Learning One Representation to Optimize All Rewards

Anonymous Author(s) Affiliation Address email

Abstract

1	We introduce the <i>forward-backward</i> (FB) representation of the dynamics of a
2	reward-free Markov decision process. It provides explicit near-optimal policies for
3	any reward specified a posteriori. During an unsupervised phase, we use reward-
4	free interactions with the environment to learn two representations via off-the-shelf
5	deep learning methods and temporal difference (TD) learning. In the test phase, a
6	reward representation is estimated either from reward observations or an explicit
7	reward description (e.g., a target state). The optimal policy for that reward is
8	directly obtained from these representations, with no planning. We assume access
9	to an exploration scheme or replay buffer for the first phase.

The unsupervised FB loss is well-principled: if training is perfect, the policies obtained are provably optimal for any reward function. With imperfect training, the sub-optimality is proportional to the unsupervised approximation error. The FB representation learns long-range relationships between states and actions, via a predictive occupancy map, without having to synthesize states as in model-based approaches.

This is a step towards learning controllable agents in arbitrary black-box stochastic environments. This approach compares well to goal-oriented RL algorithms on discrete and continuous mazes, pixel-based MsPacman, and the FetchReach virtual robot arm. We also illustrate how the agent can immediately adapt to new tasks beyond goal-oriented RL.

1 Introduction and Related Work

We consider one kind of unsupervised reinforcement learning problem: Given a Markov decision process (MDP) but no reward information, is it possible to learn and store a compact object that, for any reward function specified later, provides the optimal policy for that reward, with a minimal amount of additional computation? In a sense, such an object would encode in a compact form the solutions of all possible planning problems in the environment. This is a step towards building agents that are fully controllable after first exploring their environment in an unsupervised way.

Goal-oriented RL methods [ACR⁺17] [PAR⁺18] compute policies for a series of rewards specified in
 advance (such as reaching a set of target states), but cannot adapt in real time to new rewards, such as
 weighted combinations of target states or dense rewards.

Learning a model of the world is another possibility, but it still requires explicit planning for each

new reward; moreover, synthesizing accurate trajectories of states over long time ranges has proven

33 difficult [Tal17] KST⁺18].

Instead, we exhibit an object that is both simpler to learn than a model of the world, and contains the 34

information to recover near-optimal policies for any reward provided a posteriori, without a planning 35 phase.

36

[BBQ⁺18] learn optimal policies for all rewards that are linear combinations of a finite number of 37 feature functions provided in advance by the user. This limits applications: e.g., goal-oriented tasks 38 would require one feature per goal state, thus using infinitely many features in continuous spaces. 39 We reuse a policy parameterization from $[BBQ^{+}18]$, but introduce a novel representation with better 40 properties, based on state occupancy prediction instead of expected featurizations. We use theoretical 41 advances on successor state learning from [BTO21]. We obtain the following. 42

We prove the existence of a learnable "summary" of a reward-free discrete or continuous MDP, 43 that provides an explicit formula for optimal policies for any reward specified later. This takes the 44 form of a pair of representations $F: S \times A \times Z \to Z$ and $B: S \times A \to Z$ from state-actions into 45 a representation space $Z \simeq \mathbb{R}^d$, with policies $\pi_z(s) := \arg \max_a F(s, a, z)^{\mathsf{T}} z$. Once a reward is 46 specified, a value of z is computed from reward values and B; then π_z is used. Rewards may be 47 specified either explicitly as a function, or as target states, or by samples as in usual RL setups. 48

We provide a well-principled unsupervised loss for F and B. If FB training is perfect, then the 49 policies are provably optimal for all rewards (Theorem 1). With imperfect training, sub-optimality 50 is proportional to the FB training error (Theorems 5-6). In finite spaces, perfect training is possible with large enough dimension d (Proposition 3). Explicitly, F and B are trained so that 51 52 $F(s, a, z)^{\top}B(s', a')$ approximates the long-term probability to reach s' from s if following π_z . 53

We provide a TD-like algorithm to train F and B for this unsupervised loss, with function 54 approximation, adapted from recent methods for successor states **BTO21**. No sparse rewards are 55 used: every transition reaches some state s', so every step is exploited. As usual with TD, learning 56 seeks a fixed point but the loss itself is not observable. 57

• We prove viability of the method on several environments from mazes to pixel-based MsPacman 58 and a virtual robotic arm. For single-state rewards (learning to reach arbitrary states), we provide 59 quantitative comparisons with goal-oriented methods such as HER. (Our method is not a substitute 60 for HER: in principle they could be combined, with HER improving replay buffer management 61 for our method.) For more general rewards, which cannot be tackled a posteriori by trained 62 goal-oriented models, we provide qualitative examples. 63

We also illustrate qualitatively the sub-optimalities (long-range behavior is preserved but local 64 blurring of rewards occurs) and the representations learned. 65

[BBQ⁺18] learn optimal policies for rewards that are linear combinations of a Related work. 66 finite number of feature functions provided in advance by the user. This approach cannot tackle 67 generic rewards or goal-oriented RL: this would require introducing one feature per possible goal 68 state, requiring infinitely many features in continuous spaces. 69

Our approach does not require user-provided features describing the future tasks, thanks to using 70 successor states [BTO21] where [BBQ⁺18] use successor features. Schematically, and omitting 71 actions, successor features start with user-provided features φ , then learn ψ such that $\psi(s_0) =$ 72 $\sum_{t\geq 0} \gamma^t \mathbb{E}[\varphi(s_t) \mid s_0]$. This limits applicability to rewards that are linear combinations of φ . Here we use successor *state* probabilities, namely, we learn two representations F and B such that $F(s_0)^{\top}B(s') = \sum_{t\geq 0} \gamma^t \Pr(s_t = s' \mid s_0)$. This does not require any user-provided input. 73 74 75

Thus we learn two representations instead of one. The learned backward representation B is absent 76

from [BBQ⁺18]. B plays a different role than the user-provided features φ of [BBQ⁺18]: learning 77

of F given B is not learning of ψ given φ . The features φ of [BBQ⁺18] are split between our learned 78 B and the functions φ we can use if the reward is known to depend only on some features of the state.

79

We use a similar parameterization of policies by $F(s, a, z)^{T}z$ as in [BBQ+18], for similar reasons, 80 although z encodes a different object. 81

Successor representations where first defined in [Day93] for finite spaces, corresponding to an older 82 object from Markov chains, the fundamental matrix [KS60] Bré99 [GS97]. [SBG17] argue for their 83 84 relevance for cognitive science. For successor representations in continuous spaces, a finite number of features φ are specified first; this can be used for generalization within a family of tasks, e.g., 85 [BDM⁺17], ZSBB17, GHB⁺19] [HDB⁺19]. [BTO21] moves from successor features to successor 86 states by providing pointwise occupancy map estimates even in continuous spaces, without using the 87

sparse reward $\mathbb{1}_{s_t=s'}$. We borrow a successor state learning algorithm from **[BTO21**]. **[BTO21**] also 88

introduced simpler versions of F and B for a single, fixed policy; **BTO21** does not consider the every-optimal-policy setting.

⁹¹ There is a long literature on goal-oriented RL. For instance, [SHGS15] learn goal-dependent value

⁹² functions, regularized via an explicit matrix factorization. Goal-dependent value functions have

⁹³ been investigated in earlier works such as [FD02] and [SMD⁺11]. Hindsight experience replay

94 (HER) [ACR⁺17] improves the sample efficiency of multiple goal learning with sparse rewards. A

⁹⁵ family of rewards has to be specified beforehand, such as reaching arbitrary target states. Specifying

⁹⁶ rewards a posteriori is not possible: for instance, learning to reach target states does not extend

97 to reaching the nearest among several goals, reaching a goal while avoiding forbidden states, or 98 maximizing any dense reward.

Hierarchical methods such as options [SPS99] can be used for multi-task RL problems. However,
 policy learning on top of the options is still needed after the task is known.

For finite state spaces, [JKSY20] use reward-free interactions to build a training set that summarizes a finite environment, in the sense that any optimal policies later computed on this training set instead of the true environment are provably ε -optimal, for any reward. They prove tight bounds on the

necessary set size. Policy learning still has to be done afterwards for each reward.

105 2 Problem and Notation

We consider the following informal problem: Given a reward-free MDP (S, A, P, γ) , can we compute a convenient learnable object E such that, once a reward function $r: S \times A \to \mathbb{R}$ is specified, we can easily (with no planning) compute, from E and r, a policy π whose performance is close to maximal?

Here $\mathcal{M} = (S, A, P, \gamma)$ denotes a reward-free Markov decision process with state space S (discrete or continuous), action space A (discrete for simplicity, but this is not essential), transition probabilities P(s'|s, a) from state s to s' given action a, and discount factor $0 < \gamma < 1$ [SB18]. If S is finite, P(s'|s, a) can be viewed as a matrix; in general, for each $(s, a) \in S \times A$, P(ds'|s, a) is a probability measure on $s' \in S$. The notation P(ds'|s, a) covers all cases. All functions are assumed to be measurable.

Given $(s_0, a_0) \in S \times A$ and a policy $\pi: S \to \operatorname{Prob}(A)$, we denote $\operatorname{Pr}(\cdot|s_0, a_0, \pi)$ and $\mathbb{E}[\cdot|s_0, a_0, \pi]$ the probabilities and expectations under state-action sequences $(s_t, a_t)_{t\geq 0}$ starting with (s_0, a_0) and following policy π in the environment, defined by sampling $s_t \sim P(ds_t|s_{t-1}, a_{t-1})$ and $a_t \sim \pi(s_t)$.

Given a reward function $r: S \times A \to \mathbb{R}$, the Q-function of π for r is $Q_r^{\pi}(s_0, a_0) := \sum_{t \ge 0} \gamma^t \mathbb{E}[r(s_t, a_t)|s_0, a_0, \pi]$. We assume that rewards are bounded, so that all Q-functions are well-defined. We state the results for deterministic reward functions, but this is not essential. We abuse notation and write greedy policies as $\pi(s) = \arg \max_a Q(s, a)$ instead of $\pi(s) \in \arg \max_a Q(s, a)$. Ties may be broken any way.

3 Encoding All Optimal Policies via the Forward-Backward Representation

We first present the forward-backward (FB) representation of an MDP as a way to summarize all optimal policies via explicit formulas. The resulting learning procedure is described in Section 4.

We set a representation space $Z = \mathbb{R}^d$, used to represent both states and reward functions. We learn a pair of "forward" and "backward" representations

$$F: S \times A \times Z \to Z, \qquad B: S \times A \to Z \tag{1}$$

For each $z \in Z$, we define the policy $\pi_z(s) := \arg \max_a F(s, a, z)^{\top} z$.

The main idea is to train F and B such that $F(s_0, a_0, z)^{\top}B(s', a')$ is approximately the long-term probability to reach (s', a') if starting at (s_0, a_0) and following policy π_z . Then, by Theorem I, for any reward function x, the optimal policy π_z is the policy π_z with $z := \mathbb{E}[r(a, g)B(a, g)]$

any reward function r, the optimal policy for r is the policy π_z with $z := \mathbb{E}[r(s, a)B(s, a)]$.

Intuitively, F represents the future of a state under a certain policy. B represents the past of a state, or the ways to reach that state (Appendix B.4). If $F^{\top}B$ is large, then it is possible to reach the second state from the first. This is akin to a model of the environment, without synthesizing state trajectories. In short, if we can learn two representations F and B of state-actions such that $F^{\top}B$ approximates the long-term transitions of the policies π_z , then we can compute all optimal policies from F and B. ¹³⁷ More precisely, for any policy π and state-action (s_0, a_0) , define the *successor measure* $M^{\pi}(s_0, a_0, \cdot)$ ¹³⁸ as the measure over $S \times A$ representing the expected discounted time spent in each set $X \subset S \times A$:

$$M^{\pi}(s_0, a_0, X) := \sum_{t \ge 0} \gamma^t \Pr\left((s_t, a_t) \in X \mid s_0, a_0, \pi\right)$$
(2)

for each $X \subset S \times A$. Viewing M as a measure deals with both discrete and continuous spaces. M^{π} can be learned via a Bellman equation (Section 4). In practice, M^{π} will be represented by a function $m^{\pi}(s_0, a_0, s', a')$ taking a pair of state-actions and returning a number. Namely, assume that we can sample state-action pairs from some unknown distribution $\rho(ds, da)$, e.g., from a replay buffer or exploration policy. We define m^{π} as the density of M^{π} with respect to ρ :

$$M^{\pi}(s_0, a_0, \mathrm{d}s', \mathrm{d}a') =: m^{\pi}(s_0, a_0, s', a') \,\rho(\mathrm{d}s', \mathrm{d}a') \tag{3}$$

so that m^{π} is an ordinary function (or a distribution, see Appendix **B.6**).

The next theorem states that, if $F^{\top}B \approx m$, then all optimal policies can be read on F and B.

Theorem 1 (Forward-backward representation of an MDP). Consider an MDP with state space S and action space A. Let $Z = \mathbb{R}^d$ be some representation space. Let $F: S \times A \times Z \to Z$ and $B: S \times A \to Z$ be two functions. For each $z \in Z$, define the policy $\pi_z(s) := \arg \max_a F(s, a, z)^T z$.

Let ρ be any probability distribution on $S \times A$ (e.g., the distribution of state-actions under some exploration scheme), with full support.

Assume that F and B have been chosen (trained) to satisfy the following: for any $z \in Z$, and any state-actions (s, a) and (s', a'), the quantity $F(s, a, z)^{\top}B(s', a')$ is equal to the successor state density $m^{\pi_z}(s, a, s', a')$ of policy π_z with respect to ρ , defined by (2)–(3).

154 Then, for any bounded reward function $r: S \times A \to \mathbb{R}$, the following holds. Set

$$z_R := \mathbb{E}_{(s,a)\sim\rho} \left[r(s,a)B(s,a) \right]. \tag{4}$$

155 Then π_{z_R} is an optimal policy for reward r in the MDP. Moreover, the optimal Q-function Q^* for 156 reward r is

$$Q^{\star}(s,a) = F(s,a,z_R)^{\mathsf{T}} z_R.$$
⁽⁵⁾

In finite spaces, exact solutions F and B exist (Appendix, Prop. 3), provided the dimension d is large enough. In infinite spaces, arbitrarily good approximations can be obtained by increasing d, corresponding to a rank-d approximation of the cumulated transition probabilities m^{π} . The theoretical guarantee extends to approximate training of F and B, with optimality gap proportional to $F^{T}B - m^{\pi_{z}}$ (Appendix, Theorems 5–6 with various norms on $F^{T}B - m^{\pi}$ and r). For instance, if, for some reward r, the error $|F(s, a, z_R)|^{1}B(s', a') - m^{\pi_{z_R}}(s, a, s', a')|$ is at most ε on average over ($s', a') \sim \rho$ for every (s, a), then π_{z_R} is $3\varepsilon ||r||_{\infty} / (1 - \gamma)$ -optimal for r.

This justifies using some norm over $|F^{\top}B - m^{\pi_z}|$, averaged over $z \in \mathbb{R}^d$, as a training loss for unsupervised reinforcement learning. (Below, we use a fixed rescaled Gaussian over $z \in \mathbb{R}^d$. If prior information is available on the rewards r, the corresponding distribution of z_R may be used instead.) Note that π_z is defined via F, so the equality $F^{\top}B = m^{\pi_z}$ is a fixed point equation.

The dimension d controls how many types of rewards can be optimized well. On the downside, every reward function uncorrelated to the components B_1, \ldots, B_d of B is treated as 0, since $z_R = 0$. If Bis fixed in advance and only F is optimized, the method has similar properties to successor features based on B (Appendix B.4). But one may set a large d and let B be learned to approximate m^{π} : arguably, by Theorem [], the resulting features "linearize" optimal policies as much as possible.

The algorithm is linear in d, so d can be taken as large as the neural network models can handle. Even a relatively small d provides useful behaviors: in the experiments, $d \approx 100$ manages navigation in a pixel-based environment with a huge state space. Appendix B.2 argues theoretically that $d \approx n$ is enough for navigation in dimension n. In practice, for small d we notice some blurring of rewards between nearby states (Fig. 3), for reasons discussed in Section 4.

A similar statement holds with $\bar{m}(s, z, s', a') + F(s, a, z)^{\mathsf{T}}B(s', a')$ instead of $F(s, a, z)^{\mathsf{T}}B(s', a')$ to represent m^{π_z} (Appendix, Theorem 2). Here \bar{m} is any function that does not depend on a. Since \bar{m} has no rank restriction, the finite rank approximation only applies to the advantage function.

The features learned in F and B may have broader interest as the features that "most linearize" optimal policies.

Sketch of proof. We give an idea of the proof for finite state spaces; the general proof is in Appendix C For any reward function R and any policy π , the Q-function Q_R^{π} of π on R is equal to $M^{\pi}R$, where we view M^{π} as a matrix and R as a vector over state-actions: this is a consequence of general properties the successor measure (2). If $M^{\pi_z} \approx F(z)^{\top}B$ (viewing both F(z) and B as $d \times (S \times A)$ matrices), then $Q_R^{\pi_z} \approx F(z)^{\top}BR = F(z)^{\top}z_R$ with $z_R := BR$. Specializing to $z = z_R$, we obtain $Q_R^{\pi_{z_R}} \approx F(z_R)^{\top}z_R$. But by definition, π_{z_R} is the argmax of $F(z_R)^{\top}z_R$. Therefore, π_{z_R} is the argmax of its own Q-function for reward R, so it is optimal. The continuous case uses densities m^{π} instead of M^{π} . The remaining question is to train F and B to approximate successor densities.

¹⁹¹ 4 Learning and Using Forward-Backward Representations

Our algorithm starts with an *unsupervised learning phase*, where we learn the representations Fand B in a reward-free way, by observing state transitions in the environment, generated from any exploration scheme. Then, in a *reward estimation phase*, we estimate a policy parameter $z_R = \mathbb{E}[r(s, a)B(s, a)]$ from some reward observations, or directly set z_R if the reward is known (e.g., set $z_R = B(s, a)$ to reach a known target (s, a)). In the *exploitation phase*, we directly use the policy $\pi_{z_R}(s) = \arg \max_a F(s, a, z_R)^T z_R$.

The unsupervised learning phase. No rewards are used in this phase, and no family of tasks has to be specified manually. F and B are trained off-policy from observed transitions in the environment, to approximate the successor density: $F^{T}(s_0, a_0, z)B(s', a') \approx m^{\pi_z}(s_0, a_0, s', a')$ for every z. Training is based on the Bellman equation for the successor measure M^{π} ,

$$M^{\pi}(s_0, a_0, \{(s', a')\}) = \mathbb{1}_{s_0 = s', a_0 = a'} + \gamma \mathbb{E}_{s_1 \sim P(\mathrm{d}s_1 | s_0, a_0)} M^{\pi}(s_1, \pi(s_1), \{(s', a')\}).$$
(6)

We leverage a well-principled algorithm from [BTO21] in the single-policy setting: it learns the successor density m^{π} of a policy π without using the sparse reward $\mathbb{1}_{s_0=s', a_0=a'}$ (which would vanish in continuous spaces). This algorithm uses a parametric model $m_{\theta}^{\pi}(s_0, a_0, s', a')$. Given an observed transition (s_0, a_0, s_1) from the training set, generate an action $a_1 \sim \pi(a_1|s_1)$, and sample another state-action (s', a') from the training set, independently from (s_0, a_0, s_1) . Then update the

parameter θ by $\theta \leftarrow \theta + \eta \,\delta\theta$ with learning rate η and

$$\delta\theta := \partial_{\theta}m_{\theta}^{\pi}(s_0, a_0, s_0, a_0) + \partial_{\theta}m_{\theta}^{\pi}(s_0, a_0, s', a') \times (\gamma \, m_{\theta}^{\pi}(s_1, a_1, s', a') - m_{\theta}^{\pi}(s_0, a_0, s', a')) \tag{7}$$

This computes the density m^{π} of M^{π} with respect to the distribution ρ of state-actions in the training set. Namely, the true successor state density m^{π} is a fixed point of (7) in expectation [BTO21] (and is the only fixed point in the tabular or overparameterized case). Variants exist, such as using a target network for $m_{\theta}^{\pi}(s_1, a_1, s', a')$ on the right-hand side, as in DQN.

Thus, we first choose a parametric model F_{θ}, B_{θ} for the representations F and B, and set $m_{\theta}^{\pi_z}(s_0, a_0, s', a') := F_{\theta}(s_0, a_0, z)^{\top} B_{\theta}(s', a')$. Then we iterate the update (7) over many stateactions and values of z. This results in Algorithm 1. At each step, a value of z is picked at random, together with a batch of transitions (s_0, a_0, s_1) and a batch of state-actions (s', a') from the training set, with (s', a') independent from z and (s_0, a_0, s_1) .

For sampling z, we use a fixed distribution (rescaled Gaussians, see Appendix D). Any number of values of z may be sampled: this does not use up training samples. We use a target network with soft updates (Polyak averaging) as in DDPG. For training we also replace the greedy policy $\pi_z = \arg \max_a F(s, a, z)^{\top} z$ with a regularized version $\pi_z = \operatorname{softmax}(F(s, a, z)^{\top} z/\tau)$ with fixed temperature τ (Appendix D). Since there is unidentifiability between F and B (Appendix, Remark 4), we normalize B via an auxiliary loss in Algorithm 1].

For exploration in this phase, we use the policies being learned: the exploration policy chooses a random value of z from some distribution (e.g., Gaussian), and follows π_z for some time (Appendix, Algorithm []). However, the algorithm can also work from an existing dataset of off-policy transitions.

The reward estimation phase. Once rewards are available, we estimate a reward representation (policy parameter) z_R by weighing the representation *B* by the reward:

$$r_R := \mathbb{E}[r(s, a)B(s, a)] \tag{8}$$

where the expectation must be computed over the same distribution ρ of state-actions (s, a) used to learn F and B (see Appendix **B.5**] for using a different distribution). Thus, if the reward is black-box as in standard RL algorithms, then the exploration policy has to be run again for some time, and z_R is obtained by averaging r(s, a)B(s, a) over the states visited.

If the reward is known explicitly, this phase is unnecessary. For instance, if the reward is to reach a target state-action (s_0, a_0) while avoiding some forbidden state-actions $(s_1, a_1), ..., (s_k, a_k)$, one may directly set

$$z_R = B(s_0, a_0) - \lambda \sum B(s_i, a_i) \tag{9}$$

where the constant λ adjusts the negative reward for visiting a forbidden state. This can be used for goal-oriented RL.

If the reward is known algebraically as a function r(s, a), then z_R may be computed by averaging the function r(s, a)B(s, a) over a replay buffer from the unsupervised training phase. We may also use a reward model $\hat{r}(s, a)$ of r(s, a) trained on some reward observations from any source. An approximate value for z_R still provides an approximately optimal policy (Appendix, Prop. 7 and Thm. 9).

The exploitation phase. Once the reward representation z_R has been estimated, the Q-function is estimated as

$$Q(s,a) = F(s,a,z_R)^{\top} z_R.$$
(10)

The corresponding policy $\pi_{z_R}(s) = \arg \max_a Q(s, a)$ is used for exploitation.

Fine-tuning was not needed in our experiments, but it is possible to fine-tune the Q-function using actual rewards, by setting $Q(s, a) = F(s, a, z_R)^T z_R + q_\theta(s, a)$ where the fine-tuning model q_θ is initialized to 0 and learned via any standard Q-learning method.

Incorporating prior information on rewards in *B***.** Trying to plan in advance for all possible rewards in an arbitrary environment may be too generic and problem-agnostic, and become difficult in large environments, requiring long exploration and a large *d* to accommodate all rewards. In practice, we are often interested in rewards depending, not on the full state, but only on a part or some features of the state (e.g., a few components of the state, such as the position of an agent, or its neighbordhood, rather than the full environment).

If this is known in advance, the representation B can be trained on that part of the state only, with 254 the same theoretical guarantees (Appendix, Theorem 2). F still needs to use the full state as input. 255 This way, the FB model of the transition probabilities (2) only has to learn the future probabilities 256 of the part of interest in (s', a'), based on the full initial state (s_0, a_0) . Explicitly, if $\varphi: S \times A \to G$ 257 is a feature map to some features $g = \varphi(s, a)$, and if we know that the reward will be a function 258 R(g), then Theorem 1 still holds with B(g) everywhere instead of $B(\underline{s}, a)$, and with the successor 259 density $m^{\pi}(s_0, a_0, g)$ instead of $m^{\pi}(s_0, a_0, s', a')$ (Appendix, Theorem 2). Learning this m^{π} is done 260 by replacing $\partial_{\theta}m_{\theta}^{\pi}(s_0, a_0, s_0, a_0)$ with $\partial_{\theta}m_{\theta}^{\pi}(s_0, a_0, \varphi(s_0, a_0))$ in the first term in (7) [BTO21]. 261 Rewards can be arbitrary functions of g, so this is more general than [BBQ⁺18] which only considers 262 rewards linear in g. For instance, in MsPacman below, we let g be the 2D position (x, y) of the agent, 263 so we can optimize any reward function that depends on this position. 264

Limitations. First, this method does not solve exploration: it assumes access to a good exploration strategy. (Here we used the policies π_z with random values of z, corresponding to random rewards.)

Next, this task-agnostic approach is relevant if the reward is not known in advance, but may not bring the best performance on a particular reward. Mitigation strategies include: increasing d; using prior information on rewards by including relevant variables into B, as discussed above; and fine-tuning the Q-function at test time based on the initial $F^{\top}B$ estimate.

Indeed, as reward functions are represented by a *d*-dimensional vector $z_R = \mathbb{E}[r,B]$, some information about the reward is necessarily lost. Any reward uncorrelated to *B* is treated as 0. The necessary dimension *d* for good behavior may be large. Still, $d \approx 100$ worked in our experiments, and Appendix B.2 argues theoretically that d = O(n) is enough for navigation on an *n*-dimensional grid.

We expect this method to have an implicit bias for long-range behavior (spatially smooth rewards), while local details of the reward function may be blurred. Indeed, $F^{T}B$ is optimized to approximate the successor measure $M^{\pi} = \sum_{t} \gamma^{t} P_{\pi}^{t}$ with P_{π}^{t} the *t*-step transition kernel for each policy π . The rank-*d* approximation will favor large eigenvectors of P_{π} , i.e., small eigenvectors of the Markov chain Laplacian Id $-\gamma P_{\pi}$. These loosely correspond to long-range (low-frequency) behavior [MM07]: presumably, *F* and *B* will learn spatially smooth rewards first. Indeed, experimentally, a small *d* leads to spatial blurring of rewards and *Q*-functions (Fig. 3). Arguably, without any prior information this is a reasonable prior; [SBG17] have argued for the cognitive relevance of low-dimensional approximations of successor representations.

Variance is a potential issue in larger environments, although this did not arise in our experiments. 284 Learning m^{π} requires sampling a state-action (s_0, a_0) and an independent state-action (s', a'). In 285 large spaces, most state-action pairs will be unrelated. A possible mitigation is to use strategies such 286 as Hindsight Experience Replay $[ACR^+17]$ to select goals related to the current state-action. The 287 following may help a lot: the update of F and B decouples as an expectation over (s_0, a_0) , times an 288 expectation over (s', a'). Thus, by estimating these expectations by a moving average over a dataset, 289 it is easy to have many pairs (s_0, a_0) interact with many (s', a'). The cost is handling full $d \times d$ 290 matrices. This will be explored in future work. 291

292 5 Experiments

We first consider the task of reaching arbitrary goal states. For this, we can make quantitative comparisons to existing goal-oriented baselines. Next, we illustrate qualitatively some tasks that cannot be tackled a posteriori by goal-oriented methods, such as introducing forbidden states. Finally, we illustrate some of the representations learned.

297 5.1 Environments and Experimental Setup

We run our experiments on a selection of environments that are diverse in term of state space dimensionality, stochasticity and dynamics.

- Discrete Maze is the classical gridworld with four rooms. States are represented by one-hot unit vectors.
- Continuous Maze is a two dimensional environment with impassable walls. States are represented by their Cartesian coordinates $(x, y) \in [0, 1]^2$. The execution of one of the actions moves the agent in the desired direction, but with normal random noise added to the position of the agent.
- FetchReach is a variant of the simulated robotic arm environment from [PAR⁺18] using discrete actions instead of continuous actions. States are 10-dimensional vectors consisting of positions and velocities of robot joints.
- Ms. Pacman is a variant of the Atari 2600 game Ms. Pacman, where an episode ends when the agent is captured by a monster [RUMS18]. States are obtained by processing the raw visual input directly from the screen. Frames are preprocessed by cropping, conversion to grayscale and downsampling to 84 × 84 pixels. A state st is the concatenation of (xt-12, xt-8, xt-4, xt) frames, i.e. an 84 × 84 × 4 tensor. An action repeat of 12 is used. As Ms. Pacman is not originally a multi-goal domain, we define the goals as the 148 reachable coordinates (x, y) on the screen; these
- can be reached only by learning to avoid monsters.

For all environments, we run algorithms for 800 epochs. Each epoch consists of 25 cycles where we interleave between gathering some amount of transitions, to add to the replay buffer, and performing 40 steps of stochastic gradient descent on the model parameters. To collect transitions, we generate episodes using some behavior policy. For both mazes, we use a uniform policy while for FetchReach and Ms. Pacman, we use an ε -greedy policy with respect to the current approximation $F(s, a, z)^{\top} z$ for a sampled z. At evaluation time, ε -greedy policies are also used, with a smaller ε . More details are given in Appendix D

322 5.2 Goal-Oriented Setting: Quantitative Comparisons

We investigate the FB representation over goal-reaching tasks and compare it to goal-oriented baselines: DQN^I, and DQN with HER when needed. We define sparse reward functions. For Discrete Maze, the reward function is equal to one when the agent's state is equal exactly to the goal state. For Discrete Maze, we measured the quality of the obtained policy to be the ratio between the true expected discounted reward of the policy for its goal and the true optimal value function, on average

¹Here DQN is short for goal-oriented DQN, Q(s, a, g).

over all states. For the other environments, the reward function is equal to one when the distance of

the agent's position and the goal position is below some threshold, and zero otherwise. We assess policies by computing the average success rate, i.e the average number of times the agent successfully

reaches its goal.

332



Figure 1: Comparative performance of FB for different dimensions and DQN in the FetchReach. **Left**: success rate averaged over 20 randomly selected goals as function of the first 100 training epochs. **Right**: success rate averaged over 20 random goals after 800 training epochs.



Figure 2: Comparative performance of FB for different dimensions and DQN in Ms. Pacman. Left: success rate averaged over 20 randomly selected goals as function of the first 200 training epochs. **Right**: success rate averaged over the goal space after 800 training epochs.

Figs. 1 and 2 show the comparative performance of FB for different dimensions *d*, and DQN respectively in FetchReach and Ms. Pacman (similar results in Discrete and Continuous Mazes are provided in Appendix D). In Ms. Pacman, DQN totally fails to learn and we had to add HER to make it work. The performance of FB consistently increases with the dimension *d* and the best dimension matches the performance of the goal-oriented baseline.

In Discrete Maze, we observe a drop of performance for d = 25 (Appendix D, Fig. 8): this is due to

the spatial smoothing induced by the small rank approximation and the reward being nonzero only

if the agent is exactly at the goal. This spatial blurring is clear on heatmaps for d = 25 vs d = 75

(Fig. 3). With d = 25 the agent often stops right next to its goal.

To evaluate the sample efficiency of FB, after each epoch, we evaluate the agent on 20 randomly selected goals. Learning curves are reported in Figs. 1 and 2 (left). In all environments, we observe no loss in sample efficiency compared to the goal-oriented baseline. In Ms. Pacman, FB even learns faster than DQN+HER.

346 5.3 More Complex Rewards: Qualitative Results

We now investigate FB's ability to generalize to new tasks that cannot be solved by an already trained goal-oriented model: reaching a goal with forbidden states imposed a posteriori, reaching the nearest of two goals, and choosing between a small, close reward and a large, distant one.

First, for the task of reaching a target position $g_0 \nleftrightarrow$ while avoiding some forbidden positions $g_1, \ldots, g_k \bullet$, we set $z_R = B(g_1) - \lambda \sum_{i=1}^k B(g_i)$ and run the corresponding ε -greedy policy defined by $F(s, a, z_R)^\top z_R$. Fig. 5 shows the resulting trajectories, which succeed at solving the task for the different domains. In Ms. Pacman, the path is suboptimal (though successful) due to the sudden appearance of a monster along the optimal path. (We only plot the initial frame; see the full series of frames along the trajectory in Appendix D, Fig. 16) Fig. 4 (left) provides a contour plot of $\max_{a \in A} F(s, a, z_R)^\top z_R$ for the continuous maze and shows the landscape shape around the forbidden regions.

Next, we consider the task of reaching the closest target among two equally rewarding positions g_0 and g_1 , by setting $z_R = B(g_0) + B(g_1)$. The optimal Q-function is *not* a linear combination of the Q-functions for g_0 and g_1 . Fig. 6 shows successful trajectories generated by the policy π_{z_R} . On the contour plot of $\max_{a \in A} F(s, a, z_R)^{\top} z_R$ in Fig. 4 (right), the two rewarding positions appear as basins of attraction. Similar results for a third task are shown in Appendix D introducing a "distracting" small reward next to the initial position of the agent, with a larger reward further away.



Figure 7: Visualization of FB embedding vectors on Continuous Maze after projecting them in two-dimensional space with t-SNE. Left: the states to be mapped. Middle: the *F* embedding. Right: the *B* embedding. The walls appear as large dents; the smaller dents correspond to the number of steps needed to get past a wall.



Figure 3: Heatmap of $\max_a F(s, a, z_R)^\top z_R$ for $z_R = B(\bigstar)$ Left: d = 25. Right: d = 75.



Figure 5: Trajectories generated by the $F^{\top}B$ policies for the task of reaching a target position (star shape \bigstar while avoiding forbidden positions (red shape \bullet)

366 5.4 Embedding Visualizations



Figure 4: Contour plot of $\max_{a \in A} F(s, a, z_R)^{\top} z_R$ in Continuous Maze. Left: for the task of reaching a target while avoiding a forbidden region, **Right**: for two equally rewarding targets.



Figure 6: Trajectories generated by the $F^{\top}B$ policies for the task of reaching the closest among two equally rewarding positions (star shapes \bigstar). (Optimal *Q*-values are not linear over such mixtures.)

We visualize the learned FB state embeddings for Continuous Maze by projecting them into 2dimensional space using t-SNE [VdMH08] in Fig. 7. For the forward embeddings, we set z = 0corresponding to the uniform policy. We can see that FB partitions states according to the topology induced by the dynamics: states on opposite sides of walls are separated in the representation space and states on the same side lie together. Appendix D includes embedding visualizations for different z and for Discrete Maze and Ms. Pacman.

373 6 Conclusion

The FB representation is a learnable mathematical object that "summarizes" a reward-free MDP. It provides near-optimal policies for any reward specified a posteriori, without planning. It is learned from black-box reward-free interactions with the environment. In practice, this unsupervised method performs comparably to goal-oriented methods for reaching arbitrary goals, but is also able to tackle more complex rewards in real time. The representations learned encode the MDP dynamics and may have broader interest.

364

365

References

381 382 383	[ACR+17]	Marcin Andrychowicz, Dwight Crow, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob McGrew, Josh Tobin, Pieter Abbeel, and Wojciech Zaremba. Hindsight experience replay. In <i>NIPS</i> , 2017.
384 385 386	[ARO ⁺ 19]	Ankesh Anand, Evan Racah, Sherjil Ozair, Yoshua Bengio, Marc-Alexandre Côté, and R Devon Hjelm. Unsupervised state representation learning in atari. <i>arXiv preprint arXiv:1906.08226</i> , 2019.
387 388 389	[BBQ ⁺ 18]	Diana Borsa, André Barreto, John Quan, Daniel Mankowitz, Rémi Munos, Hado van Hasselt, David Silver, and Tom Schaul. Universal successor features approximators. <i>arXiv preprint arXiv:1812.07626</i> , 2018.
390 391 392	[BDM ⁺ 17]	André Barreto, Will Dabney, Rémi Munos, Jonathan J Hunt, Tom Schaul, David Silver, and Hado P van Hasselt. Successor features for transfer in reinforcement learning. In <i>NIPS</i> , 2017.
393 394 395	[BNVB13]	Marc G Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning environment: An evaluation platform for general agents. <i>Journal of Artificial Intelligence Research</i> , 47:253–279, 2013.
396	[Bog07]	Vladimir I Bogachev. Measure theory. Springer, 2007.
397 398	[Bré99]	Pierre Brémaud. Markov chains: Gibbs fields, Monte Carlo simulation, and queues, volume 31. 1999.
399 400	[BTO21]	Léonard Blier, Corentin Tallec, and Yann Ollivier. Learning successor states and goal- dependent values: A mathematical viewpoint. <i>arXiv preprint arXiv:2101.07123</i> , 2021.
401 402	[Day93]	Peter Dayan. Improving generalization for temporal difference learning: The successor representation. <i>Neural Computation</i> , 5(4):613–624, 1993.
403 404	[FD02]	David Foster and Peter Dayan. Structure in the space of value functions. <i>Machine Learning</i> , 49(2):325–346, 2002.
405 406 407	[GHB ⁺ 19]	Christopher Grimm, Irina Higgins, Andre Barreto, Denis Teplyashin, Markus Wulfmeier, Tim Hertweck, Raia Hadsell, and Satinder Singh. Disentangled cumulants help successor representations transfer to new tasks. <i>arXiv preprint arXiv:1911.10866</i> , 2019.
408 409	[GS97]	Charles Miller Grinstead and James Laurie Snell. <i>Introduction to probability</i> . American Mathematical Soc., 1997.
410 411 412	[HDB ⁺ 19]	Steven Hansen, Will Dabney, Andre Barreto, Tom Van de Wiele, David Warde-Farley, and Volodymyr Mnih. Fast task inference with variational intrinsic successor features. <i>arXiv preprint arXiv:1906.05030</i> , 2019.
413 414 415	[JKSY20]	Chi Jin, Akshay Krishnamurthy, Max Simchowitz, and Tiancheng Yu. Reward-free exploration for reinforcement learning. In <i>International Conference on Machine Learning</i> , pages 4870–4879. PMLR, 2020.
416	[KS60]	J. G. Kemeny and J. L. Snell. Finite Markov Chains. Van Nostrand, New York, 1960.
417 418 419	[KST ⁺ 18]	Nan Rosemary Ke, Amanpreet Singh, Ahmed Touati, Anirudh Goyal, Yoshua Ben- gio, Devi Parikh, and Dhruv Batra. Modeling the long term future in model-based reinforcement learning. In <i>International Conference on Learning Representations</i> , 2018.
420 421 422	[MM07]	Sridhar Mahadevan and Mauro Maggioni. Proto-value functions: A laplacian framework for learning representation and control in markov decision processes. <i>Journal of Machine Learning Research</i> , 8(10), 2007.
423 424 425 426	[PAR ⁺ 18]	Matthias Plappert, Marcin Andrychowicz, Alex Ray, Bob McGrew, Bowen Baker, Glenn Powell, Jonas Schneider, Josh Tobin, Maciek Chociej, Peter Welinder, et al. Multi-goal reinforcement learning: Challenging robotics environments and request for research. <i>arXiv preprint arXiv:1802.09464</i> , 2018.

[RUMS18] Paulo Rauber, Avinash Ummadisingu, Filipe Mutz, and Jürgen Schmidhuber. Hindsight 427 policy gradients. In International Conference on Learning Representations, 2018. 428 [SB18] Richard S Sutton and Andrew G Barto. Reinforcement learning: An introduction. MIT 429 press, 2018. 2nd edition. 430 [SBG17] Kimberly L Stachenfeld, Matthew M Botvinick, and Samuel J Gershman. The hip-431 pocampus as a predictive map. Nature neuroscience, 20(11):1643, 2017. 432 [SHGS15] Tom Schaul, Daniel Horgan, Karol Gregor, and David Silver. Universal value function 433 approximators. In Francis Bach and David Blei, editors, Proceedings of the 32nd 434 International Conference on Machine Learning, volume 37 of Proceedings of Machine 435 Learning Research, pages 1312–1320, Lille, France, 07–09 Jul 2015. PMLR. 436 [SMD⁺11] Richard S Sutton, Joseph Modayil, Michael Delp, Thomas Degris, Patrick M Pilarski, 437 Adam White, and Doina Precup. Horde: A scalable real-time architecture for learning 438 knowledge from unsupervised sensorimotor interaction. In The 10th International 439 Conference on Autonomous Agents and Multiagent Systems-Volume 2, pages 761–768, 440 2011. 441 [SPS99] Richard S Sutton, Doina Precup, and Satinder Singh. Between mdps and semi-mdps: A 442 framework for temporal abstraction in reinforcement learning. Artificial intelligence, 443 112(1-2):181-211, 1999. 444 [Tal17] Erik Talvitie. Self-correcting models for model-based reinforcement learning. In 445 Proceedings of the AAAI Conference on Artificial Intelligence, volume 31, 2017. 446 [VdMH08] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. Journal of 447 machine learning research, 9(11), 2008. 448 [ZSBB17] Jingwei Zhang, Jost Tobias Springenberg, Joschka Boedecker, and Wolfram Burgard. 449 Deep reinforcement learning with successor features for navigation across similar en-450 vironments. In 2017 IEEE/RSJ International Conference on Intelligent Robots and 451 Systems (IROS), pages 2371–2378. IEEE, 2017. 452

453 Checklist

454	1. For all authors
455 456	(a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
457	(b) Did you describe the limitations of your work? [Yes], Section 4
458	(c) Did you discuss any potential negative societal impacts of your work? [No] This is a
459 460	theoretical work on fully generic methods for reinforcement learning. The potential impact is the same as that of reinforcement learning in general.
461 462	(d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
463	2. If you are including theoretical results
464	(a) Did you state the full set of assumptions of all theoretical results? [Yes]
465	(b) Did you include complete proofs of all theoretical results? [Yes], Appendix \overline{C} .
466	3. If you ran experiments
467	(a) Did you include the code, data, and instructions needed to reproduce the main ex-
468	perimental results (either in the supplemental material or as a URL)? [Yes] in the
469	supplemental material.
470 471	(b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] Appendix D
472	(c) Did you report error bars (e.g., with respect to the random seed after running experi-
473	ments multiple times)? [Yes]
474 475	(d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [No]
476	4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets
477	(a) If your work uses existing assets, did you cite the creators? [Yes] Appendix D.1
478	(b) Did you mention the license of the assets? [No]
479	(c) Did you include any new assets either in the supplemental material or as a URL? [No]
480	(d) Did you discuss whether and how consent was obtained from people whose data you're
481	using/curating? [N/A]
482	(e) Did you discuss whether the data you are using/curating contains personally identifiable
483	
484	5. If you used crowdsourcing or conducted research with human subjects
485 486	(a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
487 488	(b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
489	(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]
	- F F Kaus combendation. [F art]