



Deep Reinforcement Learning for Long Term Hydropower Production Scheduling

Signe Riemer-Sørensen 
 Mathematics and Cybernetics
 SINTEF Digital, Oslo, Norway
 signe.riemer-sorensen@sintef.no

Gjert H. Rosenlund 
 Energy systems
 SINTEF Energy Research, Trondheim, Norway
 gjert.rosenlund@sintef.no

Abstract—We explore the use of deep reinforcement learning to provide strategies for long term scheduling of hydropower production. We consider a use-case where the aim is to optimise the yearly revenue given week-by-week inflows to the reservoir and electricity prices. The challenge is to decide between immediate water release at the spot price of electricity and storing the water for later power production at an unknown price, given constraints on the system. We successfully train a soft actor-critic algorithm on a simplified scenario with historical data from the Nordic power market. The presented model is not ready to substitute traditional optimisation tools but demonstrates the complementary potential of reinforcement learning in the data-rich field of hydropower scheduling.

Index Terms—machine learning, expert systems, power generation economics, hydroelectric power generation

NOMENCLATURE

a_i	Action, amount of water to convert to electricity as percentage of maximum production capacity.
i	Net inflow, normalised with r_{\max} .
f_{\max}	Maximum production relative to reservoir capacity.
k_{price}	Importance of price in the reward function.
q_{price}	Power of price in the reward function.
R_i	Reward for action a_i , arbitrary units.
r_{\max}	Volume of reservoir (Mm^3).
s	State, characterised by week number, reservoir level, weekly inflow and weekly price.
y_i	Weekly price normalised to maximum (Mm^{-3}).

I. INTRODUCTION

A. Motivation and background

Hydroelectricity is sustainable energy, but due to seasonal variations, the natural inflow to the reservoirs does not follow demand or price as shown in Figure 1. In order to maximise revenue and minimise water spillage, the power generator schedule production based on a combination of deterministic and stochastic models of the inflow and demand [1], [2]. In traditional methods, the size and complexity of the computations requires human intervention to decompose the problem into smaller parts. By utilising advancements in reinforcement learning [3], we take a step towards automating the process and solving the problem without human intervention.

In the concept of reinforcement learning a reward scheme is used to train an agent (the algorithm) to a desired behaviour

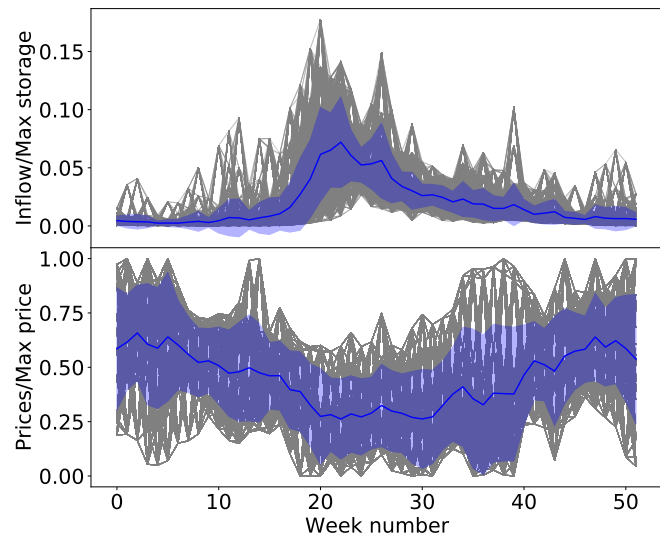


Fig. 1. 7975 samples of historic inflows and prices as a function of week number for four similar water systems in Southern Norway: Tora, Storeskar, Sula, Rinna. The solid blue line indicates the average and the shaded region is the standard deviation. Inflow is highest during spring and early summer, while prices are on average highest during the winter.

(action) in a range of situations (states) within a given system (the environment). Reinforcement learning is most efficient when the environment and its response to actions taken by the agent (rewards) are governed by some underlying patterns that are at least partly deterministic. If there are no relations between actions and rewards, there is nothing to learn (e.g. chapter 3.1 in [3]). Thus, the goal is for the algorithm to generalise over stochastic variations.

In the case of hydropower scheduling, the options for the actions are well defined (release water or save), while the feedback from the environment is a mixture of deterministic and stochastic: Prices are related to production and follow a yearly pattern but with fluctuations, inflows are weather dependent (stochastic) with general seasonal variations.

Recent years have seen a massive development within the combination of deep learning and reinforcement learning, enabling algorithms to solve complex tasks [4]–[6]. They are highly flexible and powerful methods that potentially can be adapted to many different purposes but given their

complexity, very few applications have been demonstrated outside simulated or strictly rule based systems [7].

B. Contributions and organisation

Our contribution is to modify and train a state-of-the art soft actor-critic algorithm on a toy scheduling scenario and two realistic scenarios, to demonstrate the potential of reinforcement learning in hydropower scheduling. As will be explained in Section II, the soft actor-critic algorithm is particularly well suited to systems with a high degree of stochastic fluctuations. The agent is trained in a “safe” environment made up from stochastic variations of historical data (Section IV), from which it is able to successfully learn a meaningful policy for water release. The algorithm is not (yet) a replacement for traditional optimisation and planning methods, but rather we discuss its potential as a complementary method in Section V.

II. SOFT ACTOR-CRITIC REINFORCEMENT LEARNING

Soft actor-critic reinforcement learning is based on Q-learning, where all states and actions are assigned a value [8], [9]. The aim is to learn an approximation for the value of all possible actions in each state, and for each step choose the action with the potential to provide the largest total reward. In its pure form, Q-learning is limited to discrete states and actions, it is not very flexible and scalable, and not very robust to stochastic fluctuations and time-dependent evolution [3].

In actor-critic algorithms, the Q-learning concept is expanded and the policy and value functions are separated. The value function is used in the learning process (the critic), but the actor selects the action without consulting the value function. In the soft actor-critic algorithm, the lifetime reward optimisation is combined with an entropy maximisation [10] leading to substantial improvement in terms of learning speed, final performance, sample efficiency, stability, and scalability, in particular on complex tasks or on stochastic systems [11].

III. IMPLEMENTATION

The algorithm is implemented with a policy network (the action to take for a given state), a value network (describing how advantageous each state is), a target value network, and two soft Q-networks [10] as illustrated in Figure 2.

The target value network is introduced because the target for training the Q-network depends on the value network, and the target for the value network depends on the Q-network. Consequently, this loop makes the training unstable. The solution is to use an additional network while training the Q-network, which is close to the value network but with a short time delay. Thus, the target value network becomes a memory of how the value network was a moment ago. Instead of copying the value network directly, we perform a Polyak averaging between the target value network and the value network to obtain a kind of moving average over the gained experience [12].

The theory only calls for one Q-function, but in practice, there is a tendency for the algorithm to overestimate the Q-values. This is mitigated by training two Q-networks and

always using the minimum of the two values when updating the policy and value function networks. Neither is it strictly necessary to have separate approximators for the value and Q-functions, since they are related through the policy, but in practice it helps with convergence. [11] discusses how the value function can be disposed of and introduces automatic detection of the weight of the entropy term (called the temperature).

We apply memory replay where the learning experiences are stored and reused when training the neural networks, in order to obtain a better balance between rare and common events.

We follow the pytorch implementation of [13], [14] but with significant modifications related to the activation functions and determination of accessible states.

A. The environment

While the algorithm itself is generic, the environment and reward functions are specific to the problem we want to solve. In its simplest form, the purpose of hydro power scheduling is to maximise the income, given scenarios on inflow. This becomes a balancing act between storing the water for winter (high prices), while not keeping too much, leading to spillage during spring (high inflow, snow melting).

The reservoir holds a maximum available volume for production of r_{\max} , and i accounts for the net inflow. If the capacity of the reservoir is exceeded, the water will overflow as spillage. All water volume quantities are given in Mm^3 , and in the algorithm they are normalised to maximum available water volume so they become unit-less.

Each state is defined by week number, storage, price for the week, inflow for the week, and the number of weeks it would take to empty the reservoir if running on full capacity. The latter is used to clamp the policy function to feasible actions only. The initial storage in the reservoir is randomly sampled from a specified distribution.

We consider the simplest possible viable version of the problem. The method scales well and further constraints can be implemented with relatively low effort (see Section V).

The training data (described in Section IV-A) are also loaded into the environment. The inflow is normalised to maximum reservoir volume, and the prices are normalised to the maximum price. In addition, the water remaining in the reservoir at the end of the year can be given a value e.g. the value from the last week, or the average across all scenarios.

B. The action space

The action space is defined as the amount of water to be released and converted to electricity at the spot price. Normalising with the maximum production capacity, this becomes a continuous variable with values between 0 and 1. Since the volume never can be a negative number, we have implemented a sigmoid activation function and scaling in the policy network.

C. The step function and reward

The step function determines the response by the environment of a given action performed in a given state. This includes

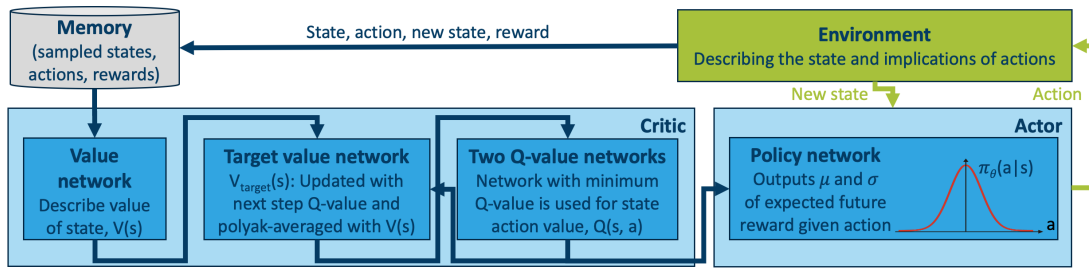


Fig. 2. The components of the soft actor-critic algorithm and their interplay.

computing the reward and the transition to the next state (update reservoir level and informing of new inflow and price values). In order to make the algorithm as transparent as possible, we have maintained a fairly simple reward determination: First the suggested production and corresponding end storage is computed. Then we check the feasibility of the action. If the end storage is larger than the capacity, the reservoir is flooding and the end storage is adjusted to maximum capacity.

The reward R_i for action a_i is given by:

$$R_i = a_i f_{\max} r_{\max} (y_i k_{\text{price}})^{q_{\text{price}}}, \quad (1)$$

where f_{\max} (unit-less) is the maximum production capacity (maximum volume of water that can be converted to electricity) relative to the total reservoir capacity. r_{\max} is the reservoir capacity (Mm^3). y_i is the weekly price normalised to maximum price (Mm^{-3}) and the importance of the price is controlled with k_{price} and q_{price} . k_{price} can also be used to relate the produced volume to actual prices and account for a non-linear production function. The reward is a unit-less number unless otherwise adjusted with k_{price} .

D. The networks

The value networks (both main and target), and the soft Q-networks are standard fully connected neural networks each with three hidden layers with relu activation functions [15] and an output layer with a linear activation function. The number of neurons in all hidden layers is a hyperparameter to be set by the user.

The policy network has two hidden layers with relu activation functions. It has two outputs, the mean and the logarithm of the standard deviation (clamped to be in a sane region). These are used for a re-parametrisation that ensures that the sampling from the policy is differential and the errors can be properly back propagated, leading to faster convergence. The action to be taken from a given state, s , is obtained from the policy function by sampling noise, from a standard normal distribution, multiply it with the standard deviation, add it to the mean, and then activate it with a sigmoid activation function to ensure an action between 0 and 1 [10].

The full update process then becomes:

- 1) Initialisation of the networks.
- 2) Initialisation of the environment.
- 3) Training loop:

- At the beginning of each episode (full year), reset the environment and sample the initial amount of water in storage.
- For each week, randomly sample the price and inflow, and obtain the action from the policy (or randomly in the early exploration phase).
- Determine the reward and the next state (step function in environment).
- Save experience to memory.
- Update networks in batches from memory:
 - Predict values of Q-functions, value and policy network (for all states in the batch and their suggested actions).
 - Evaluate the policy network to get the next states.
 - Predict the target value network.
 - Compute loss of Q-functions and value network and do one back propagation (update weights).
 - Adjust the target value function with next step Q-value.
 - Compute the loss of the target value function and do one back propagation.
 - Compute loss of policy network and do one back propagation.
 - Update target value network by Polyak averaging with the main value network.
- Save the model.

E. Hyperparameters

The network structures (number of layers) are hardcoded and not considered hyperparameters here. However, the following parameters must be chosen (the values in parentheses indicate the values applied in Section IV):

- Number of neurons in the hidden layers (100).
- Absolute values for the outer edges of the range of uniform initialisation weights ($3 \times 10^{-3} \in [0, \infty]$).
- Loss function for the value and soft Q networks (mean squared error).
- Optimisation functions for all networks (RMSprop).
- Learning rates of the value network, the soft Q networks and the policy networks (5×10^{-4} , 5×10^{-4} , 1×10^{-4}).
- Clamp values for the logarithmic standard deviation on the policy network. The results are insensitive to the exact values (-20, 2).

- Minimal logarithmic probability. The results are insensitive to the exact value (3×10^{-6}).
- Discount factor, γ ($0.99 \in [0, 1]$).
- Soft τ , determines the importance of the main value network when updating the target value network (0.0006).
- Number of exploration steps before the training begins (10000 in the artificial case and 50000 in the historic).
- Total number of weeks to train (300000, model may converge before).
- Number of experiences to save in the replay memory (all).
- Batch size to use from replay memory (100). The number of randomly selected experiences to use for every update of the networks.

We use RMSprop [16] as optimiser because its lower dependence on momentum than e.g. adam [17], makes it more adaptable and well suited to handling the non-stationary data distribution from the changing environment. In addition, the policy network has a smaller learning rate than the value functions, in order to collect experience at a faster rate than the policy adapt to the experience. While slowing down training, this prevents excessive exploration of local minima.

The hyperparameters were decided based on default values from [10], [11] and manual tuning. Unfortunately, there are no efficient procedures for obtaining the optimal set of hyperparameters in reinforcement learning. However, several of the parameters will mainly influence convergence rate, so if trained for a large enough number of episodes, the final model performance will be similar.

The environment also contains some system definitions and hyperparameters:

- Maximum capacity of the reservoir, volume (1000).
- Maximum production, volume per week (30 and 100).
- Scaling factor for the price, k_{price} (1). Numeric factor that allows for tuning the sizes of the rewards and e.g. relate them to actual prices.
- Price power, q_{price} (1). The importance of the price in the reward function. Can also be used to account for non-linear conversion.
- Initial storage (randomly sampled between 0.4 and 0.6 of full storage capacity).

IV. USE CASES

Our main purpose is to validate reinforcement learning as a method for solving the problem of seasonal hydropower scheduling. To make the implementation and tuning of hyperparameters transparent, we have designed use cases without the most complicating elements; the production of electricity is assumed to be linear to the amount of water released (but the reward function allows for a power law relationship), the station is not part of a cascaded or otherwise restricted system, and there are no pumps or hatches to operate. Consequently, they serve as a minimal viable demonstration of the algorithm.

A common measure of flexibility in a hydropower plant is the usage time defined as how long the plant would operate at maximum capacity to convert a full reservoir to electricity. The use cases are designed to reflect common usage times for

Norwegian power stations, while being sufficiently diverse for the algorithm to learn different strategies.

The first power station has a relation between maximum production capacity and total reservoir volume of 30/1000 corresponding to a usage time of 5600 hours. The average yearly inflow is taken to be the same as the reservoir capacity. The second power station is used to demonstrate the changes when both the inflow and production capacity are higher with a ratio of 100/1000 corresponding to 1680 hours of usage time, and yearly inflow of 4000. For both reservoirs, the initial storage is randomly sampled between 40% and 60% of the full storage capacity.

A. Training data

We train and apply the agent on two sets of historic data. Firstly, we train on artificial scenarios, where there is a clear structure in the price and inflow, albeit with some fluctuations. In the second case, we use historic price data from the Nordic region (NO3) for the period¹ 2008 to 2019 and inflow provided by The Norwegian Water Resources and Energy Directorate for the period² 1958 to 2018. In order to expand the inflow data without introducing additional noise, we combine inflows from four reservoirs in the NO3 region with similar weather patterns: Tora, Storeskar, Sula and Rinna (shown in Figure 1). The weekly prices and inflows are linearly scaled to the yearly minimum and maximum values. We construct each scenario by randomly sampling 52 pairs of price and inflow values from the corresponding samples for the particular weeks. Figure 4 shows examples of the scenarios in the upper two panels. In the artificial scenarios, the water remaining in the reservoir at the end of the year is valued at its price for the last week and in the historic scenarios it fetches maximum price. The agent is only rewarded for remaining water at the end of the year if the reservoir filling is between 0.4 and 0.6. This is to encourage some storage for the next year. In an operational setting, the best treatment of the end-value condition has to be further investigated.

B. Results

Training the algorithm on the use cases described in Section IV for 300000 episodes takes a day on a 3.1GHz processor with 16Gb RAM available. This training only needs to be done once for each reservoir/system. Applying the trained model requires less than a minute to provide a plan for 52 weeks. After deployment, the model can be updated continuously with experience gathered from new data (timescale of minutes). For more complex systems (e.g. cascading reservoirs), the initial training is expected to take longer. Figure 3 shows an example of the total rewards as a function of training episode. The transition from random exploration to policy driven actions leads to an overall increase in rewards but with continued exploration.

¹provided by the NordPool Group nordpoolgroup.com

²nve.no/hydrologi/hydrologiske-data/historiske-data/historiske-vannforingsdata-til-produksjonsplanlegging

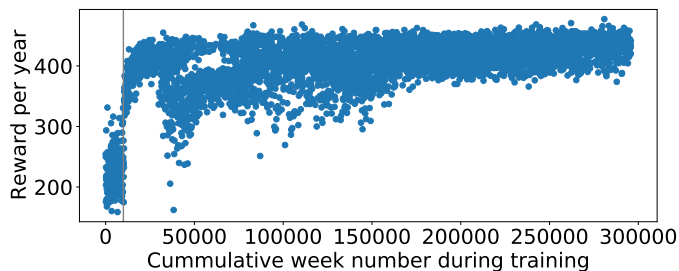


Fig. 3. Example of total reward as a function of training episode. The vertical line indicates the transition from random exploration to policy driven actions.

Figure 4 shows the resulting actions for five example scenarios for each fictive reservoir. In all cases, the agent was trained for 3000000 episodes corresponding to 5700 years of training data. As seen in Figure 3, the model converged earlier and the training could have been stopped. The agent has successfully learned a meaningful policy for water release given the limited state information.

When the price and inflow data are dominated by fluctuations, the resulting actions will also have a large range of fluctuations. This is partly due to the lack of generalisability for the value function and partly due to the entropy optimisation, where the soft actor-critic agent will automatically optimise the entropy and increase the random exploration for regions of the parameter space with large fluctuations. This behaviour is part of the reason why soft actor-critic is ideal for stochastic dominated environments, because it protects against the actor learning to always perform the same action and fail miserably when conditions change. In Figure 4 we show the actions from the trained agent without the additional stochasticity, in order to see how the policy exploits what it has learned. This approach is known to lead to better performance than including the stochasticity at runtime [18].

In the artificial case, the agent has learned to release most of the water when there is a high revenue albeit with some variation to protect against stochastic variation.

In the historic case, the agent releases water in the beginning of the year when the price is high, but less water in second half of the year in order to obtain the reward for having water in storage at the end of the year. For the system with production capacity of 30/1000, the storage is not fully exploited as there is no incitement in the reward function or environment for storing more than 60% of the capacity. For the system with higher inflow and larger production capacity, the additional capacity becomes a buffer to hold the inflow.

V. DISCUSSION AND OUTLOOK

We have demonstrated that the agent is able to learn a policy that is sensible for human interpretation. However, this reflects a minimal viable case and at the current level, the policy is not ready for real-life deployment. The resulting production plan is somewhat different from the suggestions of traditional stochastic optimisation tools, indicating that the reinforcement learning may be able to exploit the data differently than

classical models. Consequently, we do not advocate that the classical models can be replaced by reinforcement learning, but rather that they can be complementary.

Reinforcement learning has previously been attempted used in hydropower scheduling [19], [20] and in combining multiple connected reservoirs [21]. Relative to [20], [21] which both apply Q-learning and consequently discrete state/action spaces, soft actor-critic algorithm allows for continuous state/action spaces and is more stable towards stochastic variation. Another major difference is the step towards higher degree of realism with the implementation of price/market in the environment. The framework is flexible and the algorithm is easily applied to different time scales or extended with additional constraints e.g. cascading reservoirs, minimum production requirements, non-linear production functions, or ramp-up costs. In addition, the neural network structure can be changed to take temporal aspects into account (recurrent neural networks) or structured information in the form of graphs.

For cascading reservoirs or systems with external water flow restrictions, the constraints must be implemented in the environment. This can either be as hard constraints on allowed actions or via combination into a single reward function (e.g. as in [21]). For cascading systems, the agent can be given control of the entire system, in which case the action space must be expanded accordingly to include production and active spillage in all units. Due to the increased size of the action space, it is expected that the model training will take longer and the training data must naturally reflect the topology.

The agent can be further informed by combining the training data with weather forecasts or traditional models for e.g. forecasting inflow and demand. In addition, one can take advantage of transfer learning [22], where models can be pre-trained on data from one region and then transferred to similar regions reducing the training time at each location.

We see a huge application potential for reinforcement learning in situations where the current solution to handle system complexity is aggregation and dis-aggregation of models. When the entire system becomes too complex and computationally heavy for a single model, the individual reservoirs are aggregated in the model and their joint production scheduled. Reinforcement learning could then be applied for the dis-aggregation of the model and the scheduling for the individual reservoirs given constraints from the aggregated model. It is straight forward to couple reservoirs and add an overall production criteria in the framework.

VI. CONCLUSION

We have explored reinforcement learning to provide strategies for long term scheduling for hydropower production. A soft actor-critic agent is able to learn a meaningful policy on both artificial and real-world training scenarios. While the method is not yet perfect, its flexibility and ability to generalise even complex scenarios gives it a potential to complement traditional optimisation methods for hydropower scheduling.

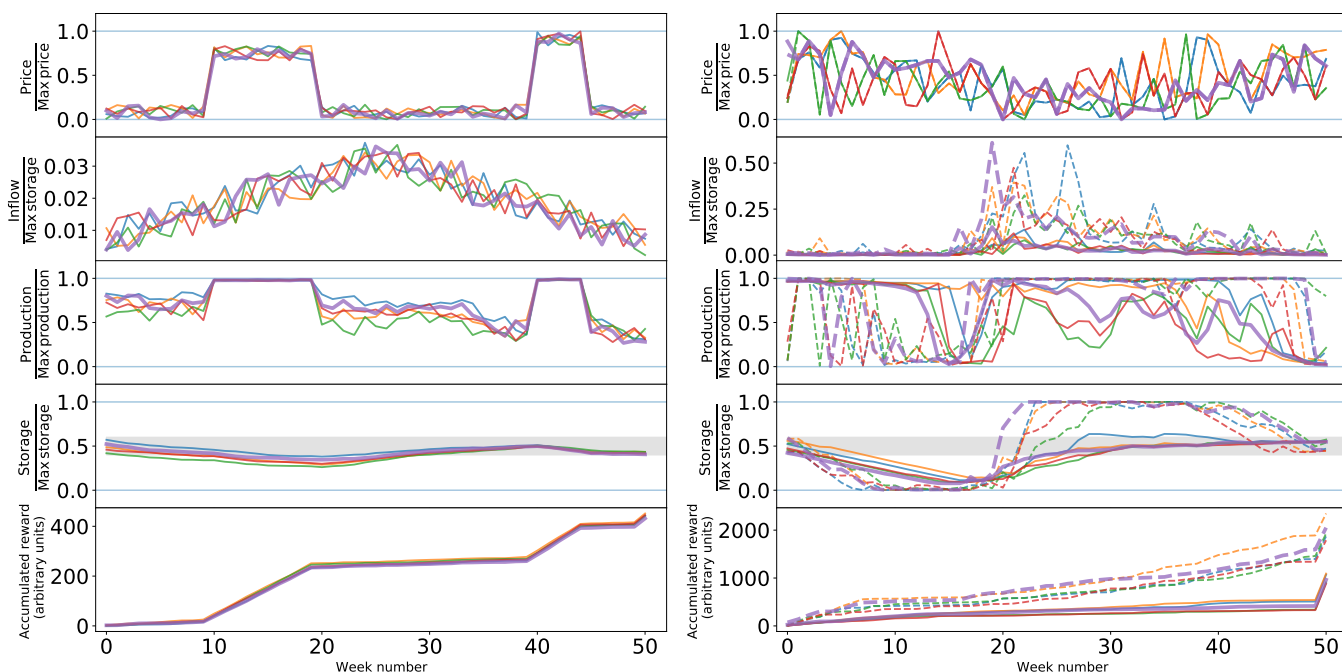


Fig. 4. Examples of price, inflows, actions, storage and accumulated reward per week for five different scenarios randomly sampled from the training data. *Left*: An artificial price scenario with instantaneous price changes and clearly defined seasonal variation of inflow. *Right*: Prices and inflows based on historic Nordic data with some seasonal trends but also a large variation. The solid lines represent a production/storage ratio of 30/1000 and yearly inflow of 1000, and the dashed lines represents 100/1000 and yearly inflow of 4000. The high inflow scenario has more water available and consequently higher rewards.

ACKNOWLEDGMENT

The authors thanks Eivind Bøhn for enlightening discussions, and the reviewers for their constructive feedback.

REFERENCES

- [1] O. Wolfgang, A. Haugstad, B. Mo, A. Gjelsvik, I. Wangensteen, and G. Doorman, "Hydro reservoir handling in Norway before and after deregulation," *Energy*, vol. 34, no. 10, pp. 1642 – 1651, 2009, 11th Conference on Process Integration, Modelling and Optimisation for Energy Saving and Pollution Reduction.
- [2] A. Helseth, M. Fodstad, and B. Mo, "Optimal medium-term hydropower scheduling considering energy and reserve capacity markets," *IEEE Transactions on Sustainable Energy*, vol. 7, no. 3, pp. 934–942, 2016.
- [3] R. S. Sutton and A. G. Barto, *Reinforcement learning : an introduction*. MIT Press, Cambridge, MA, 2018. [Online]. Available: <http://incompleteideas.net/book/the-book.html>
- [4] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, and et al., "Human-level control through deep reinforcement learning," *Nature*, vol. 518, pp. 529 EP –, 02 2015.
- [5] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. van den Driessche, and et al., "Mastering the game of go with deep neural networks and tree search," *Nature*, vol. 529, pp. 484 EP –, 01 2016.
- [6] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, and et al., "Continuous control with deep reinforcement learning," *Unpublished*, 2015. [Online]. Available: <https://arxiv.org/abs/1509.02971>
- [7] A. Irpan, "Deep reinforcement learning doesn't work yet!" <https://www.alexirpan.com/2018/02/14/rl-hard.html>, 2018.
- [8] C. J. C. H. Watkins, "Learning from delayed rewards," PhD Thesis, University of Cambridge, 1989.
- [9] C. J. C. H. Watkins and P. Dayan, "Q-learning," *Machine Learning*, vol. 8, no. 3-4, pp. 279–292, may 1992.
- [10] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, "Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor," *Unpublished*, 2018. [Online]. Available: <https://arxiv.org/abs/1801.01290>
- [11] T. Haarnoja, A. Zhou, K. Hartikainen, G. Tucker, S. Ha, J. Tan, and et al., "Soft actor-critic algorithms and applications," *Unpublished*, 2018. [Online]. Available: <https://arxiv.org/abs/1812.05905>
- [12] B. Polyak and A. Juditsky, "Acceleration of stochastic approximation by averaging," *SIAM Journal on Control and Optimization*, vol. 30, no. 4, pp. 838–855, 1992.
- [13] V. V. Kumar, "Soft actor critic demystified," Jan. 2018. [Online]. Available: <https://github.com/vaishak2future/sac>
- [14] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, and et al., "Pytorch: An imperative style, high-performance deep learning library," in *Advances in Neural Information Processing Systems 32*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, Eds. Curran Associates, Inc., 2019, pp. 8024–8035.
- [15] R. H. R. Hahnloser, R. Sarpeshkar, M. A. Mahowald, R. J. Douglas, and H. S. Seung, "Digital selection and analogue amplification coexist in a cortex-inspired silicon circuit," *Nature*, vol. 405, no. 6789, pp. 947–951, 2000.
- [16] T. Tieleman and G. Hinton, "Lecture 6.5 - rmsprop, coursera: Neural networks for machine learning," 2012. [Online]. Available: https://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides_lec6.pdf
- [17] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *Unpublished*, 2014. [Online]. Available: <https://arxiv.org/abs/1412.6980>
- [18] J. Achiam, "Spinning up in deep rl," <https://spinningup.openai.com/en/latest/algorithms/sac.html>, 2018.
- [19] P. Côté and R. Leconte, "Comparison of stochastic optimization algorithms for hydropower reservoir operation with ensemble streamflow prediction," *Journal of Water Resources Planning and Management*, vol. 142, no. 2, p. 04015046, 2016.
- [20] A. Castelletti, S. Galelli, M. Restelli, and R. Soncini-Sessa, "Tree-based reinforcement learning for optimal water reservoir operation," *Water Resources Research*, vol. 46, no. 9, 2010.
- [21] J.-H. Lee and J. W. Labadie, "Stochastic optimization of multireservoir systems via reinforcement learning," *Water Resources Research*, vol. 43, no. 11, 2007.
- [22] S. J. Pan and Q. Yang, "A survey on transfer learning," *IEEE Transactions on Knowledge and Data Engineering*, vol. 22, no. 10, pp. 1345–1359, 2010.