both the mean and the variability of solar radiation are removed and the clearness can be modelled as a stationary autoregressive Gaussian time series, with a bimodal marginal distribution, with one primary mode representing clear sky and a less peaked secondary mode representing a range of cloud cover conditions. From Glasbey's model [8] we derived a Markov chain discretizing the clearness to 14 levels, because conditions for the existence of an optimal Markovian policy are more easily met when the state space is countable (see [7, Proposition 4.4.3]). The Markov model holds from 8 am to 4 pm. During the remaining hours no production of solar energy is assumed.

In the MDP framework one is required to specify a state vector, an action space and a reward function. We denote the state at epoch t as $s_t = (c_t, b_t, \gamma_t, z_t)$, where c_t is the clearness, b_t is the SOC, γ_t is the load and z_t is the price of energy. The (discretized) action space is A : an action $a_t \in A$ specifies how much energy should be charged or discharged from the storage. Prices and loads are assumed to be deterministic, in order to reduce the numerical complexity of the algorithm and to avoid the curse of dimensionality, whereas b_t is stochastic, but, given s_t , b_{t+1} is completely determined. The reward function $r_t(s_t, a_t)$ depends on both the state and the action. We will consider two forms of reward functions: the first one takes into

account only energy losses (w_t) in the network; the second one contains the losses, but multiplied by the price of the energy.

$$r_t(s_t, a_t) = -w_t \quad (1)$$

or

$$r_t(s_t, a_t) = -z_t w_t. \quad (2)$$

Then, for the period from 8 am to 4 pm, subdivided into $N - 1 (= 32)$ 15-minute time intervals, an optimal finite-horizon policy is readily obtained from the backward induction algorithm: for every epoch t and every state s_t , the algorithm provides a function which maps the state into the action space, so that the expected total reward for the entire decision-making horizon is maximal.

This optimal policy must be connected to the optimal policy during the dark hours, that is, from 4 pm to 8 am the next day in our case. As there are no stochastic sources in this period, we have used deterministic dynamic programming, starting at 7:45 am and proceeding backwards until 4 pm the previous day, moving across $M - 1 (= 64)$ 15-minute time intervals.

In a mathematical form, the above amounts to maximizing an expected total reward

$$\mathbb{E}_{s_1} \left\{ \sum_{t=1}^{N-1} r_t(s_t, a_t) \right\} + \sum_{t=1}^{M-1} \tilde{r}_t(\tilde{s}_t, \tilde{a}_t) + \tilde{r}_M \quad (3)$$

for any initial state s_1 , with respect to the sequence of actions $(a_1, \dots, a_{N-1}, \tilde{a}_1, \dots, \tilde{a}_{M-1})$, where the tilde accent identifies the quantities subject to no stochastic input, and \tilde{r}_M is the constant terminal reward at 8 am the next day. Every action in the action sequence is constrained in two ways: first, by an upper and a lower SOC limit and by an upper limit on the current flowing from and to the storage; second, by the admissibility of the power flow. Details are given in Section III.

A. Formalization of the Optimization Problem

The state s_t is aleatory but observable and is function of clearness, SOC of the storage device and power requested by loads:

$$s_t = (c_t, b_t, \gamma_t). \quad (4)$$

Let $a_t \in A_{s_t}$ be the action to be chosen at epoch t . As previously stated, a revenue $r_t = f(s_t, a_t)$ is associated to every action a_t from state s_t .

Within every period the variables are assumed to be constant. Under this hypothesis at epoch t the revenue can be calculated.

The action a_t makes the system go from state s_t to state s_{t+1} and is chosen according to a transition kernel made up of four terms:

- $p_t(c_{t+1}|c_t)$ the probability that there is a certain clearness c_{t+1} , given that the clearness at epoch t was c_t ;
- $p_t(b_{t+1}|a_t, s_t)$ the probability that the storage system reaches the state of charge b_{t+1} , given that the starting state is s_t and the the control action is a_t ;

- $p_t(\gamma_{t+1}|\gamma_t)$ the probability that the load has a certain value γ_{t+1} , given that the load at epoch t was γ_t ;
- $p_t(z_{t+1}|z_t)$ the probability that the energy price is z_{t+1} , given that the price at epoch t was z_t .

Thus, the transition kernel, which is the probability that the system goes from state s_t to state s_{t+1} due to the application of action a_t , can be written, in the most general form, as:

$$\begin{aligned} p_t(s_{t+1}|s_t, a_t) &= p_t(c_{t+1}|c_t) p_t(b_{t+1}|a_t, c_t, b_t, \gamma_t) \\ &\quad p_t(\gamma_{t+1}|\gamma_t) p_t(z_{t+1}|z_t) \\ &= p_t(c_{t+1}|c_t) p_t(b_{t+1}|a_t, s_t) \\ &\quad p_t(\gamma_{t+1}|\gamma_t) p_t(z_{t+1}|z_t). \end{aligned} \quad (5)$$

It should be noted that b_{t+1} is obtained as a deterministic function of (s_t, a_t) , although the sequence b_t itself is stochastic. In fact, $p_t(b_{t+1}|a_t, s_t)$ is one if b_{t+1} is the SOC implied by action a_t and zero elsewhere. The load γ_t and the price z_t are assumed to be deterministic sequences. Thus, under these hypotheses the transition kernel is only influenced by the clearness. Compared to [14], although in our case decision epochs are at 15-minute intervals, we operate under similar assumptions: the market price is not random, the fixed load takes the place of committed wind power, the solar power takes the place of the actual wind power.

A control policy $\pi \in \Pi$ is made up of the decision functions $d_t : S \mapsto A$, $\pi = (d_1, \dots, d_{N-1})$, where S is the state space and A is the action space. A decision function d_t indicates, for every system state s_t , the value of the action a_t to be taken.

The overall value of the control policy $u_t^\pi(s_t)$ is the random variable representing the expected revenue obtained by adopting control policy π for a finite-horizon problem starting at epoch t with state s_t and ending at epoch N . The aim is to find $u_1^\pi(s_1)$. This is the so-called value-to-go, which represents the expected reward of the chosen control policy π , given that the initial state is s_1 . It can be derived by applying the backward recursive formula, for $t = N - 1, \dots, 1$

$$\begin{aligned} u_t^\pi(s_t) &= r_t(s_t, d_t(s_t)) + \\ &\quad + \mathbb{E}\{u_{t+1}^\pi(s_{t+1}) | s_t, d_t\} \\ &= r_t(s_t, d_t(s_t)) + \\ &\quad + \sum_{s_{t+1} \in S} p_t(s_{t+1}|s_t, d_t(s_t)) u_{t+1}^\pi(s_{t+1}), \\ u_N^\pi(s_N) &= r_N(s_N). \end{aligned} \quad (6)$$

The reward $r_N(s_N)$ is the value-to-go resulting from the sequence of actions taken during the dark hours.

B. Action Space

The energy stored in the battery at epoch t is $E_t = b_t E$. If control action a_t is taken, at epoch $t + 1$ the total amount of energy in the storage system is

$$E_{t+1} = E_t + a_t E = (b_t + a_t) E \quad (7)$$

Deterministic dynamic programming

```

1:  $\tilde{r}_M \leftarrow \text{const.}$ 
2:  $\tilde{u}_M^*(i) \leftarrow \tilde{r}_M, \forall i$ 
3: set PV production to 0  $\forall t$ 
4: for  $t = M - 1$  to 1 do
5:   update network data at epoch  $t$ 
6:   for  $i = 1$  to  $n_{\text{SOC}}$  states do
7:     for  $j = 1$  to  $n_{\text{SOC}}$  states do
8:       evaluate  $\Delta \tilde{E}_t^{ij}$ , i.e., the energy drawn from or given to the energy storage device during the transitions from  $i$ -th SOC at  $t$  to  $j$ -th SOC at time  $t + 1$ 
9:       evaluate  $\tilde{a}_t^{ij}$  in (7)
10:      chk = check feasibility of transition  $i \rightarrow j$ 
11:      if chk is false then
12:         $\tilde{r}_t^{ij} := \tilde{r}_t(i, \tilde{a}_t^{ij}) \leftarrow -\infty$ 
13:      else
14:        evaluate  $\tilde{r}_t^{ij} > -\infty$ 
15:      end if
16:    end for
17:    evaluate  $\tilde{u}_t^*(i) = \max_j \left( \tilde{u}_{t+1}^*(j) + \tilde{r}_t^{ij} \right)$ 
18:    evaluate  $(\tilde{a}_t^*)^*$  associated with  $\tilde{u}_t^*(i)$  and  $(\Delta \tilde{E}_t^*)^*$ 
19:  end for
20: end for

```

Fig. 1. Pseudo-code describing the deterministic dynamic programming algorithm implemented for the presented application

Thus, the amount of energy drawn from, or supplied to, the storage system is

$$a_t = \frac{\Delta E_t}{E}, \quad (8)$$

which can be positive or negative and is bounded by the physical and dynamical limits of the storage system. In general, these limits depend on the current state of charge, $a_t \in [a^{\min}(b_t), a^{\max}(b_t)]$.

C. Algorithm

The optimization algorithm is divided in two stages: a deterministic one, which runs during night-time hours, when there is no PV production, and a stochastic stage, which runs when there is PV production.

1) *Deterministic*: Recall that $M - 1$ is the total number of “night-time” epochs, i.e., those hours when radiation is not sufficient to produce a significant amount of energy. We represented this situation by setting the clearness to zero from 4 pm to 7:45 am the next day, thus obtaining 64 quarter-of-an-hour decision epoch. Then, starting from 8 am—which is the terminal time, indexed by M —the algorithm goes back to 4 pm. The pseudo-code is shown in Fig. 1, where the tilde is used to distinguish the revenue, the exchanged energy, the actions and the value-to-go of the deterministic part of the algorithm.

2) *Stochastic*: For $t = N$ the optimal value-to-go is evaluated for every possible state s_N

$$u_N^*(s_N) = r_N^*(s_N), \quad \forall s_N \in S. \quad (9)$$

For t from $N - 1$ to 1 the current epoch index is decreased by one unit and the optimal value of $u_t(s_t)$ is evaluated:

Stochastic dynamic programming

```

1: for  $t = N - 1$  to 1 do
2:   update network data at epoch  $t$ 
3:   for  $i = 1$  to  $n_{\text{SOC}}$  states do
4:     for  $k = 1$  to  $n_{\text{clr}}$  clearnesses do
5:       update PV production data for clearness  $k$ 
6:       for  $j = 1$  to  $n_{\text{SOC}}$  states do
7:         evaluate  $\Delta E_t^{ij}$ , i.e., the energy drawn from or given to the energy storage device during the transitions from  $i$ -th SOC at  $t$  to  $j$ -th SOC at time  $t + 1$ 
8:         evaluate  $a_t^{ij}$  in (7)
9:         chk = check feasibility of transition  $i \rightarrow j$ 
10:        if chk is false then
11:           $r_t^{ikj} := r_t((i, k), a_t^{ij}) \leftarrow -\infty$ 
12:        else
13:          evaluate  $r_t^{ikj} > -\infty$ 
14:        end if
15:         $w_k^{ij} \leftarrow 0$ 
16:        if  $t + 1 = N$  then
17:           $h = 1$  {at epoch  $N$  night-time hours period starts}
18:           $u_{t+1}^*(j, h) \leftarrow \tilde{u}_1^*(j)$ 
19:           $w_k^{ij} \leftarrow p(h|k) u_{t+1}^*(j, h)$  { $p(h|k) = 1, \forall k$ }
20:        else
21:          for  $h = 1$  to  $n_{\text{clr}}$  clearnesses do
22:             $w_k^{ij} \leftarrow w_k^{ij} + p(h|k) u_{t+1}^*(j, h)$ 
23:          end for
24:        end if
25:        end for
26:        evaluate  $u_t^*(i, k) = \max_j \left( r_t^{ikj} + w_k^{ij} \right)$ 
27:        evaluate  $(a_t^*)^*$  associated with  $u_t^*(i, k)$  and  $(\Delta E_t^*)^*$ 
28:      end for
29:    end for
30:  end for

```

Fig. 2. Pseudo-code describing the stochastic dynamic programming algorithm implemented for the presented application.

$$u_t^*(s_t) = \max_{a \in A_{s_t}} \left\{ r_t(s_t, a) + \sum_{s_{t+1} \in S} p_t(s_{t+1}|s_t, a) u_{t+1}^*(s_{t+1}) \right\}. \quad (10)$$

In general, there can be more than one value of a maximizing (10): these values are collected in the set $A_{s_t, t}^*$. When $t = 1$ the value $u_1^*(s_1)$ —i.e., the value-to-go associated with an optimal control policy π^* —is obtained. In this case π^* is a set of sequences of decision functions such as $(d_1^*, \dots, d_{N-1}^*)$, where any d_t^* selects a particular action from $A_{s_t, t}^*$.

The pseudo-code of the stochastic part of the algorithm is shown in Fig. 2.

III. CASE STUDY

The test network used for the simulations is the European medium voltage (MV) benchmark network of the CIGRÉ [17], the single line diagram of which is depicted in Fig. 3. The network is made up of two feeders, framed by dashed lines in Fig. 3, both operating at 20 kV and fed by the high voltage grid by means of two separate 25 MVA-rated transformers. The network is made up of 15 lines and 15 nodes and feeds 14 loads, both residential and commercial. It is worth noting that the feeders can be connected by closing switch S1, thus transforming the radial distribution configuration in a meshed one. In the present simulations all switches (S1, S2, and S3) are

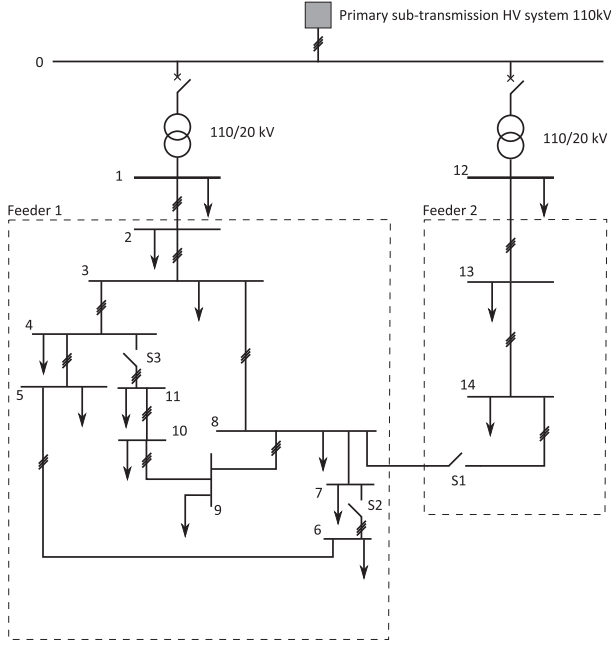


Fig. 3. Single line diagram of the MV feeders used for the simulation.

supposed to be closed. This configuration would guarantee, by itself, a more balanced sharing of the power flows in the grid and is one of the few structural countermeasures that a distribution system operator can exploit to avoid new installations when coping with the growing impact of renewable energy sources. For more detailed information concerning network parameters and configuration refer to [17]. The network has been implemented using MATPOWER [18]. The power flow algorithm provided by this toolbox has been used.

The price used in the simulations is depicted in Fig. 4 and is derived from the so-called PUN, the hourly-based Italian Power Exchange clearing price as it comes from the day-ahead market closure, in 2010. As reported in [3], the price profile has been built using the median values of energy price for each hour and can be considered a paradigm for energy price.

The clearness has been discretized to $n_{chr} = 14$ classes. The transition matrix for the clearness is reported in (11), shown at the bottom of the page, where the generic entry (i, j) is $p_t(j|i)$.

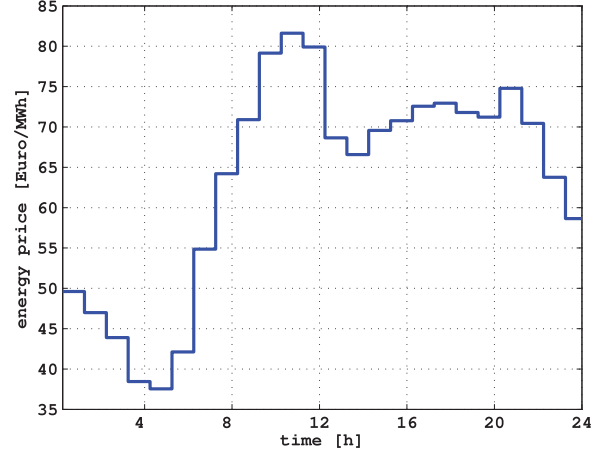


Fig. 4. Mean energy price used in the simulation and derived from Italian Power Exchange (IPEX) in 2010 [3].

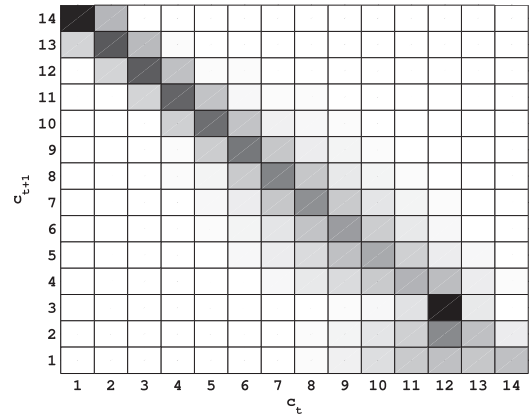


Fig. 5. Graphical representation of the transition matrix for the clearness, as reported in (11).

It is worth noting that each row of (11) sums up to 1. A graphical representation of the transition matrix is reported in Fig. 5, where white color is associated to 0 and black to 0.749, the maximum value of the transition matrix (11). The stationary distribution associated to the transition matrix for the clearness is plotted in Fig. 6.

$$\begin{pmatrix}
 72.8 & 25.1 & 1.10 & 0.4 & 0.3 & 0.2 & 0.1 & 0.0 & 0.0 & 0.0 & 0.0 & 0.0 & 0.0 & 0.0 \\
 14.1 & 59.3 & 23.5 & 1.40 & 0.7 & 0.5 & 0.3 & 0.1 & 0.0 & 0.0 & 0.0 & 0.0 & 0.0 & 0.0 \\
 0.4 & 14.6 & 58.1 & 22.1 & 2.0 & 1.3 & 0.9 & 0.4 & 0.1 & 0.0 & 0.0 & 0.0 & 0.0 & 0.0 \\
 0.1 & 0.6 & 15.2 & 55.9 & 20.9 & 3.2 & 2.1 & 1.30 & 0.5 & 0.2 & 0.0 & 0.0 & 0.0 & 0.0 \\
 0.0 & 0.2 & 1.0 & 15.8 & 52.6 & 20.1 & 4.7 & 3.10 & 1.6 & 0.6 & 0.2 & 0.0 & 0.0 & 0.0 \\
 0.0 & 0.1 & 0.5 & 1.90 & 16.6 & 48.3 & 19.7 & 6.40 & 3.9 & 1.8 & 0.6 & 0.1 & 0.0 & 0.0 \\
 0.0 & 0.1 & 0.3 & 1.20 & 3.6 & 17.6 & 43.5 & 19.4 & 7.8 & 4.2 & 1.7 & 0.5 & 0.1 & 0.0 \\
 0.0 & 0.0 & 0.2 & 0.7 & 2.4 & 6.0 & 19.1 & 38.9 & 18.7 & 8.2 & 4.0 & 1.4 & 0.4 & 0.1 \\
 0.0 & 0.0 & 0.1 & 0.3 & 1.4 & 4.1 & 9.1 & 20.8 & 34.5 & 17.3 & 7.7 & 3.3 & 1.0 & 0.3 \\
 0.0 & 0.0 & 0.0 & 0.1 & 0.7 & 2.5 & 6.5 & 12.4 & 22.2 & 30.3 & 15.3 & 6.5 & 2.5 & 0.9 \\
 0.0 & 0.0 & 0.0 & 0.0 & 0.2 & 1.0 & 3.3 & 7.6 & 12.7 & 18.7 & 26.3 & 22.5 & 5.9 & 1.8 \\
 0.0 & 0.0 & 0.0 & 0.0 & 0.0 & 0.1 & 0.3 & 0.9 & 1.9 & 2.9 & 10.2 & 74.9 & 7.8 & 0.8 \\
 0.0 & 0.0 & 0.0 & 0.0 & 0.0 & 0.1 & 0.5 & 1.90 & 4.6 & 8.2 & 15.6 & 41.8 & 21.4 & 5.9 \\
 0.0 & 0.0 & 0.0 & 0.0 & 0.0 & 0.1 & 0.4 & 1.80 & 5.3 & 11.5 & 18.2 & 21.6 & 19.7 & 21.5
 \end{pmatrix} \times 10^{-2} \quad (11)$$

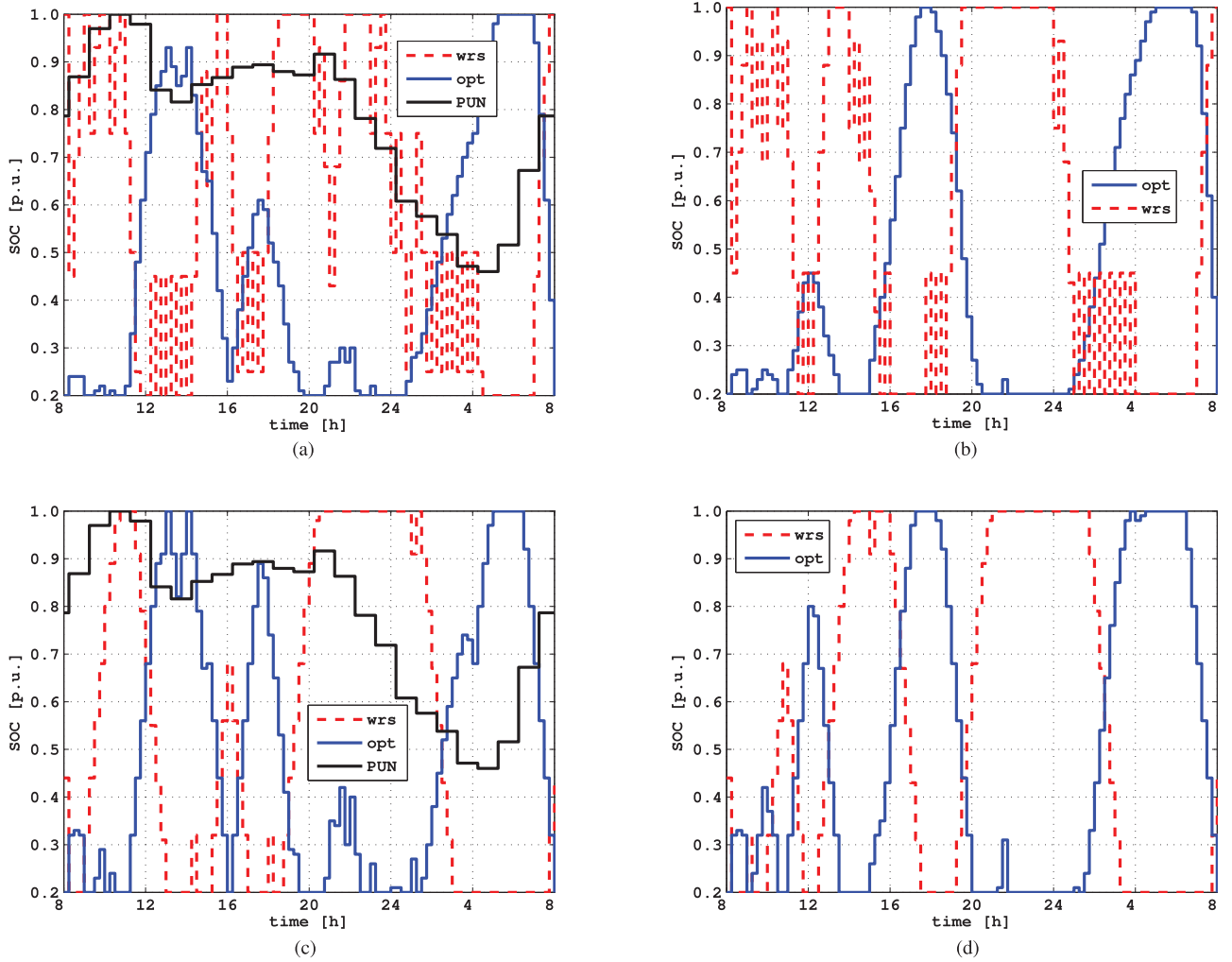


Fig. 7. SOC paths for one day of the simulations. In each plot the optimal policy (opt) is drawn with a solid blue line while the worst policy (wrs) is drawn with a dashed red line. The solid black line in the left-hand side plots is the normalized PUN, superimposed in order to highlight the influence of energy price on the policies. Top row: ESSS no. 6 = (1 MWh, 1 MW); (a) revenue based on energy price, (b) revenue based on losses minimization. Bottom row: ESSS no. 1 = (250 kWh, 125 kW); (c) revenue based on energy price, (d) revenue based on losses minimization.

$$\bar{y}^{\text{algo}} = \frac{1}{5000} \sum_{k=1}^{5000} \left(\frac{(y_k^{\text{algo}} - y_k^{\text{nostg}}) \times 100}{y_k^{\text{nostg}}} \right),$$

$\text{algo} \in \{\text{opt}, \text{rnd}, \text{wrs}\}.$

Maximum values of the standard deviations for each policy are reported in the captions of the tables. From the analysis of the results, it emerged that, in this particular case, the differences between best and worst policies in terms of revenue—although the values of the revenue reached by optimal policy are strictly higher than those obtained by the worst policy—are not particularly striking. This fact is due to some concurrent causes: 1) the overlap between load and PV production, 2) the small size of the energy storage device with respect to the overall consumption, and 3) the fact that energy price is high during daylight hours, when PV production is high. Moreover, the meshed grid, by itself, reduces energy losses.

The most significant results are obtained when using a losses reduction revenue function. The beneficial effect coming from a

correct usage of the storage system increase when going farther from the HV/MV transformers and when using a high-capacity and high-power storage system. Losses reduction comes to a maximum figure of 1.62%. This value can be regarded as significant, as reported in [19, footnote no. 48, p. 27]. It is verified that the random policy always lays in between the optimal and the worst policy, thus confirming that every other possible management strategy would lead to no better results than the optimal one. It can be also noticed that only the optimal policy produces a reduction, while random and worst policies do not only produce no beneficial effect but can even increase both losses and costs.

Despite this little difference in the revenue, the optimal and the worst policy are extremely different. In fact, the latter tends to be the complement of the former. In order to highlight these differences, SOC paths during the same day, but under different hypotheses, are shown in Fig. 7. The day chosen is the no. 832, which corresponds to the day with the lowest number of SOC cycle of the worst algorithm for ESSS no. 6. The clearness associated to this day is reported in Fig. 8.

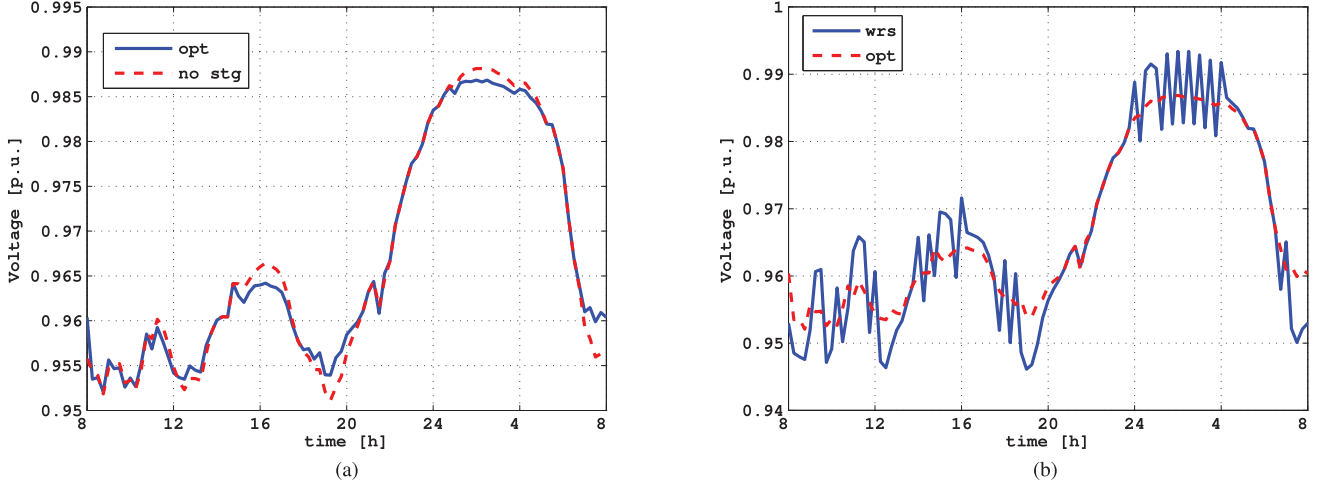


Fig. 9. Voltage profiles at bus no. 6 for day no. 3683 and ESSS no. 6 with storage system connected at bus 5. The optimization criteria is losses minimization. (a) comparison between the optimal policy (solid blue line) and the base case with no storage (dashed red line); (b) comparison between the optimal policy (dashed red line) and the worst policy (solid blue line).

ESSS to the next one, the difference between worst and optimal policy is such that

$$(x_{i+2}^{\text{wrs}} - x_{i+2}^{\text{opt}}) - (x_i^{\text{wrs}} - x_i^{\text{opt}}) > 0. \quad (14)$$

This means that there is some sort of “regularity” in the optimal policy. It can also be noted that in the set of ESSSs (1–3–5), $x_1^{\text{wrs}} - x_1^{\text{opt}}$ is, generally, negative. Thus, it can be inferred that the combination of low capacity and low power is not suited for this kind of application. The only exceptions are for node 12 in Tab. V and for nodes 12 and 13 in Tab. IV and derive from network configuration and PV plant location.

Finally, a comparison between voltage profiles is shown in Fig. 9. The voltage profiles of bus no. 6 for day no. 3683 and ESSS no. 6 with storage system connected at bus 5 (optimization only on losses). This particular ESSS, along with the placement of the energy storage system has been chosen because it gives the best result. The particular day is when the maximum spread between highest and lowest voltage occurs, while the bus selected is that which displays the minimum voltage. This selection has been made because of the huge amount of data (more than 37 million of voltage values for each optimization criterion). It can be noticed that the optimal management of the energy storage system enhance the voltage profiles.

V. CONCLUSION

The application of the MDP framework to the optimal scheduling of a storage device in the presence of solar generation is a tool which can contribute effectively to network management. The superiority of this method is intrinsic, because it takes into account the whole chain of consequences of every possible action across the entire decisional horizon, instead of considering only the here-and-now reward at every epoch. The Markovian assumption for the clearness stochastic process offers a simplified yet accurate enough representation of the uncertainty of the solar source, so that the optimal strategy is valid (on average) over all possible instances of the clearness over a day and it can be applied to all days with a similar law

of the clearness pattern. Then, by estimating different transition matrices for every season in the year from long enough series of observed radiation data, a year-long optimal scheduling strategy for the storage can be obtained.

The specific testing environment chosen for this article has demonstrated that, although the improvement measured by the selected objective function can be modest (be it in terms of network losses or energy prices), this methodology can optimize other features of the system, such as the number of cycles of the storage, so that its life will be increased. Incidentally, the direct usage of the number of cycles as objective function within this framework would not be possible, because the optimization algorithm is based on a chain of evaluations of reward functions spanning only one epoch at a time. Then, either relying directly on the main objective function or on performance measures arising as by-products of the optimization, this methodology could be used to assess the effectiveness of a range of system layouts, such as the placement of the storage, its size and its power rating.

Further work should consider different systems, such as networks equipped with a larger solar field or multiple storages. A problem that could arise is that, after executing the policy optimization for multiple layouts, it might not be easy to discriminate among them, especially if the associated rewards are close to each other. Therefore, appropriate choice criteria of the reward function or of by-product performance measures should be studied.

ACKNOWLEDGMENT

The authors would like to thank the anonymous reviewers for their careful reading of the manuscript and for suggesting several modifications which greatly helped to improve the paper.

REFERENCES

- [1] M. Farrokhifar, S. Grillo, and E. Tironi, “Loss minimization in medium voltage distribution grids by optimal management of energy storage devices,” in *Proc. IEEE PowerTech*, 2013, pp. 1–5.

- [2] M. Farrokhifar, S. Grillo, and E. Tironi, "Optimal placement of energy storage devices for loss reduction in distribution networks," in *Proc. IEEE Innov. Smart Grid Technol. Eur.*, 2013, pp. 1–5.
- [3] S. Grillo, M. Marinelli, S. Massucco, and F. Silvestro, "Optimal management strategy of a battery-based storage system to improve renewable energy integration in distribution networks," *IEEE Trans. Smart Grid*, vol. 3, no. 2, pp. 950–958, Jun. 2012.
- [4] V. Marano, G. Rizzo, and F. A. Tiano, "Application of dynamic programming to the optimal management of a hybrid power plant with wind turbines, photovoltaic panels and compressed air energy storage," *Appl. Energy*, vol. 97, pp. 849–859, Sep. 2012.
- [5] H. Morais, P. Kádár, P. Faria, Z. A. Vale, and H. Khodr, "Optimal scheduling of a renewable micro-grid in an isolated load area using mixed-integer linear programming," *Renew. Energy*, vol. 35, no. 1, pp. 151–156, Jan. 2010.
- [6] Z. Ziadi, S. Taira, M. Oshiro, and T. Funabashi, "Optimal power scheduling for smart grids considering controllable loads and high penetration of photovoltaic generation," *IEEE Trans. Smart Grid*, vol. 5, no. 5, pp. 2350–2359, Sep. 2014.
- [7] M. L. Puterman, *Markov Decision Processes*. Hoboken, NJ, USA: Wiley, 1994.
- [8] C. A. Glasbey, "Nonlinear autoregressive time series with multivariate Gaussian mixtures as marginal distributions," *Appl. Stat.*, vol. 50, pp. 143–154, 2001.
- [9] M. A. Murtaza and M. Tahir, "Optimal data transmission and battery charging policies for solar powered sensor networks using Markov decision processes," in *Proc. IEEE Wireless Commun. Netw. Conf.*, 2013, pp. 992–997.
- [10] J. Donadee and M. Ilić, "Stochastic optimization of grid to vehicle frequency regulation capacity bids," *IEEE Trans. Smart Grid*, vol. 5, no. 2, pp. 1061–1069, Mar. 2014.
- [11] P. M. van de Ven, N. Hegde, L. Massoulié, and T. Salonidis, "Optimal control of end-user energy storage," *IEEE Trans. Smart Grid*, vol. 4, no. 2, pp. 789–797, Jun. 2013.
- [12] J. Qin, R. Sevlian, D. Varodayan, and J. Rajagopal, "Optimal electric energy storage operation," in *Proc. IEEE Power Energy Soc. Gen. Meeting*, 2012, pp. 1–6.
- [13] N. Gast, D.-C. Tomozei, and J.-Y. Le Boudec, "Optimal generation and storage scheduling in the presence of renewable forecast uncertainties," *IEEE Trans. Smart Grid*, vol. 5, no. 3, pp. 1328–1339, May 2014.
- [14] L. A. Hannah and D. B. Dunson, "Approximate dynamic programming for storage problems," in *Proc. 28th Int. Conf. Mach. Learn.*, 2011, pp. 337–344.
- [15] V. Musolino, L. Piegari, and E. Tironi, "Technical and economical evaluation of storage systems in naval applications," in *Proc. Int. Conf. Clean Electr. Power (ICCEP)*, 2013, pp. 120–127.
- [16] G. Celli, S. Mocci, F. Pilo, and M. Loddo, "Optimal integration of energy storage in distribution networks," in *Proc. IEEE PowerTech Conf.*, 2009, pp. 1–7.
- [17] K. Strunz *et al.*, "Benchmark systems for network integration of renewable and distributed energy resources," CIGRÉ TF C6.04.02, Tech. Rep. 575, Apr. 2014.
- [18] R. D. Zimmerman, C. E. Murillo-Sánchez, and R. J. Thomas, "MATPOWER: Steady-state operations, planning, and analysis tools for power systems research and education," *IEEE Trans. Power Syst.*, vol. 26, no. 1, pp. 12–19, Feb. 2014.
- [19] VV. AA., "Smart Distribution System: Promozione selettiva degli investimenti nei sistemi innovativi di distribuzione di energia elettrica. Orientamenti iniziali. (Smart Distribution System: Selective promotion of investments in innovative electric distribution systems. Initial guidelines.);" Autorità per l'Energia Elettrica, il Gas e il Sistema Idrico (AEEGSI), Tech. Rep. 255/2015/R/eel, May 2015, in Italian [Online]. Available: <http://www.autorita.energia.it/it/docs/dc/15/255-15.jsp>



Samuele Grillo (S'05–M'09) received the Laurea degree in electronic engineering, and the Ph.D. degree in power systems from the University of Genova, Italy, in 2004 and 2008, respectively. He is currently an Assistant Professor (tenure track) with the Dipartimento di Elettronica, Informazione e Bioingegneria, Politecnico di Milano, Milan, Italy. His research interests include smart grids, integration of distributed renewable sources, and energy storage devices in power networks, optimization, and control techniques applied to power systems.



Antonio Piegari received the Ph.D. in statistics from Padua University, Padua, Italy. He is currently a Researcher with CNR-IMATI. His research interests include industrial statistics. In particular, he likes to address problems in engineering and technology using stochastic modeling tools, with a preference for Markov and point processes. His main contributions can be found in the application of stochastic modeling to reliability and to electrical power systems.



Enrico Tironi received the M.S. degree in electrical engineering from the Politecnico di Milano, Milan, Italy, in 1972. In 1972, he joined the Department of Electrical Engineering, Politecnico di Milano, where he is currently a Full Professor with the Dipartimento di Elettronica, Informazione e Bioingegneria. His research interests include power electronics, power quality, distributed generation, and energy storage systems. He is a member of the Italian Standard Authority (C.E.I.), the Italian Electrical Association, and the Italian National Research Council Group of

Electrical Power System.