Unrolled Policy Iteration Via Graph Filters

Sergio Rozada

King Juan Carlos University sergio.rozada@urjc.es

Samuel Rey

King Juan Carlos University samuel.rev.escudero@urjc.es

Miguel Alcocer

King Juan Carlos University m.alcocer.2022@alumnos.urjc.es

Gonzalo Mateos

University of Rochester qmateosb@ur.rochester.edu

Antonio G. Marques

King Juan Carlos University antonio.garcia.marques@urjc.es

Abstract

Dynamic programming (DP) is a cornerstone for solving Markov decision processes (MDPs) through Bellman's optimality equations. Classical algorithms such as policy iteration exploit their fixed-point structure but become costly in large state—action spaces or with long-term dependencies. We propose BellNet, a parametric model that unrolls and truncates policy iterations, trained to minimize the Bellman error from random value function initializations. By interpreting the MDP transition matrix as the adjacency of a weighted directed graph, we leverage graph signal processing to re-parameterize BellNet as a cascade of nonlinear graph filters, offering a concise and transferable representation of policy and value iteration with explicit control of inference complexity. Experiments in grid environments show that BellNet approximates optimal policies in far fewer iterations than classical methods and generalizes, without retraining, to related unseen tasks.

1 Introduction

Dynamic programming (DP), widely applied across engineering domains [1], is often cast as a Markov decision process (MDP) [2]. Its central task is to solve Bellman's equations (BEQs) for the value function (VF), which encodes long-term cumulative rewards. Since BEQs are fixed-point equations, classical DP methods rely on iterative algorithms [2, 3], where state transitions induce a natural digraph structure. While effective, these algorithms can become expensive. The number of required iterations scales rapidly with the size of the state—action space and worsens in long-horizon problems. To mitigate this challenges, in this work we leverage algorithm unrolling [4, 5] and graph signal processing (GSP) [6, 7] to design learnable neural architectures that reduce the iteration burden. Algorithm unrolling reformulates iterative updates as a finite sequence of layers, retaining the interpretability of model-based methods while introducing trainable components [8, 9].

We introduce BellNet, which unrolls the steps of policy and value iteration into a deep parametric model. Drawing from GSP, we show that each policy iteration update can be expressed as a polynomial in the transition matrix, which can be viewed as the adjacency of a weighted digraph, followed by a nonlinearity. Thus, unrolled policy iteration corresponds to a cascade of nonlinear graph filters [10]. This perspective yields architectures that (i) approximate policy iteration with fewer steps, and (ii) generalize across related environments without retraining. To summarize, our contributions are:

39th Conference on Neural Information Processing Systems (NeurIPS 2025) Workshop: New Perspectives in Advancing Graph Machine Learning.

- C1 We introduce BellNet, an unrolled version of policy iteration structured as a cascade of nonlinear graph filters;
- C2 We put forth a learning problem, where the filter coefficients are trained to minimize the so-termed Bellman error from random VF initializations; and
- **C3** We experimentally show in a grid-world setting that the learned BellNet model converges in significantly fewer iterations and generalizes well to similar environments.

Prior work. In reinforcement learning (RL), unrolling has been used in image-based settings [11], and to learn the MDP topology by interpreting the transition matrix as a graph [12, 13]. Unlike our work, existing approaches (a) focus on value iteration; (b) address RL rather than DP, thus they estimate transition probabilities instead of exploiting the graph structure to design the unrolled architecture; and (c) target single tasks instead of enabling generalization across MDPs. Prior RL works have used GSP tools to improve algorithmic efficiency. For instance, [14] postulates the VFs lie in a low-dimensional subspace induced by the state transition digraph; [15] estimates the optimal policy on a subset of states and extends it via graph interpolation; and [16] applies graph reduction to simplify the decision process. While effective, these methods are task-specific. In contrast, BellNet is task-agnostic and applicable across different MDPs. Finally, a growing body of work in GSP investigates the properties of graph filters, e.g., permutation equivariance, stability, or transferability [17–22]. We empirically show that BellNet inherits some of these desirable properties, although a deeper theoretical analysis is left for future work.

2 Preliminaries: Fundamentals of DP and GSP

DP. In DP, we consider an MDP defined by the tuple $(\mathcal{S}, \mathcal{A}, \mathbf{P}, \mathbf{R})$, where \mathcal{S} and \mathcal{A} are discrete state and action spaces, $\mathbf{P} \in [0,1]^{|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}|}$ is a known transition probability matrix whose rows, indexed by state-action pairs (s,a), define distributions over next states s', and $\mathbf{R} \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}|}$ contains the rewards. Solving the MDP amounts to finding a policy $\pi: \mathcal{S} \mapsto [0,1]^{|\mathcal{A}|}$ that maximizes the VFs, defined as expected cumulative rewards. A policy maps each state s to a distribution over actions a, and the VF under π is given by $Q^{\pi}(s,a) = \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^{t} r_{t} \mid s_{0} = s, a_{0} = a \right]$, where $\gamma \in [0,1]$ is a discount factor and the instantaneous reward r_{t} is the entry of \mathbf{R} indexed by the state-action pair at time t. We arrange policy probabilities in the matrix $\mathbf{\Pi} \in [0,1]^{|\mathcal{S}| \times |\mathcal{A}|}$ and the VFs in $\mathbf{Q}_{\pi} \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}|}$. For convenience, we henceforth use the vectorizations $\mathbf{r} = \text{vec}(\mathbf{R})$ and $\mathbf{q}_{\pi} = \text{vec}(\mathbf{Q}_{\pi})$. The BEQs characterize the VFs \mathbf{q}_{π} for a fixed policy π [23]. Denoting $\mathbf{P}_{\pi} = \mathbf{P}(\mathbf{I} \odot \mathbf{\Pi}^{\top})^{\top}$, where \odot is the Khatri-Rao product and \mathbf{I} the identity matrix, we have that

$$\mathbf{q}_{\pi} = \mathbf{r} + \gamma \mathbf{P}_{\pi} \mathbf{q}_{\pi}. \tag{1}$$

This fixed-point linear system of equations can be solved iteratively. Iterating until convergence is referred to as *policy evaluation* in DP parlance. Greedy maximization of \mathbf{Q}_{π} with respect to actions (columns) produces a new policy $\mathbf{\Pi}'$, i.e.,

$$\Pi'_{i,j} = \begin{cases} 1 & \text{if } j = \arg\max_k Q_{ik}, \\ 0 & \text{otherwise.} \end{cases}$$
(2)

This step, known as *policy improvement*, produces a policy Π' that is guaranteed to outperform Π in terms of the attained VFs [3]. Crucially, if $\Pi' = \Pi$, then $\Pi = \Pi^*$ is optimal, i.e., attains the maximum VFs $\mathbf{Q}_{\pi} = \mathbf{Q}^*$ for all state-action pairs. This underpins *policy iteration*, an iterative method that alternates policy evaluation and policy improvement to compute the optimal VFs. Interestingly, for the optimal VFs \mathbf{Q}^* , it also holds that

$$\mathbf{q}^{\star} = \mathbf{r} + \gamma \mathbf{P} \mathbf{v}^{\star} \quad \text{with} \quad v_i^{\star} = \max_k Q_{ik}^{\star}.$$
 (3)

This defines a nonlinear fixed-point system that can be solved iteratively through a procedure known as *value iteration*. Value iteration is equivalent to performing one step of the policy evaluation iteration followed by policy improvement [23].

GSP. A graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ is defined by a set of N nodes \mathcal{V} and a set of edges $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$. The connectivity of \mathcal{G} is captured by the sparse adjacency matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$, where $A_{ij} \neq 0$ if and only if $(i,j) \in \mathcal{E}$, and the entry A_{ij} denotes the weight of the edge from node i to node j. A graph signal

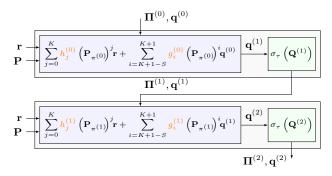


Figure 1: BellNet: A cascade of learnable graph filters and softmax nonlinearities that unrolls policy iteration.

is a function defined on the set of nodes, represented as a vector $\mathbf{x} \in \mathbb{R}^N$, where x_i denotes the signal value at node i. Graph filters are linear, topology-aware operators that process graph signals. They can be expressed as matrix polynomials of the adjacency matrix \mathbf{A} [10, 24], namely

$$\mathbf{H} = \sum_{j=0}^{N-1} h_j \mathbf{A}^j,\tag{4}$$

where $\mathbf{h} = [h_0, \dots, h_{N-1}]^{\top}$ is the vector of filter coefficients. Since each power \mathbf{A}^j encodes information about the j-hop neighborhood of \mathcal{G} , the output $\mathbf{y} = \mathbf{H}\mathbf{x}$ can be interpreted as a diffusion (or aggregation) of the input signal \mathbf{x} across neighborhoods of increasing size, with the coefficients h_j weighting the contribution from each j-hop component [25].

3 Unrolling DP via GSP

Algorithm unrolling is a foundational technique for infusing model-based inductive bias into datadriven learning [8]. Given an iterative algorithm, unrolling builds a parametric mapping, typically a neural network, by assigning each iteration to a corresponding block, such as a network layer. The operations of the original algorithm are preserved and reinterpreted as layer-wise computations, enabling the model to learn algorithm-specific behavior from data. Next, we unroll policy iteration and draw GSP connections in the process. Each unrolled block consists of two main steps: policy evaluation, which involves solving (1), and policy improvement, where (2) is applied.

Policy evaluation. The BEQ (1) is a linear system of equations. However, solving it directly is often impractical due to the large size of state-action spaces. Instead, one typically iterates by applying the right-hand-side (rhs) of (1) repeatedly until convergence is reached. Additional simplifications exploit structural properties of the MDP, such as linear dynamics [26, 27], low-rank structure [28–30], or kernel-based representations [31, 32]. Here, instead, we propose leveraging the graph structure of the MDP. To elucidate this connection, we expand the BEQ recursion as follows

$$\mathbf{q}^{(k)} = \mathbf{r} + \gamma \mathbf{P}_{\pi} \mathbf{q}^{(k-1)} = \mathbf{r} + \gamma \mathbf{P}_{\pi} \mathbf{r} + \gamma^{2} (\mathbf{P}_{\pi})^{2} \mathbf{q}^{(k-2)}$$

$$= \dots = \sum_{j=0}^{k-1} \gamma^{j} (\mathbf{P}_{\pi})^{j} \mathbf{r} + \gamma^{k} (\mathbf{P}_{\pi})^{k} \mathbf{q}^{(0)}.$$
(5)

This expression comprises an exponentially decaying bias $\mathbf{b}^{(k)} := \gamma^k \left(\mathbf{P}_{\pi}\right)^k \mathbf{q}^{(0)}$ that depends on $\mathbf{q}^{(0)}$; and a graph filter $\mathbf{H}^{(k)} := \sum_{j=0}^{k-1} \gamma^j \left(\mathbf{P}_{\pi}\right)^j$ applied to \mathbf{r} . The latter characterization follows since $\mathbf{H}^{(k)}$ is a polynomial of \mathbf{P}_{π} , which represents the adjacency matrix of a weighted digraph \mathcal{G} . The nodes are state-action pairs while the edge weights correspond to the Markovian transition probabilities \mathbf{P} and the current policy $\mathbf{\Pi}$. From this viewpoint, the powers of the discount factor γ act as the filter coefficients in (4), i.e., $h_j = \gamma^j$. Consequently, policy evaluation can be interpreted as applying a graph filter to the reward. Due to the fixed-point theorem [3], an infinite-order filter is guaranteed to recover the true VF for policy π , so that $\mathbf{q}_{\pi} = \mathbf{H}^{(\infty)}\mathbf{r} + \mathbf{b}^{(\infty)} = \sum_{j=0}^{\infty} \gamma^j \left(\mathbf{P}_{\pi}\right)^j \mathbf{r}$.

Moreover, our GSP perspective enables concrete simplifications of the proposed model. While the graph filter underlying policy evaluation is, in principle, of infinite degree, an equivalent filter with limited degree exists.

Proposition 1. Value function \mathbf{q}^{π} for a fixed policy π can be expressed as a finite-order graph filter

$$\mathbf{q}_{\pi} = \sum_{j=0}^{\infty} \gamma^{j} (\mathbf{P}_{\pi})^{j} \mathbf{r} = \sum_{j=0}^{K} \bar{h}_{j} (\mathbf{P}_{\pi})^{j} \mathbf{r}, \tag{6}$$

with $K \leq |\mathcal{S}||\mathcal{A}|$. If \mathbf{P}_{π} is diagonalizable, then $K \leq |\mathcal{S}|$.

Proof. By the Cayley–Hamilton theorem [33], any matrix polynomial of $\mathbf{P}_{\pi} \in [0,1]^{|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}||\mathcal{A}|}$ can be reparameterized as a polynomial of degree at most $K = |\mathcal{S}||\mathcal{A}|$. If \mathbf{P}_{π} is diagonalizable, the degree of its minimal polynomial is at most $\mathrm{rank}(\mathbf{P}_{\pi}) = \mathrm{rank}(\mathbf{P}) = |\mathcal{S}|$, so any polynomial of \mathbf{P}_{π} can be expressed with order at most $K = |\mathcal{S}|$.

Beyond exact policy evaluation, our approach also encompasses approximate policy evaluation via early stopping after a fixed number of iterations. Recall that value iteration corresponds to a single application of the rhs of (1). In any case, early stopping is equivalent to a graph filter of some order K and fixed coefficients $h_j = \gamma^j$. This also introduces a non-vanishing bias term that must be accounted for [cf. (5)]. Furthermore, the estimate $\hat{\mathbf{q}}_{\pi}$ may not converge to the true VF \mathbf{q}_{π} , so it must be reused to initialize the next policy evaluation under the updated policy $\mathbf{\Pi}'$. Identity (6) can be extended to this case by explicitly incorporating the bias term as

$$\hat{\mathbf{q}}_{\pi} = \sum_{j=0}^{K} h_j(\mathbf{P}_{\pi})^j \mathbf{r} + h_{K+1} \left(\mathbf{P}_{\pi}\right)^{K+1} \mathbf{q}^{(0)}.$$
 (7)

Although the filter in (7) is well motivated, certain problems may benefit from greater expressive capacity. We therefore consider a more general layer that replaces the single bias term $h_{K+1}(\mathbf{P}_{\pi})^{K+1}\mathbf{q}^{(0)}$ with a second graph filter acting on $\mathbf{q}^{(0)}$, yielding

$$\hat{\mathbf{q}}_{\pi} = \sum_{j=0}^{K} h_{j} \left(\mathbf{P}_{\pi} \right)^{j} \mathbf{r} + \sum_{i=K+1-S}^{K+1} g_{i} \left(\mathbf{P}_{\pi} \right)^{i} \mathbf{q}^{(0)}.$$
 (8)

Here, $\{h_j\}_{j=0}^K$ are the coefficients of the filter applied to \mathbf{r} , and $\{g_i\}_{i=K-S+1}^{K+1}$ those of the additional filter applied to $\mathbf{q}^{(0)}$. Restricting the second filter to $i \in [K+1-S, K+1]$ (with $0 \le S \le K+1$) ensures that this parametrization recovers (7) as a special case when S=0 (with $g_{K+1}=h_{K+1}$).

Policy improvement. As defined in (2), policy improvement is a nonlinear row-wise max operation applied to \mathbf{Q}_{π} , analogous to max-pooling, selecting the maximum in each row. For differentiability, we replace the max operation with a softmax, as detailed in the next section.

4 BellNet: Learning Policy Iteration

Through the GSP lens, policy iteration is a cascade of nonlinear graph filtering operations that converge to the optimal VFs of the MDP. This perspective motivates BellNet, our proposed unrolling of policy iteration to solve BEQs. BellNet is a deep architecture composed of L+1 layers. Each layer takes as input a VF vector $\mathbf{q} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ and its associated softmax policy $\mathbf{\Pi} \in [0,1]^{|\mathcal{S}|\times|\mathcal{A}|}$, and outputs an enhanced VF vector and policy. The mapping between the input and output of the l-th layer is implemented by graph filters with learnable coefficients $\mathbf{h}^{(l)}$ and $\mathbf{g}^{(l)}$. Formally, let $\bar{\mathbf{q}}$ be the (possibly random) initial estimate of the VF, and let $\mathcal{H} = \{(\mathbf{h}^{(l)}, \mathbf{g}^{(l)})\}_{l=0}^L$ collect all the learnable coefficients. Then, BellNet, repented by the mapping $\Phi(\cdot; \mathcal{H})$, implements $\{\hat{\mathbf{q}}, \hat{\mathbf{\Pi}}\} := \Phi(\bar{\mathbf{q}}; \mathcal{H})$ with $\hat{\mathbf{q}} = \mathbf{q}^{(L+1)}$, $\hat{\mathbf{\Pi}} = \mathbf{\Pi}^{(L+1)}$, $\mathbf{q}^{(0)} = \bar{\mathbf{q}}$, and layer-wise updates:

$$\begin{split} \mathbf{q}^{(l+1)} &= \sum_{j=0}^{K} h_{j}^{(l)} \left(\mathbf{P}_{\pi^{(l)}} \right)^{j} \mathbf{r} + \sum_{i=K+1-S}^{K+1} g_{i}^{(l)} \left(\mathbf{P}_{\pi^{(l)}} \right)^{i} \mathbf{q}^{(l)} \\ \mathbf{\Pi}^{(l+1)} &= \sigma_{\tau}(\mathbf{Q}^{(l+1)}), \quad \text{with} \quad [\sigma_{\tau}(\mathbf{Q})]_{ij} = \frac{e^{Q_{ij}/\tau}}{\sum_{k=1}^{|A|} e^{Q_{ik}/\tau}}, \end{split}$$

for $l=0,\ldots,L-1$, where $\mathbf{Q}^{(l)}=\operatorname{unvec}(\mathbf{q}^{(l)}),$ σ_{τ} is a row-wise softmax operator with temperature parameter τ , and $\mathbf{h}^{(l)}=[h_0^{(l)},\ldots,h_{K+1}^{(l)}]$ and $\mathbf{g}^{(l)}=[g_{K+1-S}^{(l)},\ldots,g_{K+1}^{(l)}]$ are the filter coefficients of the l-th layer. Each layer implements two reduced-order, parallel graph filters, sums their respective outputs, and then applies a softmax nonlinearity. The BellNet model is illustrated in Fig. 1. Setting $L=\infty,$ $K=\infty,$ and S=0, with $h_j^{(l)}=\gamma^j$ and $g_i^{(l)}=\gamma^i,$ and replacing the softmax with the max operator, recovers policy iteration. Likewise, setting K=0 and S=0 yields value iteration.

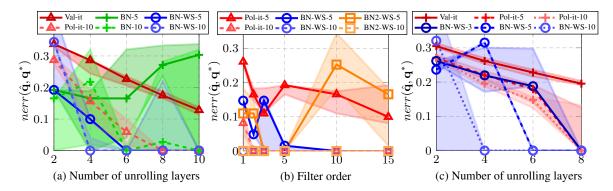


Figure 2: Evaluation of BellNet across different scenarios. We report the median error of the estimated $\hat{\mathbf{q}}$, computed as in (10), over 15 realizations. a) Shows the error as L increases; b) illustrates the error as K increases; and c) evaluates the transferability capacity of BellNet.

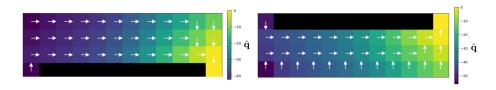


Figure 3: Cliff walking environment (left) and its mirrored version (right). Cliff regions are shown in black; arrows indicate the policy learned by BellNet, and the color map represents the corresponding VFs. BellNet is trained on the top environment, while the policy in the bottom environment is inferred without retraining.

Learning. To complete the approach, we formulate the optimization adopted to learn the filter coefficients \mathcal{H} . The loss function is inspired by temporal difference (TD) methods [34, 35]. We solve a sequential optimization problem that minimizes the Bellman error [36, 37], which is the discrepancy between the left and rhs of (3). Specifically, with n being an iteration index, we solve

$$\mathcal{H}_{[n+1]} = \operatorname{argmin}_{\mathcal{H}} \|\mathbf{r} + \gamma \mathbf{P} \mathbf{v}_{[n]} - \mathbf{\Phi}(\bar{\mathbf{q}}, \mathcal{H})\|_{2}^{2}, \tag{9}$$

where $\mathbf{v}_{[n]}$ is defined as per (3) from $\{\mathbf{q}_{[n]}, \mathbf{\Pi}_{[n]}\} := \mathbf{\Phi}(\bar{\mathbf{q}}, \mathcal{H}_{[n]})$. Note that $\{\mathbf{q}_{[n]}, \mathbf{\Pi}_{[n]}\}$ depends on the current iterate $\mathcal{H}_{[n]}$ and not on the optimized coefficients \mathcal{H} . By slight abuse of notation, $\mathbf{\Phi}$ in (9) refers only to the VF output $\hat{\mathbf{q}}$. We also highlight that: (a) as customary in TD, for each n we update the filter coefficients via gradient descent; (b) our DP method does not require data samples, but the transition probability matrix \mathbf{P} instead; and (c) BellNet is initialized with an arbitrary VF $\bar{\mathbf{q}}$ and trained to converge to the optimal VF and policy regardless of $\bar{\mathbf{q}}$.

Transferability. Graph filters are permutation-equivariant and transferable to larger graphs from a convergent sequence [18], making them particularly well suited to generalize across related problems. In our DP context, this property can be leveraged to train BellNet on a single MDP and deploy it on other similar or larger MDPs. Doing so yields solutions faster than evaluating policies from scratch, as we demonstrate numerically in Section 5. Moreover, the vanilla BellNet described so far operates with a fixed unrolling depth and distinct parameters per block. An attractive alternative is to *share* weights across blocks. Although weight sharing admittedly reduces expressiveness, it markedly decreases the number of learnable parameters [5, 38]. Crucially, it allows the same block to be reused as many times as desired during inference–exceeding the original training depth to enable efficient, scalable transfer as well as to delineate favorable complexity versus policy approximation tradeoffs.

5 Numerical Results and Concluding Remarks

We assess the performance of BellNet in the cliff walking environment, a grid-world setup where the goal is to reach a target location in the minimum number of steps without falling off the grid. The state space \mathcal{S} corresponds to positions on the grid, and the action space consists of moving up, down, left, or right. Two instances of this environment are depicted in Fig. 3. Simulations are conducted

using the Gymnasium library [39, 40]. Experiments were run on an AMD EPYC 9634 (168 threads) server with two NVIDIA RTX 4090 GPUs (24 GB each).

We compute the true VF \mathbf{q}^* using policy iteration with sufficiently many policy evaluation and improvement steps, and report the normalized error defined as

$$nerr(\hat{\mathbf{q}}, \mathbf{q}^*) = \|\hat{\mathbf{q}}/\|\hat{\mathbf{q}}\|_2 - \mathbf{q}^*/\|\mathbf{q}^*\|_2\|_2^2.$$
 (10)

Figure 2 depicts the median error along with the interquartile range (between the 25th and 75th percentiles), computed over 15 random realizations. We compare the performance of BellNet with and without weight sharing (denoted "BN-WS" and "BN" in the legend), as well as value iteration ("Val-it") and policy iteration ("Pol-it"), across multiple scenarios. Furthermore, "BN-WS" denotes BellNet with a single graph filter (i.e., S=0), whereas "BN2-WS" corresponds to the case with two graph filters and S=K.

Test case 1 (Depth). We first examine how increasing the number of unrolling layers (equivalently, the number of policy improvement steps for "Val-it" and "Pol-it") influences performance. Figure 2a shows results using filter orders 5 and 10 for "BN", and 10 policy evaluation updates in "Pol-it". Apparently, the weight sharing strategy leads to better performance with lower variance, whereas distinct filter coefficients results in more unstable behavior. Moreover, BellNet consistently outperforms policy iteration, recovering the optimal policy with only 4 layers compared to 10 required by "Pol-it".

Test case 2 (Filter order). Next, we investigate the role of the filter order in the performance of BellNet. Figure 2b shows the error of "Pol-it" and "BN-WS" as the number of policy evaluation steps and, correspondingly, the filter order, increases as indicated on the x-axis. We evaluate "Pol-it" with 5 and 10 policy improvement steps, and use the same number of unrolling layers for "BN-WS" and "BN2-WS". Comparing "BN-WS" and "BN2-WS", we find that both architectures achieve similar performance when the number of unrolling layers is sufficiently large, although "BN2-WS" exhibits greater instability at L=5. As expected, we also find that a higher filter order improves the performance of "BN-WS", with a smaller order being sufficient when the number of unrolling layers increases. Interestingly, this is not the case for "Pol-it", where the number of policy improvement steps has a greater impact than the number of evaluation steps in this setting. Consistent with the previous experiment, these results highlight how BellNet outperforms "Pol-it" when the number of policy improvement steps is moderately small.

Test case 3 (Transferability). The last experiment inspects BellNet's transferability properties. We train BellNet in the original grid-world setting used in previous test cases, and then use it to predict the optimal policy in a modified environment where the positions of the cliffs, origin, and destination have changed. As shown in Fig. 2c, BellNet successfully predicts the optimal policy in the new environment without requiring retraining. For comparison, we compute the optimal policy in the modified environment using value iteration and policy iteration with 5 and 10 policy evaluation steps. We observe that the error in the estimated $\hat{\bf q}$ decreases as the number of layers (indicated on the x-axis) increases, or, when higher-order filters are used. Notably, BellNet outperforms the classical baselines when both the number of unrolling layers and the filter order are sufficiently large. Overall, these preliminary results show that BellNet not only offers a novel approach to estimating the optimal policy, but also generalizes effectively to other related environments not seen during training.

Work partially funded by the Spanish AEI (10.13039/501100011033) grant PID2022-136887NB-I00, and the Community of Madrid via the Ellis Madrid Unit and grants URJC/CAM F861 and F1180 and TEC-2024/COM-89. A related version of this work appeared at CAMSAP 2025 under the title *Unrolling Dynamic Programming via Graph Filters*. Minor edits of this document were made with the assistance of ChatGPT.

References

- [1] E. V. Denardo. Dynamic Programming: Models and Applications. Courier Corporation, 2012.
- [2] M. L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, 2014.
- [3] D. Bertsekas. *Dynamic Programming and Optimal Control: Volume I*, volume 4. Athena Scientific, 2012.
- [4] K. Gregor and Y. LeCun. Learning fast approximations of sparse coding. In *Int. Conf. Mach. Learning*, pages 399–406, 2010.
- [5] S. Chen, Y. C. Eldar, and L. Zhao. Graph unrolling networks: Interpretable neural networks for graph signal denoising. *IEEE Trans. Signal Process.*, 69:3699–3713, 2021.
- [6] A. Ortega, P. Frossard, J. Kovačević, J. M. F. Moura, and P. Vandergheynst. Graph signal processing: Overview, challenges, and applications. *Proc. IEEE*, 106(5):808–828, 2018.
- [7] G. Leus, A. G. Marques, J. M. F. Moura, A. Ortega, and D. I. Shuman. Graph signal processing: History, development, impact, and outlook. *IEEE Signal Process. Mag.*, 40(4):49–60, 2023.
- [8] V. Monga, Y. Li, and Y. C. Eldar. Algorithm unrolling: Interpretable, efficient deep learning for signal and image processing. *IEEE Signal Process. Mag.*, 38(2):18–44, 2021.
- [9] S. Hadou, N. NaderiAlizadeh, and A. Ribeiro. Robust stochastically-descending unrolled networks. *IEEE Trans. Signal Process.*, 72:5484–5499, 2024.
- [10] E. Isufi, F. Gama, D. I. Shuman, and S. Segarra. Graph filters for signal processing and machine learning on graphs. *IEEE Trans. Signal Process.*, 72:4745–4781, 2024.
- [11] A. Tamar, Y. Wu, G. Thomas, S. Levine, and P. Abbeel. Value iteration networks. In *Conf. Neural Inform. Process. Syst.*, volume 29, 2016.
- [12] S. Niu, S. Chen, H. Guo, C. Targonski, M. Smith, and J. Kovačević. Generalized value iteration networks: Life beyond lattices. In *AAAI Conf. Artificial Intell.*, volume 32, 2018.
- [13] A. Deac, P.-L. Bacon, and J. Tang. Graph neural induction of value iteration. *arXiv preprint* arXiv:2009.12604, 2020.
- [14] L. Liu, A. Chattopadhyay, and U. Mitra. On solving MDPs with large state space: Exploitation of policy structures and spectral properties. *IEEE Trans. Commun.*, 67(6):4151–4165, 2019.
- [15] L. Liu and U. Mitra. Policy sampling and interpolation for wireless networks: A graph signal processing approach. In *Proc. IEEE Global Commun. Conf.*, pages 1–6, 2019.
- [16] M. Levorato, S. Narang, U. Mitra, and A. Ortega. Reduced dimension policy iteration for wireless network control via multiscale analysis. In *Proc. IEEE Global Commun. Conf.*, pages 3886–3892, 2012.
- [17] F. Gama, J. Bruna, and A. Ribeiro. Stability properties of graph neural networks. *IEEE Trans. Signal Process.*, 68:5680–5695, 2020.
- [18] L. Ruiz, F. Gama, and A. Ribeiro. Graph neural networks: Architectures, stability, and transferability. *Proc. IEEE*, 109(5):660–682, 2021. doi: 10.1109/JPROC.2021.3055400.
- [19] R. Levie, W. Huang, L. Bucci, M. Bronstein, and G. Kutyniok. Transferability of spectral graph convolutional neural networks. *J. Mach. Learning Res.*, 22(272):1–59, 2021.
- [20] J. Cervino, L. Ruiz, and A. Ribeiro. Learning by transference: Training graph neural networks on growing graphs. *IEEE Trans. Signal Process.*, 71:233–247, 2023.
- [21] S. Rey, M. Navarro, V. M. Tenorio, S. Segarra, and A. G. Marques. Redesigning graph filter-based GNNs to relax the homophily assumption. In *IEEE Int. Conf. Acoust., Speech and Signal Process.*, pages 1–5, 2025.

- [22] Z. Wang, J. Cervino, and A. Ribeiro. A manifold perspective on the statistical generalization of graph neural networks. In *Int. Conf. Learn. Representations*, 2025.
- [23] R. Bellman. Dynamic programming. Science, 153(3731):34–37, 1966.
- [24] S. Rey, V. M. Tenorio, and A. G. Marques. Robust graph filter identification and graph denoising from signal observations. *IEEE Trans. Signal Process.*, 71:3651–3666, 2023.
- [25] S. Segarra, A. G. Marques, and A. Ribeiro. Optimal graph-filter design and applications to distributed linear network operators. *IEEE Trans. Signal Process.*, 65(15):4117–4131, 2017.
- [26] M. G. Lagoudakis and R. Parr. Least-squares policy iteration. J. Mach. Learning Res., 4(Dec): 1107–1149, 2003.
- [27] A. Geramifard et al. A tutorial on linear function approximators for dynamic programming and reinforcement learning. *Foundations and Swerves*® *in Machine Learning*, 6(4):375–451, 2013.
- [28] A. Agarwal, S. Kakade, A. Krishnamurthy, and W. Sun. FLAMBE: Structural complexity and representation learning of low rank MDPs. In *Conf. Neural Inform. Process. Syst.*, volume 33, pages 20095–20107, 2020.
- [29] S. Rozada, S. Paternain, and A. G. Marques. Tensor and matrix low-rank value-function approximation in reinforcement learning. *IEEE Trans. Signal Process.*, 72:1634–1649, 2024.
- [30] S. Rozada, J. L. Orejuela, and A. G. Marques. Solving finite-horizon MDPs via low-rank tensors. *arXiv preprint arXiv:2501.10598*, 2025.
- [31] D. Ormoneit and S. Sen. Kernel-based reinforcement learning. *Mach. Learn.*, 49(2):161–178, 2002.
- [32] Y. Akiyama and K. Slavakis. Nonparametric Bellman mappings for value iteration in distributed reinforcement learning. arXiv preprint arXiv:2503.16192, 2025.
- [33] R. A. Horn and C. R. Johnson. *Matrix Analysis*. Cambridge University Press, 2012.
- [34] R. S. Sutton. Reinforcement Learning: An Introduction. A Bradford Book, 2018.
- [35] M. Geist and O. Pietquin. Algorithmic survey of parametric value function approximation. *IEEE Trans. Neural Netw. Learning Syst.*, 24(6):845–867, 2013.
- [36] D. Choi and B. Van Roy. A generalized Kalman filter for fixed point approximation and efficient temporal-difference learning. *Discrete Event Dyn. Syst.*, 16(2):207–239, 2006.
- [37] K. Asadi, S. Sabach, Y. Liu, O. Gottesman, and R. Fakoor. TD convergence: An optimization perspective. In *Conf. Neural Inform. Process. Syst.*, volume 36, pages 49169–49186, 2023.
- [38] Steven J. Nowlan and Geoffrey E. Hinton. Simplifying neural networks by soft weight-sharing. Neural Comput., 4(4):473–493, 1992.
- [39] G. Brockman et al. OpenAI Gym. arXiv preprint arXiv:1606.01540, 2016.
- [40] M. Towers et al. Gymnasium: A standard interface for reinforcement learning environments. *arXiv preprint arXiv:2407.17032*, 2024.