
The Value Equivalence Principle for Model-Based Reinforcement Learning

Christopher Grimm

Computer Science & Engineering
University of Michigan
crgrimm@umich.edu

André Barreto, Satinder Singh, David Silver

DeepMind
{andrebarreto,baveja,davidsilver}@google.com

Abstract

Learning models of the environment from data is often viewed as an essential component to building intelligent reinforcement learning (RL) agents. The common practice is to separate the learning of the model from its use, by constructing a model of the environment’s dynamics that correctly predicts the observed state transitions. In this paper we argue that the limited representational resources of model-based RL agents are better used to build models that are directly useful for value-based planning. As our main contribution, we introduce the principle of value equivalence: two models are value equivalent with respect to a set of functions and policies if they yield the same Bellman updates. We propose a formulation of the model learning problem based on the value equivalence principle and analyze how the set of feasible solutions is impacted by the choice of policies and functions. Specifically, we show that, as we augment the set of policies and functions considered, the class of value equivalent models shrinks, until eventually collapsing to a single point corresponding to a model that perfectly describes the environment. In many problems, directly modelling state-to-state transitions may be both difficult and unnecessary. By leveraging the value-equivalence principle one may find simpler models without compromising performance, saving computation and memory. We illustrate the benefits of value-equivalent model learning with experiments comparing it against more traditional counterparts like maximum likelihood estimation. More generally, we argue that the principle of value equivalence underlies a number of recent empirical successes in RL, such as Value Iteration Networks, the Predictron, Value Prediction Networks, TreeQN, and MuZero, and provides a first theoretical underpinning of those results.

1 Introduction

Reinforcement learning (RL) provides a conceptual framework to tackle a fundamental challenge in artificial intelligence: how to design agents that learn while interacting with the environment [36]. It has been argued that truly general agents should be able to learn a model of the environment that allows for fast re-planning and counterfactual reasoning [32]. Although this is not a particularly contentious statement, the question of *how* to learn such a model is far from being resolved. The common practice in model-based RL is to conceptually separate the learning of the model from its use. In this paper we argue that the limited representational capacity of model-based RL agents is better allocated if the future *use* of the model (*e.g.*, value-based planning) is also taken into account during its construction [22, 15, 13].

Our primary contribution is to formalize and analyze a clear principle that underlies this new approach to model-based RL. Specifically, we show that, when the model is to be used for value-based planning, requirements on the model can be naturally captured by an equivalence relation induced by a set

of policies and functions. This leads to the *principle of value equivalence*: two models are value equivalent with respect to a set of functions and a set of policies if they yield the same updates under corresponding Bellman operators. The policies and functions then become the mechanism through which one incorporates information about the intended use of the model during its construction. We propose a formulation of the model learning problem based on the value equivalence principle and analyze how the set of feasible solutions is impacted by the choice of policies and functions. Specifically, we show that, as we augment the set of policies and functions considered, the class of value equivalent models shrinks, until eventually collapsing to a single point corresponding to a model that perfectly describes the environment.

We also discuss cases in which one can meaningfully restrict the class of policies and functions used to tailor the model. One common case is when the construction of an optimal policy through value-based planning only requires that a model predicts a subset of value functions. We show that in this case the resulting value equivalent models can perform well under much more restrictive conditions than their traditional counterparts. Another common case is when the agent has limited representational capacity. We show that in this scenario it suffices for a model to be value equivalent with respect to appropriately-defined bases of the spaces of representable policies and functions. This allows models to be found with less memory or computation than conventional model-based approaches that aim at predicting all state transitions, such as maximum likelihood estimation. We illustrate the benefits of value-equivalent model learning in experiments that compare it against more conventional counterparts. More generally, we argue that the principle of value equivalence underlies a number of recent empirical successes in RL and provides a first theoretical underpinning of those results [40, 34, 26, 16, 33].

2 Background

As usual, we will model the agent’s interaction with the environment using a *Markov Decision Process* (MDP) $\mathcal{M} \equiv \langle \mathcal{S}, \mathcal{A}, r, p, \gamma \rangle$ where \mathcal{S} is the state space, \mathcal{A} is the action space, $r(s, a, s')$ is the reward associated with a transition to state s' following the execution of action a in state s , $p(s'|s, a)$ is the transition kernel and $\gamma \in [0, 1)$ is a discount factor [30]. For convenience we also define $r(s, a) = \mathbb{E}_{S' \sim p(\cdot|s, a)}[r(s, a, S')]$.

A *policy* is a mapping $\pi : \mathcal{S} \mapsto \mathcal{P}(\mathcal{A})$, where $\mathcal{P}(\mathcal{A})$ is the space of probability distributions over \mathcal{A} . We define $\mathbb{I} \equiv \{\pi \mid \pi : \mathcal{S} \mapsto \mathcal{P}(\mathcal{A})\}$ as the set of all possible policies. The agent’s goal is to find a policy $\pi \in \mathbb{I}$ that maximizes the *value* of every state, defined as

$$v_\pi(s) \equiv \mathbb{E}_\pi \left[\sum_{i=0}^{\infty} \gamma^i r(S_{t+i}, A_{t+i}) \mid S_t = s \right], \quad (1)$$

where S_t and A_t are random variables indicating the state occupied and the action selected by the agent at time step t and $\mathbb{E}_\pi[\cdot]$ denotes expectation over the trajectories induced by π .

Many methods are available to carry out the search for a good policy [36, 39]. Typically, a crucial step in these methods is the computation of the value function of candidate policies—a process usually referred to as *policy evaluation*. One way to evaluate a policy π is through its *Bellman operator*:

$$\mathcal{T}_\pi[v](s) \equiv \mathbb{E}_{A \sim \pi(\cdot|s), S' \sim p(\cdot|s, A)} [r(s, A) + \gamma v(S')], \quad (2)$$

where v is any function in the space $\mathbb{V} \equiv \{f \mid f : \mathcal{S} \mapsto \mathbb{R}\}$. It is known that $\lim_{n \rightarrow \infty} (\mathcal{T}_\pi)^n v = v_\pi$, that is, starting from any $v \in \mathbb{V}$, the repeated application of \mathcal{T}_π will eventually converge to v_π [30].

In RL it is generally assumed that the agent does not know p and r , and thus cannot directly compute (2). In *model-free* RL this is resolved by replacing v_π with an action-value function and estimating the expectation on the right-hand-side of (2) through sampling [35]. In *model-based* RL, the focus of this paper, the agent learns approximations $\tilde{r} \approx r$ and $\tilde{p} \approx p$ and use them to compute (2) with p and r replaced by \tilde{p} and \tilde{r} [36].

3 Value equivalence

Given a state space \mathcal{S} and an action space \mathcal{A} , we call the tuple $m \equiv (r, p)$ a *model*. Note that a model plus a discount factor γ induces a Bellman operator (2) for every policy $\pi \in \mathbb{I}$. In this paper we

are interested in computing an approximate model $\tilde{m} = (\tilde{r}, \tilde{p})$ such that the induced operators $\tilde{\mathcal{T}}_\pi$, defined analogously to (2), are good approximations of the true \mathcal{T}_π . Our main argument is that models should only be distinguished with respect to the policies and functions they will actually be applied to. This leads to the following definition:

Definition 1 (Value equivalence). *Let $\Pi \subseteq \mathbb{I}$ be a set of policies and let $\mathcal{V} \subseteq \mathbb{V}$ be a set of functions. We say that models m and \tilde{m} are value equivalent with respect to Π and \mathcal{V} if and only if*

$$\mathcal{T}_\pi v = \tilde{\mathcal{T}}_\pi v \text{ for all } \pi \in \Pi \text{ and all } v \in \mathcal{V},$$

where \mathcal{T}_π and $\tilde{\mathcal{T}}_\pi$ are the Bellman operators induced by m and \tilde{m} , respectively.

Two models are value equivalent with respect to Π and \mathcal{V} if the effect of the Bellman operator induced by any policy $\pi \in \Pi$ on any function $v \in \mathcal{V}$ is the same for both models. Thus, if we are only interested in Π and \mathcal{V} , value-equivalent models are functionally identical. This can be seen as an equivalence relation that partitions the space of models conditioned on Π and \mathcal{V} :

Definition 2 (Space of value-equivalent models). *Let Π and \mathcal{V} be defined as above and let \mathcal{M} be a space of models. Given a model m , the space of value-equivalent models $\mathcal{M}_m(\Pi, \mathcal{V}) \subseteq \mathcal{M}$ is the set of all models $\tilde{m} \in \mathcal{M}$ that are value equivalent to m with respect to Π and \mathcal{V} .*

Let \mathbb{M} be a space of models containing at least one model m^* which perfectly describes the interaction of the agent with the environment. More formally, m^* induces the true Bellman operators \mathcal{T}_π defined in (2). Given a space of models $\mathcal{M} \subseteq \mathbb{M}$, often one is interested in models $m \in \mathcal{M}$ that are value equivalent to m^* . We will thus simplify the notation by defining $\mathcal{M}(\Pi, \mathcal{V}) \equiv \mathcal{M}_{m^*}(\Pi, \mathcal{V})$.

3.1 The topology of the space of value-equivalent models

The space $\mathcal{M}(\Pi, \mathcal{V})$ contains all the models in \mathcal{M} that are value equivalent to the true model m^* with respect to Π and \mathcal{V} . Since any two models $m, m' \in \mathcal{M}(\Pi, \mathcal{V})$ are equally suitable for value-based planning using Π and \mathcal{V} , we are free to use other criteria to choose between them. For example, if m is much simpler to represent or learn than m' , it can be preferred without compromises.

Clearly, the principle of value equivalence can be useful if leveraged in the appropriate way. In order for that to happen, it is important to understand the space of value-equivalent models $\mathcal{M}(\Pi, \mathcal{V})$. We now provide intuition for this space by analyzing some of its core properties. We refer the reader to Figure 1 for an illustration of the concepts to be discussed in this section. We start with a property that follows directly from Definitions 1 and 2:

Property 1. *Given $\mathcal{M}' \subseteq \mathcal{M}$, we have that $\mathcal{M}'(\Pi, \mathcal{V}) \subseteq \mathcal{M}(\Pi, \mathcal{V})$.*

The proofs of all theoretical results are in Appendix A.1. Property 1 states that, given a set of policies Π and a set of functions \mathcal{V} , reducing the size of the space of models \mathcal{M} also reduces the space of value-equivalent models $\mathcal{M}(\Pi, \mathcal{V})$. One immediate consequence of this property is that, if we consider the space of all policies \mathbb{I} and the space of all functions \mathbb{V} , we have one of two possibilities: either we end up with a perfect model or we end up with no model at all. Or, more formally:

Property 2. *$\mathcal{M}(\mathbb{I}, \mathbb{V})$ either contains m^* or is the empty set.*

Property 1 describes what happens to $\mathcal{M}(\Pi, \mathcal{V})$ when we vary \mathcal{M} with fixed Π and \mathcal{V} . It is also interesting to ask what happens when we fix the former and vary the latter. This leads to the next property:

Property 3. *Given $\Pi' \subseteq \Pi$ and $\mathcal{V}' \subseteq \mathcal{V}$, we have that $\mathcal{M}(\Pi, \mathcal{V}) \subseteq \mathcal{M}(\Pi', \mathcal{V}')$.*

According to Property 3, as we increase the size of Π or \mathcal{V} the size of $\mathcal{M}(\Pi, \mathcal{V})$ decreases. Although this makes intuitive sense, it is reassuring to know that value equivalence is a sound principle for

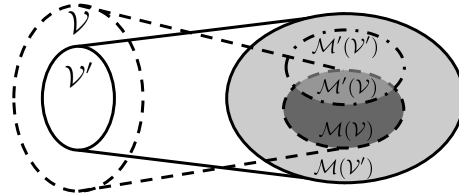


Figure 1: Understanding the space of value-equivalent models for a fixed Π , $\mathcal{M}' \subseteq \mathcal{M}$ and $\mathcal{V}' \subseteq \mathcal{V}$. Denote $\mathcal{M}(\mathcal{V}) \equiv \mathcal{M}(\mathcal{V}, \Pi)$. **Property 1:** $\mathcal{M}'(\mathcal{V}) \subseteq \mathcal{M}(\mathcal{V})$ and $\mathcal{M}'(\mathcal{V}') \subseteq \mathcal{M}(\mathcal{V}')$. **Property 3:** $\mathcal{M}(\mathcal{V}) \subseteq \mathcal{M}(\mathcal{V}')$ and $\mathcal{M}'(\mathcal{V}) \subseteq \mathcal{M}'(\mathcal{V}')$. **Property 4:** if $m^* \in \mathcal{M}$, then $m^* \in \mathcal{M}(\mathcal{V})$.

model selection, since by adding more policies to Π or more values to \mathcal{V} we can progressively restrict the set of feasible solutions. Thus, if \mathcal{M} contains the true model, we eventually pin it down. Indeed, in this case the true model belongs to *all* spaces of value equivalent models, as formalized below:

Property 4. *If $m^* \in \mathcal{M}$, then $m^* \in \mathcal{M}(\Pi, \mathcal{V})$ for all Π and all \mathcal{V} .*

3.2 A basis for the space of value-equivalent models

As discussed, it is possible to use the sets Π and \mathcal{V} to control the size of $\mathcal{M}(\Pi, \mathcal{V})$. But what exactly is the effect of Π and \mathcal{V} on $\mathcal{M}(\Pi, \mathcal{V})$? How much does $\mathcal{M}(\Pi, \mathcal{V})$ decrease in size when we, say, add one function to \mathcal{V} ? In this section we address this and similar questions.

We start by showing that, whenever a model is value equivalent to m^* with respect to discrete Π and \mathcal{V} , it is automatically value equivalent to m^* with respect to much larger sets. In order to state this fact more concretely we will need two definitions. Given a discrete set \mathcal{H} , we define $\text{span}(\mathcal{H})$ as the set formed by all linear combinations of the elements in \mathcal{H} . Similarly, given a discrete set \mathcal{H} in which each element is a function defined over a domain \mathcal{X} , we define the *pointwise span* of \mathcal{H} as

$$p\text{-span}(\mathcal{H}) \equiv \left\{ h : h(x) = \sum_i \alpha_{xi} h_i(x) \right\}, \text{ with } \alpha_{xi} \in \mathbb{R} \text{ for all } x \in \mathcal{X}, i \in \{1, \dots, |\mathcal{H}|\} \quad (3)$$

where $h_i \in \mathcal{H}$. Pointwise span can alternatively be characterized by considering each element in the domain separately: $g \in p\text{-span}(\mathcal{H}) \iff g(x) \in \text{span}\{h(x) : h \in \mathcal{H}\}$ for all $x \in \mathcal{X}$. Equipped with these concepts we present the following result:

Proposition 1. *For discrete Π and \mathcal{V} , we have that $\mathcal{M}(\Pi, \mathcal{V}) = \mathcal{M}(p\text{-span}(\Pi) \cap \Pi, \text{span}(\mathcal{V}))$.*

Proposition 1 provides one possible answer to the question posed at the beginning of this section: the contraction of $\mathcal{M}(\Pi, \mathcal{V})$ resulting from the addition of one policy to Π or one function to \mathcal{V} depends on their effect on $p\text{-span}(\Pi)$ and $\text{span}(\mathcal{V})$. For instance, if a function v can be obtained as a linear combination of the functions in \mathcal{V} , adding it to this set will have no effect on the space of equivalent models $\mathcal{M}(\Pi, \mathcal{V})$. More generally, Proposition 1 suggests a strategy to *build* the set \mathcal{V} : one should find a set of functions that form a basis for the space of interest. When \mathcal{S} is finite, for example, having \mathcal{V} be a basis for $\mathbb{R}^{|\mathcal{S}|}$ means that the value equivalence principle will apply to every function $v \in \mathbb{R}^{|\mathcal{S}|}$. The same reasoning applies to Π . In fact, because $p\text{-span}(\Pi)$ grows independently pointwise, it is relatively simple to build a set Π that covers the space of policies one is interested in. In particular, when \mathcal{A} is finite, it is easy to define a set Π for which $p\text{-span}(\Pi) \supseteq \Pi$: it suffices to have for every state-action pair $(s, a) \in \mathcal{S} \times \mathcal{A}$ at least one policy $\pi \in \Pi$ such that $\pi(a|s) = 1$. This means that we can apply the value equivalence principle to the entire set Π using $|\mathcal{A}|$ policies only.

Combining Proposition 1 and Property 2 we see that by defining Π and \mathcal{V} appropriately we can focus on the subset of \mathcal{M} whose models perfectly describe the environment:

Remark 1. *If $\Pi \subseteq p\text{-span}(\Pi)$ and $\mathcal{V} = \text{span}(\mathcal{V})$, then $\mathcal{M}(\Pi, \mathcal{V}) = m^*$ or $\mathcal{M}(\Pi, \mathcal{V}) = \emptyset$.*

We have shown how Π and \mathcal{V} have an impact on the number of value equivalent models in $\mathcal{M}(\Pi, \mathcal{V})$; to make the discussion more concrete, we now focus on a specific model space \mathcal{M} and analyze the rate at which this space shrinks as we add more elements to Π and \mathcal{V} . Before proceeding we define a set of functions \mathcal{H} as *pointwise linearly independent* if $h \notin p\text{-span}(\mathcal{H} \setminus \{h\})$ for all $h \in \mathcal{H}$.

Suppose both \mathcal{S} and \mathcal{A} are finite. In this case a model can be defined as $m = (\mathbf{r}, \mathbf{P})$, where $\mathbf{r} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ and $\mathbf{P} \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}| \times |\mathcal{A}|}$. A policy can then be thought of as a vector $\pi \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$. We denote the set of all transition matrices induced by transition kernels as \mathbb{P} . To simplify the analysis we will consider that \mathbf{r} is known and we are interested in finding a model $\hat{\mathbf{P}} \in \mathbb{P}$. In this setting, we write $\mathbb{P}(\Pi, \mathcal{V})$ to denote the set of transition matrices that are value equivalent to the true transition matrix \mathbf{P}^* . We define the dimension of a set \mathcal{X} as the lowest possible Hamel dimension of a vector-space enclosing some translated version of it: $\dim[\mathcal{X}] = \min_{\mathcal{W}, c \in \mathcal{W}(\mathcal{X})} \mathcal{H}\text{-dim}[\mathcal{W}]$ where $\mathcal{W}(\mathcal{X}) = \{(\mathcal{W}, c) : \mathcal{X} + c \subseteq \mathcal{W}\}$, \mathcal{W} is a vector-space, c is an offset and $\mathcal{H}\text{-dim}[\cdot]$ denotes the Hamel dimension. Recall that the Hamel dimension of a vector-space is the size of the smallest set of mutually linearly independent vectors that spans the space (this corresponds to the usual notion of dimension, that is, the minimal number of coordinates required to uniquely specify each point). So, under no restrictions imposed by Π and \mathcal{V} , we have that $\dim[\mathbb{P}] = (|\mathcal{S}| - 1)|\mathcal{S}||\mathcal{A}|$. We now show how fast the size of $\mathbb{P}(\Pi, \mathcal{V})$ decreases as we extend the ranges of Π and \mathcal{V} :

Proposition 2. Let Π be a set of m pointwise linearly independent policies $\pi_i \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ and let \mathcal{V} be a set of k linearly independent vectors $v_i \in \mathbb{R}^{|\mathcal{S}|}$. Then,

$$\dim[\mathbb{P}(\Pi, \mathcal{V})] \leq |\mathcal{S}| (|\mathcal{S}||\mathcal{A}| - mk).$$

Interestingly, Proposition 2 shows that the elements of Π and \mathcal{V} interact in a multiplicative way: when there are m pointwise linearly independent policies, enlarging \mathcal{V} with a single function v that is linearly independent of its counterparts will decrease the bound on the dimension of $\mathbb{P}(\Pi, \mathcal{V})$ by a factor of m . This makes intuitive sense if we note that by definition $m \leq |\mathcal{A}|$: for an expressive enough Π , each $v \in \mathcal{V}$ will provide information about the effect of all actions in $a \in \mathcal{A}$. Conversely, because $\text{span}(\mathcal{V}) = k \leq |\mathcal{S}|$, we can only go so far in pinning down the model when $m < |\mathcal{A}|$ —which also makes sense, since in this case we cannot possibly know about the effect of all actions, no matter how big \mathcal{V} is. Note that when $m = |\mathcal{A}|$ and $k = |\mathcal{S}|$ the space $\mathbb{P}(\Pi, \mathcal{V})$ reduces to $\{P^*\}$.

4 Model learning based on the value-equivalence principle

We now discuss how the principle of value equivalence can be incorporated into model-based RL. Often in model-based RL one learns a model $\tilde{m} = (\tilde{r}, \tilde{p})$ without taking the space $\mathcal{M}(\Pi, \mathcal{V})$ into account. The usual practice is to cast the approximations $\tilde{r} \approx r$ and $\tilde{p} \approx p$ as optimization problems over a model-space \mathcal{M} that do not involve the sets Π and \mathcal{V} . Given a space \mathcal{R} of possible approximations \tilde{r} , we can formulate the approximation of the rewards as $\text{argmin}_{\tilde{r} \in \mathcal{R}} \ell_r(r, \tilde{r})$, where ℓ_r is a loss function that measures the dissimilarity between r and \tilde{r} . The approximation of the transition dynamics can be formalized in an analogous way: $\text{argmin}_{\tilde{p} \in \mathcal{P}} \ell_p(p, \tilde{p})$, where \mathcal{P} is the space of possible approximations \tilde{p} .

A common choice for ℓ_r is

$$\ell_{r, \mathcal{D}}(r, \tilde{r}) \equiv \mathbb{E}_{(S, A) \sim \mathcal{D}} [(r(S, A) - \tilde{r}(S, A))^2], \quad (4)$$

where \mathcal{D} is a distribution over $\mathcal{S} \times \mathcal{A}$. The loss ℓ_p is usually defined based on the principle of *maximum likelihood estimation* (MLE):

$$\ell_{p, \mathcal{D}}(p, \tilde{p}) \equiv \mathbb{E}_{(S, A) \sim \mathcal{D}} [\text{D}_{\text{KL}}(p(\cdot | S, A) || \tilde{p}(\cdot | S, A))], \quad (5)$$

where D_{KL} is the Kullback-Leibler (KL) divergence. Since we normally do not have access to r and p , the losses (4) and (5) are usually minimized using transitions sampled from the environment [38]. There exist several other criteria to approximate p based on state transitions, such as maximum *a posteriori* estimation, maximum entropy estimation, and Bayesian posterior inference [13]. Although we focus on MLE for simplicity, our arguments should extend to these other criteria as well.

Both (4) and (5) have desirable properties that justify their widespread adoption [24]. However, we argue that ignoring the future use of \tilde{r} and \tilde{p} may not always be the best choice [22, 15]. To illustrate this point, we now show that, by doing so, one might end up with an approximate model when an exact one were possible. Let $\mathcal{P}(\Pi, \mathcal{V})$ be the set of value equivalent transition kernels in \mathcal{P} . Then,

Proposition 3. *The maximum-likelihood estimate of p^* in \mathcal{P} may not belong to a $\mathcal{P}(\Pi, \mathcal{V}) \neq \emptyset$.*

Proposition 3 states that, even when there exist models in \mathcal{P} that are value equivalent to p^* with respect to Π and \mathcal{V} , the minimizer of (5) may not be in $\mathcal{P}(\Pi, \mathcal{V})$. In other words, even when it is possible to perfectly handle the policies in Π and the values in \mathcal{V} , the model that achieves the smallest MLE loss will do so only approximately. This is unsurprising since the loss (5) is agnostic of Π and \mathcal{V} , providing instead a model that represents a compromise across all policies Π and all functions \mathcal{V} .

We now define a value-equivalence loss that explicitly takes into account the sets Π and \mathcal{V} :

$$\ell_{\Pi, \mathcal{V}}(m^*, \tilde{m}) \equiv \sum_{\pi \in \Pi} \sum_{v \in \mathcal{V}} \|\mathcal{T}_\pi v - \tilde{\mathcal{T}}_\pi v\|, \quad (6)$$

where $\tilde{\mathcal{T}}_\pi$ are Bellman operators induced by \tilde{m} and $\|\cdot\|$ is a norm. Given (6), the problem of learning a model based on the value equivalence principle can be formulated as $\text{argmin}_{\tilde{m} \in \mathcal{M}} \ell_{\Pi, \mathcal{V}}(m^*, \tilde{m})$.

As noted above, we usually do not have access to \mathcal{T}_π , and thus the loss (6) will normally be minimized based on sample transitions. Let $\mathcal{S}_\pi \equiv \{(s_i^\pi, a_i^\pi, r_i^\pi, \hat{s}_i^\pi) | i = 1, 2, \dots, n^\pi\}$ be n^π sample transitions associated with policy $\pi \in \Pi$. We assume that the initial states s_i^π were sampled according to some

distribution \mathcal{D}' over \mathcal{S} and the actions were sampled according to the policy π , $a_i^\pi \sim \pi(\cdot | s_i^\pi)$ (note that \mathcal{D}' can be the distribution resulting from a direct interaction of the agent with the environment). When $\|\cdot\|$ appearing in (6) is a p -norm, we can write its empirical version as

$$\ell_{\Pi, \mathcal{V}, \mathcal{D}'}(m^*, \tilde{m}) \equiv \sum_{\pi \in \Pi} \sum_{v \in \mathcal{V}} \sum_{s \in \mathcal{S}'_\pi} \left[\frac{\sum_{i=1}^{n_\pi} \mathbb{1}\{s_i^\pi = s\} (r_i^\pi + \gamma v(\hat{s}_i^\pi))}{\sum_{i=1}^{n_\pi} \mathbb{1}\{s_i^\pi = s\}} - \tilde{\mathcal{T}}_\pi v[s] \right]^p, \quad (7)$$

where \mathcal{S}'_π is a set containing only the initial states $s_i^\pi \in \mathcal{S}_\pi$ and $\mathbb{1}\{\cdot\}$ is the indicator function. We argue that, when we know policies Π and functions \mathcal{V} that are sufficient for planning, the appropriate goal for model-learning is to minimize the value-equivalence loss (6). As shown in Proposition 3, the model \tilde{m} that minimizes (4) and (5) may not achieve zero loss on (6) even when such a model exists in \mathcal{M} . In general, though, we should not expect there to be a model $\tilde{m} \in \mathcal{M}$ that leads to zero value-equivalence loss. Even then, value equivalence may lead to a better model than conventional counterparts (see Figure 2 for intuition and Appendix A.1.2 for a concrete example).

4.1 Restricting the sets of policies and functions

The main argument of this paper is that, rather than learning a model that suits all policies Π and all functions \mathcal{V} , we should instead focus on the sets of policies Π and functions \mathcal{V} that are necessary for planning. But how can we know these sets *a priori*? We now show that it is possible to exploit structure on both the *problem* and the *solution* sides.

First, we consider structure in the *problem*. Suppose we had access to the true model m^* . Then, given an initial function v , a value-based planning algorithm that makes use of m^* will generate a sequence of functions $\vec{V}_v \equiv \{v_1, v_2, \dots\}$ [10]. Clearly, if we replace m^* with any model in $\mathcal{M}(\Pi, \vec{V}_v)$, the behavior of the algorithm starting from v remains unaltered. This allows us to state the following:

Proposition 4. *Suppose $v \in \mathcal{V}' \implies \mathcal{T}_\pi v \in \mathcal{V}'$ for all $\pi \in \Pi$. Let $p\text{-span}(\Pi) \supseteq \Pi$ and $\text{span}(\mathcal{V}) = \mathcal{V}'$. Then, starting from any $v' \in \mathcal{V}'$, any $\tilde{m} \in \mathcal{M}(\Pi, \mathcal{V})$ yields the same solution as m^* .*

Because \mathcal{T}_π are contraction mappings, it is always possible to define a $\mathcal{V}' \subset \mathcal{V}$ such that the condition of Proposition 4 holds: we only need to make \mathcal{V}' sufficiently large to encompass v and the operators' fixed points. But in some cases there exist more structured \mathcal{V}' : in Appendix A.1 we give an example of a finite state-space MDP in which a sequence $v_1, v_2 = \mathcal{T}_\pi v_1, v_3 = \mathcal{T}_{\pi'} v_2, \dots$ that reaches a specific k -dimensional subspace of $\mathbb{R}^{|\mathcal{S}|}$ stays there forever. The value equivalence principle provides a mechanism to exploit this type of structure, while conventional model-learning approaches, like MLE, are oblivious to this fact. Although in general we do not have access to \mathcal{V}' , in some cases this set will be revealed through the very process of enforcing value equivalence. For example, if \tilde{m} is being learned online based on a sequence $v_1, v_2 = \mathcal{T}_\pi v_1, v_3 = \mathcal{T}_{\pi'} v_2, \dots$, as long as the sequence reaches a $v_i \in \mathcal{V}'$ we should expect \tilde{m} to eventually specialize to \mathcal{V}' [13, 33].

Another possibility is to exploit geometric properties of the *value functions* \vec{V}_v . It is known that the set of all value functions of a given MDP forms a polytope $\check{\mathcal{V}} \subset \mathcal{V}$ [11]. Even though the sequence \vec{V}_v an algorithm generates may not be strictly inside the polytope $\check{\mathcal{V}}$, this set can still serve as a reference in the definition of \mathcal{V} . For example, based on Proposition 1, we may want to define a \mathcal{V} that spans as much of the polytope $\check{\mathcal{V}}$ as possible [5]. This suggests that the functions in \mathcal{V} should be actual value functions v_π associated with policies $\pi \in \Pi$. In Section 5 we show experiments that explore this idea.

We now consider structure in the *solution*. Most large-scale applications of model-based RL use function approximation. Suppose the agent can only represent policies $\pi \in \tilde{\Pi}$ and value functions $v \in \tilde{\mathcal{V}}$. Then, a value equivalent model $\tilde{m} \in \mathcal{M}(\tilde{\Pi}, \tilde{\mathcal{V}})$ is as good as any model. To build intuition, suppose the agent uses state aggregation to approximate the value function. In this case two models

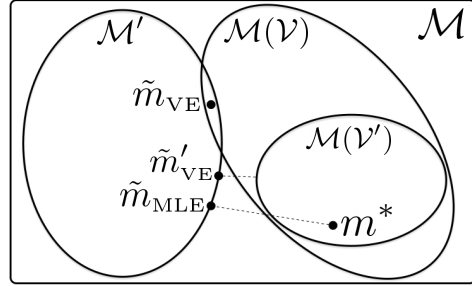


Figure 2: When the hypothesis space \mathcal{M}' and the space of value-equivalent models $\mathcal{M}(\mathcal{V})$ intersect, the resulting model \tilde{m}_{VE} has zero loss (6), while the corresponding MLE model \tilde{m}_{MLE} may not (Proposition 3). But even when \mathcal{M}' and $\mathcal{M}(\mathcal{V}')$ do not intersect, the resulting \tilde{m}'_{VE} can outperform \tilde{m}_{MLE} with the appropriate choices of Π and \mathcal{V} , as we illustrate in the experiments of Section 5.

with the same transition probabilities between clusters of states are indistinguishable from the agent’s point of view. It thus makes sense to build \mathcal{V} using piecewise-constant functions that belong to the space of function representable by the agent, $v \in \tilde{\mathcal{V}}$. The following remark generalises this intuition:

Remark 2. *Suppose the agent represents the value function using a linear function approximation: $\tilde{\mathcal{V}} = \{\tilde{v} \mid \tilde{v}(s) = \sum_{i=1}^d \phi_i(s)w_i\}$, where $\phi_i : \mathcal{S} \mapsto \mathbb{R}$ are fixed features and $w \in \mathbb{R}^d$ are learnable parameters. In addition, suppose the agent can only represent policies $\pi \in \tilde{\Pi}$. Then, Proposition 1 implies that if we use the features themselves as the functions adopted with value equivalence, $\mathcal{V} = \{\phi_i\}_{i=1}^d$, we have that $\mathcal{M}(\tilde{\Pi}, \{\phi_i\}_{i=1}^d) = \mathcal{M}(\tilde{\Pi}, \tilde{\mathcal{V}})$. In other words, models that are value equivalent with respect to the features are indiscernible to the agent.*

According to the remark above, when using linear function approximation, a model that is value equivalent with respect to the approximator’s features will perform no worse than any other model. This prescribes a concrete way to leverage the value equivalence principle in practice, since the set of functions \mathcal{V} is automatically defined by the choice of function approximator. Note that, although the remark is specific to linear value function approximation, it applies equally to linear and non-linear models (this is in contrast with previous work showing the equivalence between model-free RL using linear function approximation and model-based RL with a linear model for expected features [27, 38]). The principle of finding a basis for $\tilde{\mathcal{V}}$ also extends to non-linear value function approximation, though in this case it is less clear how to define a set \mathcal{V} that spans $\tilde{\mathcal{V}}$. One strategy is to *sample* the functions to be included in \mathcal{V} from the set $\tilde{\mathcal{V}}$ of (non-linear) functions the agent can represent. Despite its simplicity, this strategy can lead to good performance in practice, as we show next.

5 Experiments

We now present experiments illustrating the usefulness of the value equivalence principle in practice. Specifically, we compare models computed based on value equivalence (VE) with models resulting from maximum likelihood estimation (MLE). All our experiments followed the same protocol: (i) we collected sample transitions from the environment using a policy that picks actions uniformly at random, (ii) we used this data to learn an approximation \tilde{r} using (4) as well as approximations \tilde{p} using either MLE (5) or VE (7), (iii) we learned a policy $\tilde{\pi}$ based on $\tilde{m} = (\tilde{r}, \tilde{p})$, and (iv) we evaluated $\tilde{\pi}$ on the actual environment. The specific way each step was carried out varied according to the characteristics of the environment and function approximation used; see App. A.2 for details.

One of the central arguments of this paper is that the value equivalence principle can yield a better allocation of the limited resources of model-based agents. In order to verify this claim, we varied the representational capacity of the agent’s models \tilde{m} and assessed how well MLE and VE performed under different constraints. As discussed, VE requires the definition of two sets: Π and \mathcal{V} . It is usually easy to define a set of policies Π such that $p\text{-span}(\Pi) \supseteq \mathbb{P}$; since all the environments used in our experiments have a finite action space \mathcal{A} , we accomplished that by defining $\Pi = \{\pi^a\}_{a \in \mathcal{A}}$ where $\pi^a(a|s) = 1$ for all $s \in \mathcal{S}$. We will thus restrict our attention to the impact of the set of functions \mathcal{V} .

As discussed, one possible strategy to define \mathcal{V} is to use actual value functions in an attempt to span as much as possible of the value polytope $\tilde{\mathcal{V}}$ [5]. Figure 3 shows results of VE when using this strategy. Specifically, we compare VE’s performance with MLE’s on two well known domains: “four rooms” [37] and “catch” [25]. For each domain, we show two types of results: we either fix the capacity of the model \tilde{p} and vary the size of \mathcal{V} or vice-versa (in the Appendix we show results with all possible combinations of model capacities and sizes of \mathcal{V}). Note how the models produced by VE outperform MLE’s counterparts across all scenarios, and especially so under stricter restrictions on the model. This corroborates our hypothesis that VE yields models that are tailored to future use.

Another strategy to define \mathcal{V} is to use functions from $\tilde{\mathcal{V}}$, the space of functions representable by the agent, in order to capture as much as possible of this space. In Figure 4 we compare VE using this strategy with MLE. Here we use as domains catch and “cart-pole” [4] (but see Appendix for the same type of result on the four-rooms environment). As before, VE largely outperforms MLE, in some cases with a significant improvement in performance. We call attention to the fact that in cart-pole we used neural networks to represent both the transition models \tilde{p} and the value functions \tilde{v} , which indicates that VE can be naturally applied with nonlinear function approximation.

It is important to note the broader significance of our experiments. While our theoretical analysis of value equivalence focused on the case where \mathcal{M} contained a value equivalent model, this is not

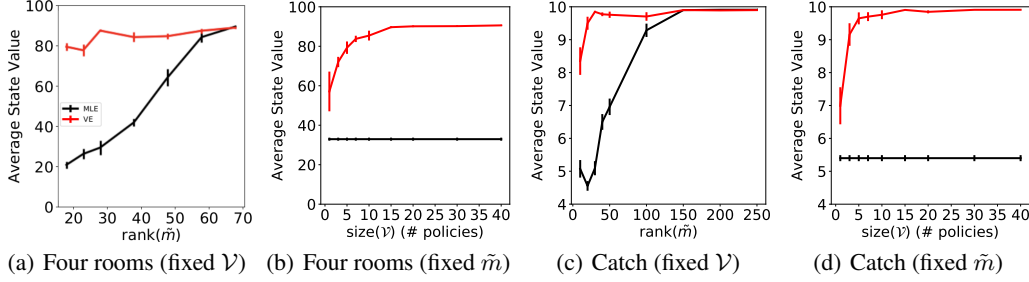


Figure 3: Results with \mathcal{V} composed of true value functions of randomly-generated policies. The models \tilde{p} are rank-constrained transition matrices $\tilde{P} = DK$, with $D \in \mathbb{R}^{|\mathcal{S}| \times k}$, $K \in \mathbb{R}^{k \times |\mathcal{S}|}$, and $k < |\mathcal{S}|$. Error bars are one standard deviation over 30 runs.

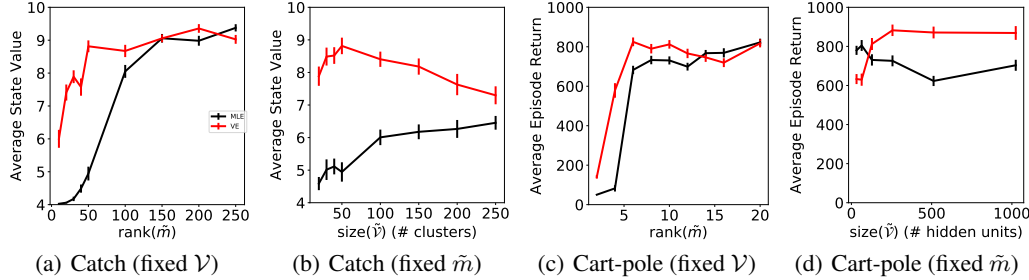


Figure 4: Results with \mathcal{V} composed of functions sampled from the agent’s representational space $\tilde{\mathcal{V}}$. **(a–b)** Functions in \mathcal{V} are the features of the linear function approximation (state aggregation), as per Remark 2. Models \tilde{p} are rank-constrained transition matrices (cf. Figure 3). **(c–d)** Functions in \mathcal{V} are randomly-generated neural networks. Models \tilde{p} are neural networks with rank-constrained linear transformations between layers (Appendix A.2). Error bars are one standard deviation over 30 runs.

guaranteed in practice. Our experiments illustrate that, in spite of lacking such a guarantee, we see a considerable gap in performance between VE and MLE, indicating that VE models still offer a strong benefit. Our goal here was to provide insight into the value equivalence principle; in the next section we point to prior work to demonstrate the utility of value equivalence in large-scale settings.

6 The value-equivalence principle in practice

Recently, there have been several successful empirical works that can potentially be understood as applications of the value-equivalence principle, like Silver et al.’s [34] *Predictron*, Oh et al.’s [26] *Value Prediction Networks*, Farquhar et al.’s [16] *TreeQN*, and Schrittwieser et al.’s [33] *MuZero*. Specifically, the model-learning aspect of these prior methods can be understood, with some abuse of notation, as a value equivalence principle of the form $\mathcal{T}v = \tilde{\mathcal{T}}v$, where \mathcal{T} is a Bellman operator applied with the true model m^* and $\tilde{\mathcal{T}}$ is a Bellman operator applied with an approximate model \tilde{m} .

There are many possible forms for the operators \mathcal{T} and $\tilde{\mathcal{T}}$. First, value equivalence can be applied to an uncontrolled Markov reward process; the resulting operator \mathcal{T}_π is analogous to having a single policy in Π . Second, it can be applied over n steps, using a Bellman operator \mathcal{T}_π^n that rolls the model forward n steps: $\mathcal{T}_\pi^n[v](s) = \mathbb{E}_\pi[R_{t+1} + \dots + \gamma^{n-1}R_{t+n} + \gamma^n v_\pi(S_{t+n}) | S_t = s]$, or a λ -weighted average \mathcal{T}_π^λ [6]. Third, a special case of the n -step operator $\mathcal{T}_{a_1, \dots, a_n}$ can be applied to an open-loop action sequence $\{a_1, \dots, a_n\}$. Fourth, it can be applied to the Bellman optimality operator, \mathcal{T}_{G_v} , where G_v is the “greedy” policy induced by v defined as $G_v(a|s) = \mathbb{1}\{a = \operatorname{argmax}_{a'} \mathbb{E}[R + \gamma v(S') | s, a']\}$. This idea can also be extended to an n -step greedy search operator, $\mathcal{T}_{G_v^n}[v](s) = \max_{a_1, \dots, a_n} \mathbb{E}[R_{t+1} + \dots + \gamma^{n-1}R_{t+n} + \gamma^n v(S_{t+n}) | S_t = s, A_t = a_1, \dots, A_{t+n} = a_n]$. Finally, instead of applying value equivalence over a fixed set of value functions \mathcal{V} , we can have a set \mathcal{V}_t that varies over time—for example, \mathcal{V}_t can be a singleton with an estimate of the value function of the current greedy policy.

The two operators \mathcal{T} and $\tilde{\mathcal{T}}$ can also differ. For example, on the environment side we can use the optimal value function, which can be interpreted as $\mathcal{T}^\infty v = v^*$ [40, 34], while the approximate operator can be $\tilde{\mathcal{T}}_\pi^\lambda$ [34] or $\tilde{\mathcal{T}}_{G_{v_t}}^n$ [40]. We can also use approximate values $\tilde{\mathcal{T}}v \approx \mathcal{T}v'$ where $v' \approx v$,

for example by applying n -step operators to approximate value functions, $\tilde{\mathcal{T}}^n v \approx \mathcal{T}^n v' = \mathcal{T}^n \mathcal{T}^k v = \mathcal{T}^{n+k} v$ [26, 33] or $\tilde{\mathcal{T}}^n v \approx \mathcal{T}^n v' = \mathcal{T}^n \tilde{\mathcal{T}}^k v$ [16], or even to approximate policies, $\tilde{\mathcal{T}}^n v_a \approx \mathcal{T}^n v'_a$ where $v_a = \pi(a|s) \approx \pi'(a|s) = v'_a$ for all $a \in \mathcal{A}$ [33]. The table below characterises the type of value equivalence principle used in prior work. We conjecture that this captures the essential idea underlying each method for model-learning, acknowledging that we ignore many important details.

Algorithm	Operator $\tilde{\mathcal{T}}$	Policies Π	Functions \mathcal{V}
Predictron [34]	$\tilde{\mathcal{T}}^\lambda v_t$	None	Value functions for pseudo-rewards
VIN [40]	$\tilde{\mathcal{T}}_{G_{v_t}}^n v_t$	G_{v_t}	Value function
TreeQN [16]	$\tilde{\mathcal{T}}_{G_{v_t}}^n v_t$	$G_{v_t}^n$	Value function
VPN [26]	$\tilde{\mathcal{T}}_{a_1 \dots a_n}^n v_t$	$\{a_1, \dots, a_n\} \sim \pi_t$	Value function
MuZero [33]	$\tilde{\mathcal{T}}_{a_1 \dots a_n}^n v_t$	$\{a_1, \dots, a_n\} \sim \pi_t$	Distributional value bins, policy components

All of these methods, with the exception of VIN, sample the Bellman operator, rather than computing full expectations (*c.f.* (7)). In addition, all of the above methods jointly learn the state representation alongside a value-equivalent model based upon that representation. Only MuZero includes both many policies and many functions, which may be sufficient to approximately span the policy and function space required to plan in complex environments; this perhaps explains its stronger performance.

7 Related work

Farahmand et al.’s [14, 15] *value-aware model learning* (VAML) is based on a premise similar to ours. They study a robust variant of (6) that considers the worst-case choice of $v \in \mathcal{V}$ and provide the gradient when the value-function approximation \tilde{v} is linear and the model \tilde{p} belongs to the exponential family. Later, Farahmand [13] also considered the case where the model is learned iteratively. Both versions of VAML come with finite sample-error upper bound guarantees [14, 15, 13]. More recently, Asadi et al. [2] showed that minimizing the VAML objective is equivalent to minimizing the Wasserstein metric. Abachi et al. [1] applied the VAML principle to policy gradient methods. The theory of VAML is complementary to ours: we characterise the space of value-equivalent models, while VAML focuses on the solution and analysis of the induced optimization problem.

Joseph et al. [22] note that minimizing prediction error is not the same as maximizing the performance of the resulting policy, and propose an algorithm that optimizes the parameters of the model rather than the policy’s. Ayoub et al. [3] proposes an algorithm that keeps a set of models that are consistent with the most recent value function estimate. They derive regret bounds for the algorithm which suggest that value-targeted regression estimation is both sufficient and efficient for model-based RL.

More broadly, other notions of equivalence between MDPs have been proposed in the literature [12, 28, 20, 31, 17, 23, 41, 29, 8, 42]. Any notion of equivalence over states can be recast as a form of state aggregation; in this case the functions mapping states to clusters can (and probably should) be used to enforce value equivalence (Remark 2). But the principle of value equivalence is more general: it can be applied with function approximations other than state aggregation and can be used to exploit structure in the problem even when there is no clear notion of state abstraction (Appendix A.1.2).

In this paper we have assumed that the agent has access to a well-defined notion of state $s \in \mathcal{S}$. More generally, the agent only receives observations from the environment and must construct its own state function—that is, a mapping from histories of observations to features representing states. This is an instantiation of the problem known as *representation learning* [43, 21, 9, 18, 45, 44, 19, 7]. An intriguing question which arises in this context is whether a model learned through value equivalence induces a space of “compatible” state representations, which would suggest that the loss (6) could also be used for representation learning. This may be an interesting direction for future investigations.

8 Conclusion

We introduced the principle of value equivalence: two models are value equivalent with respect to a set of functions and a set of policies if they yield the same updates of the former on the latter. Value equivalence formalizes the notion that models should be tailored to their future use and provides a mechanism to incorporate such knowledge into the model learning process. It also unifies some important recent work in the literature, shedding light on their empirical success. Besides helping to explain some past initiatives, we believe the concept of value equivalence may also give rise to theoretical and algorithmic innovations that leverage the insights presented.

Broader impact

The bulk of the research presented in this paper consists of foundational theoretical results about the learning of models for model-based reinforcement learning agents. While applications of these agents can have social impacts depending upon their use, our results merely serve to illuminate desirable properties of models and facilitate the subsequent training of agents using them. In short, this work is largely theoretical and does not present any foreseeable societal impact, except in the general concerns over progress in artificial intelligence.

Acknowledgements

We would like to thank Gregory Farquhar and Eszter Vertes for the great discussions regarding the value equivalence principle. We also thank the anonymous reviewers for their comments and suggestions to improve the paper.

References

- [1] Romina Abachi, Mohammad Ghavamzadeh, and Amir massoud Farahmand. Policy-Aware Model Learning for Policy Gradient Methods. *arXiv preprint cs.AI:2003.00030*, 2020.
- [2] Kavosh Asadi, Evan Cater, Dipendra Misra, and Michael L. Littman. Equivalence Between Wasserstein and Value-Aware Model-Based Reinforcement Learning. In *FAIM Workshop on Prediction and Generative Modeling in Reinforcement Learning*, 2018.
- [3] Alex Ayoub, Zeyu Jia, Csaba Szepesvári, Mengdi Wang, and Lin Yang. Model-Based Reinforcement Learning with Value-Targeted Regression. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2020.
- [4] Andrew G. Barto, Richard S. Sutton, and Charles W Anderson. Neuronlike Adaptive Elements That Can Solve Difficult Learning Control Problems. *IEEE Transactions on Systems, Man, and Cybernetics*, pages 834–846, 1983.
- [5] Marc Bellemare, Will Dabney, Robert Dadashi, Adrien Ali Taiga, Pablo Samuel Castro, Nicolas Le Roux, Dale Schuurmans, Tor Lattimore, and Clare Lyle. A Geometric Perspective on Optimal Representations for Reinforcement Learning. In *Advances in Neural Information Processing Systems*, pages 4360–4371, 2019.
- [6] Dimitri P. Bertsekas and John N. Tsitsiklis. *Neuro-Dynamic Programming*. Athena Scientific, 1996.
- [7] Ondrej Biza, Robert Platt, Jan-Willem van de Meent, and Lawson LS Wong. Learning Discrete State Abstractions with Deep Variational Inference. *arXiv preprint arXiv:2003.04300*, 2020.
- [8] Pablo Samuel Castro. Scalable Methods for Computing State Similarity in Deterministic Markov Decision Processes. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 10069–10076, 2020.
- [9] Dane Corneil, Wulfram Gerstner, and Johanni Brea. Efficient Model-Based Deep Reinforcement Learning with Variational State Tabulation. *arXiv preprint arXiv:1802.04325*, 2018.
- [10] Will Dabney, André Barreto, Mark Rowland, Robert Dadashi, John Quan, Marc G. Bellemare, and David Silver. The Value-Improvement Path: Towards Better Representations for Reinforcement Learning, 2020.
- [11] Robert Dadashi, Adrien Ali Taiga, Nicolas Le Roux, Dale Schuurmans, and Marc G. Bellemare. The Value Function Polytope in Reinforcement Learning. In *Proceedings of the International Conference on Machine Learning (ICML)*, volume 97, pages 1486–1495, 2019.
- [12] Thomas Dean and Robert Givan. Model Minimization in Markov Decision Processes. In *AAAI/IAAI*, pages 106–111, 1997.

- [13] Amir-massoud Farahmand. Iterative Value-Aware Model Learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 9090–9101, 2018.
- [14] Amir-Massoud Farahmand, André Barreto, and Daniel Nikovski. Value-Aware Loss Function for Model Learning in Reinforcement Learning. In *Proceedings of the European Workshop on Reinforcement Learning (EWRL)*, 2013.
- [15] Amir-Massoud Farahmand, André Barreto, and Daniel Nikovski. Value-Aware Loss Function for Model-Based Reinforcement Learning. In *Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 54, pages 1486–1494, 2017.
- [16] G Farquhar, T Rocktäschel, M Igl, and S Whiteson. TreeQN and ATreeC: Differentiable Tree-Structured Models for Deep Reinforcement Learning. In *6th International Conference on Learning Representations, ICLR 2018-Conference Track Proceedings*, volume 6. ICLR, 2018.
- [17] Norm Ferns, Prakash Panangaden, and Doina Precup. Metrics for Finite Markov Decision Processes. In *UAI*, volume 4, pages 162–169, 2004.
- [18] Vincent François-Lavet, Yoshua Bengio, Doina Precup, and Joelle Pineau. Combined Reinforcement Learning via Abstract Representations. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 3582–3589, 2019.
- [19] Carles Gelada, Saurabh Kumar, Jacob Buckman, Ofir Nachum, and Marc G Bellemare. Deep-MDP: Learning Continuous Latent Space Models for Representation Learning. In *International Conference on Machine Learning*, pages 2170–2179, 2019.
- [20] Robert Givan, Thomas Dean, and Matthew Greig. Equivalence Notions and Model Minimization in Markov Decision Processes. *Artificial Intelligence*, 147(1-2):163–223, 2003.
- [21] Maximillian Igl, Luisa Zintgraf, Tuan Anh Le, Frank Wood, and Shimon Whiteson. Deep Variational Reinforcement Learning for POMDPs. In *ICML 2018: Proceedings of the Thirty-Fifth International Conference on Machine Learning*, July 2018. URL <http://www.cs.ox.ac.uk/people/shimon.whiteson/pubs/iglicml18.pdf>.
- [22] Joshua Joseph, Alborz Geramifard, John W Roberts, Jonathan P How, and Nicholas Roy. Reinforcement Learning with Misspecified Model Classes. In *2013 IEEE International Conference on Robotics and Automation*, pages 939–946. IEEE, 2013.
- [23] Lihong Li, Thomas J. Walsh, and Michael L. Littman. Towards a Unified Theory of State Abstraction for MDPs. In *Proceedings of the International Symposium on Artificial Intelligence and Mathematics*, pages 531–539, 2006.
- [24] Russell B. Millar. *Maximum Likelihood Estimation and Inference*. Hoboken: Wiley, 2011.
- [25] Volodymyr Mnih, Nicolas Heess, Alex Graves, and Koray Kavukcuoglu. Recurrent Models of Visual Attention. In *Advances in Neural Information Processing Systems (NIPS)*, pages 2204–2212, 2014.
- [26] Junhyuk Oh, Satinder Singh, and Honglak Lee. Value Prediction Networks. In *Advances in Neural Information Processing Systems*, pages 6118–6128, 2017.
- [27] Ronald Parr, Lihong Li, Gavin Taylor, Christopher Painter-Wakefield, and Michael L. Littman. An Analysis of Linear Models, Linear Value-Function Approximation, and Feature Selection for Reinforcement Learning. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 752–759, 2008.
- [28] Pascal Poupart and Craig Boutilier. Value-Directed Compression of POMDPs. In *Advances in Neural Information Processing Systems (NIPS)*, pages 1547–1554. MIT Press, 2002.
- [29] Pascal Poupart and Craig Boutilier. Value-Directed Belief State Approximation for POMDPs. *CoRR*, abs/1301.3887, 2013. URL <http://arxiv.org/abs/1301.3887>.
- [30] Martin L. Puterman. *Markov Decision Processes—Discrete Stochastic Dynamic Programming*. John Wiley & Sons, Inc., 1994.

- [31] Balaraman Ravindran and Andrew G Barto. Approximate Homomorphisms: A Framework for Non-Exact Minimization in Markov Decision Processes. 2004.
- [32] Stuart J. Russell and Peter Norvig. *Artificial Intelligence: A Modern Approach*. Pearson Education, 3 edition, 2003.
- [33] Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Simon Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, et al. Mastering Atari, Go, Chess and Shogi by Planning with a Learned Model. *arXiv preprint arXiv:1911.08265*, 2019.
- [34] David Silver, Hado van Hasselt, Matteo Hessel, Tom Schaul, Arthur Guez, Tim Harley, Gabriel Dulac-Arnold, David Reichert, Neil Rabinowitz, Andre Barreto, et al. The Predictron: End-to-End Learning and Planning. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 3191–3199. JMLR. org, 2017.
- [35] Richard S. Sutton. Learning to Predict by the Methods of Temporal Differences. *Machine Learning*, 3:9–44, 1988.
- [36] Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, 2018. URL <https://mitpress.mit.edu/books/reinforcement-learning-second-edition>. 2nd edition.
- [37] Richard S Sutton, Doina Precup, and Satinder Singh. Between MDPs and Semi-MDPs: A Framework for Temporal Abstraction in Reinforcement Learning. *Artificial intelligence*, 112 (1-2):181–211, 1999.
- [38] Richard S. Sutton, Csaba Szepesvári, Alborz Geramifard, and Michael Bowling. Dyna-Style Planning with Linear Function Approximation and Prioritized Sweeping. In *Proceedings of the Conference on Uncertainty in Artificial Intelligence (UAI)*, page 528–536, 2008.
- [39] Csaba Szepesvári. *Algorithms for Reinforcement Learning*. Synthesis Lectures on Artificial Intelligence and Machine Learning. Morgan & Claypool Publishers, 2010.
- [40] Aviv Tamar, Yi Wu, Garrett Thomas, Sergey Levine, and Pieter Abbeel. Value Iteration Networks. In *Advances in Neural Information Processing Systems*, pages 2154–2162, 2016.
- [41] Jonathan Taylor, Doina Precup, and Prakash Panagaden. Bounding Performance Loss in Approximate MDP Homomorphisms. In *Advances in Neural Information Processing Systems (NIPS)*, pages 1649–1656, 2009.
- [42] Elise van der Pol, Thomas Kipf, Frans A Oliehoek, and Max Welling. Plannable Approximations to MDP Homomorphisms: Equivariance under Actions. In *International Conference on Autonomous Agents and MultiAgent Systems*, 2020.
- [43] Manuel Watter, Jost Springenberg, Joschka Boedecker, and Martin Riedmiller. Embed to Control: A Locally Linear Latent Dynamics Model for Control from Raw Images. In *Advances in Neural Information Processing Systems (NIPS)*, pages 2746–2754, 2015.
- [44] Amy Zhang, Zachary C Lipton, Luis Pineda, Kamyar Azizzadenesheli, Anima Anandkumar, Laurent Itti, Joelle Pineau, and Tommaso Furlanello. Learning Causal State Representations of Partially Observable Environments. *arXiv preprint arXiv:1906.10437*, 2019.
- [45] Marvin Zhang, Sharad Vikram, Laura Smith, Pieter Abbeel, Matthew Johnson, and Sergey Levine. SOLAR: Deep Structured Representations for Model-Based Reinforcement Learning. In *International Conference on Machine Learning (ICML)*, pages 7444–7453. PMLR, 2019.

The Value-Equivalence Principle for Model-Based Reinforcement Learning

Supplementary Material

Christopher Grimm
Computer Science & Engineering
University of Michigan
crgrimm@umich.edu

André Barreto, Satinder Singh, David Silver
DeepMind
{andrebarreto,baveja,davidsilver}@google.com

In this supplement we give details of our theoretical results and experiments that had to be left out of the main paper due to space constraints. We prove our theoretical results and provide a detailed description of our experimental procedure. Importantly, we present an illustrative example showing how value equivalence (VE) may lead to a better solution for a Markov decision process (MDP) than maximum-likelihood estimate (MLE). *We show this to be true both in the exact case, when there exist a value-equivalent model in the model class considered, and in the approximate case, when such a model does not exist in the model class.* Our appendix is organized as follows:

- Section A.1.1 contains derivations of the properties and propositions presented in the main text.
- Section A.1.2 contains a sequence of examples using a toy MDP that illustrate points made in the discussion surrounding Propositions 3 and 4. Moreover, we include an additional result which illustrates a situation in which approximate VE models can outperform the MLE model.
- Section A.2 provides a detailed outline of the pipeline used across our experiments in the main text. We also report several additional results that had to be left out of the main paper due to space constraints.

The numbering of equations, figures and citations resume from what is used in the main paper.

A Appendix

A.1 Proofs of theoretical results and illustrative examples

A.1.1 Proofs

Property 1. *Given $\mathcal{M}' \subseteq \mathcal{M}$, we have that $\mathcal{M}'(\Pi, \mathcal{V}) \subseteq \mathcal{M}(\Pi, \mathcal{V})$.*

Proof. This result directly follows from Definitions 1 and 2. □

Property 2. *$\mathcal{M}(\Pi, \mathcal{V})$ either contains m^* or is the empty set.*

Proof. $\mathcal{M}(\Pi, \mathcal{V}) \subseteq \mathbb{M}(\Pi, \mathcal{V}) = \{m^*\}$ (Property 1). □

Property 3. *Given $\Pi' \subseteq \Pi$ and $\mathcal{V}' \subseteq \mathcal{V}$, we have that $\mathcal{M}(\Pi, \mathcal{V}) \subseteq \mathcal{M}(\Pi', \mathcal{V}')$.*

Proof. We will show the result by contradiction. Suppose there is a model $\tilde{m} \in \mathcal{M}(\Pi, \mathcal{V})$ such that $\tilde{m} \notin \mathcal{M}(\Pi', \mathcal{V}')$. This means that there exists a $\pi \in \Pi'$ and a $v \in \mathcal{V}'$ for which $\tilde{\mathcal{T}}_\pi v \neq \mathcal{T}_\pi v$. But since $\Pi' \subseteq \Pi$ and $\mathcal{V}' \subseteq \mathcal{V}$, it must be the case that $\pi \in \Pi$ and $v \in \mathcal{V}$, which contradicts the claim that $\tilde{m} \in \mathcal{M}(\Pi, \mathcal{V})$. □

Property 4. *If $m^* \in \mathcal{M}$, then $m^* \in \mathcal{M}(\Pi, \mathcal{V})$ for all Π and all \mathcal{V} .*

Proof. $m^* \in \mathcal{M}(\Pi, \mathcal{V}) \subseteq \mathcal{M}(\Pi, \mathcal{V})$ (Property 3). □

Proposition 1. For discrete Π and \mathcal{V} , we have that $\mathcal{M}(\Pi, \mathcal{V}) = \mathcal{M}(p\text{-span}(\Pi) \cap \mathbb{I}, \text{span}(\mathcal{V}))$.

Proof. Let $\pi \in p\text{-span}(\Pi) \cap \mathbb{I}$. Based on (3), we know that there exists an $\alpha_s \in \mathbb{R}^{|\Pi|}$ such that $\pi(\cdot|s) = \sum_i \alpha_{si} \pi_i(\cdot|s)$, where $\pi_i \in \Pi$. Thus, for $\tilde{m} \in \mathcal{M}(\Pi, \mathcal{V})$, we can write

$$\begin{aligned} \tilde{\mathcal{T}}_\pi[v](s) &= \mathbb{E}_{A \sim \pi(\cdot|s), S' \sim \tilde{p}(\cdot|s, A)} [\tilde{r}(s, A) + \gamma v(S')] \\ &= \int \pi(a|s) \mathbb{E}_{S' \sim \tilde{p}(\cdot|s, a)} [\tilde{r}(s, a) + \gamma v(S')] da \\ &= \int \sum_i \alpha_{si} \pi_i(a|s) \mathbb{E}_{S' \sim \tilde{p}(\cdot|s, a)} [\tilde{r}(s, a) + \gamma v(S')] da \\ &= \sum_i \alpha_{si} \int \pi_i(a|s) \mathbb{E}_{S' \sim \tilde{p}(\cdot|s, a)} [\tilde{r}(s, a) + \gamma v(S')] da \\ &= \sum_i \alpha_{si} \mathbb{E}_{A \sim \pi_i(\cdot|s), S' \sim \tilde{p}(\cdot|s, a)} [\tilde{r}(s, a) + \gamma v(S')] \\ &= \sum_i \alpha_{si} \mathcal{T}^{\pi_i}[v](s). \end{aligned}$$

Let $v \in \text{span}(\mathcal{V})$. We know there is a $\beta \in \mathbb{R}^{|\mathcal{V}|}$ such that $v = \sum_i \beta_i v_i$, with $v_i \in \mathcal{V}$.

$$\begin{aligned} \tilde{\mathcal{T}}_\pi[v](s) &= \mathbb{E}_{A \sim \pi(\cdot|s), S' \sim \tilde{p}(\cdot|s, A)} [\tilde{r}(s, A) + \gamma \sum_i \beta_i v_i(S')] \\ &= \sum_i \beta_i \mathbb{E}_{A \sim \pi(\cdot|s), S' \sim \tilde{p}(\cdot|s, A)} [\tilde{r}(s, A) + \gamma v_i(S')] \\ &= \sum_i \beta_i \tilde{\mathcal{T}}_\pi[v_i](s). \end{aligned}$$

□

In order to prove Proposition 2 we will need four lemmas which we state and prove below.

Lemma 1. For arbitrary matrices $\mathbf{A} \in \mathbb{R}^{k \times n}$, $\mathbf{C} \in \mathbb{R}^{m \times \ell}$, we can construct a vector-space $\mathcal{B} = \{\mathbf{B} \in \mathbb{R}^{n \times m} : \mathbf{ABC} = \mathbf{0}\}$ where $\mathbf{0}$ denotes a $k \times \ell$ matrix of zeros. It follows that

$$\mathcal{H}\text{-dim}[\mathcal{B}] = nm - \text{rank}(\mathbf{A}) \cdot \text{rank}(\mathbf{C}). \quad (8)$$

Proof. We begin by converting the condition $\mathbf{ABC} = \mathbf{0}$ into a matrix-vector product. Let \mathbf{a}^i and \mathbf{c}^j denote the i 'th row of \mathbf{A} and j 'th column of \mathbf{C} respectively. Observe that $(\mathbf{ABC})_{ij} = \mathbf{a}^i \mathbf{B} \mathbf{c}^j = \sum_{x,y} \mathbf{a}_x^i \mathbf{c}_y^j \mathbf{B}_{xy}$, which implies that

$$\mathbf{ABC} = \mathbf{0} \iff \sum_{x,y} \mathbf{a}_x^i \mathbf{c}_y^j \mathbf{B}_{xy} = 0 \quad \forall i \in [k], j \in [\ell] \quad (9)$$

where $[k]$ denotes $\{1, \dots, k\}$.

For each (i, j) pair, the above expression is suggestive of a dot-product between two $n \times m$ vectors: a combination of \mathbf{a}^i and \mathbf{c}^j , and a ‘‘flattened’’ version of \mathbf{B} . Define the former combination of vectors as $\mathbf{d}^{ij} = [\mathbf{a}_1^i \mathbf{c}_1^j, \mathbf{a}_1^i \mathbf{c}_2^j, \dots, \mathbf{a}_n^i \mathbf{c}_m^j]^\top \in \mathbb{R}^{nm \times 1}$, and stack them as rows as: $\mathbf{D} = [\mathbf{d}^{11}, \mathbf{d}^{12}, \dots, \mathbf{d}^{nm}]^\top \in \mathbb{R}^{k\ell \times nm}$. To flatten \mathbf{B} , simply define $\mathbf{b} = [\mathbf{B}_{11}, \mathbf{B}_{12}, \dots, \mathbf{B}_{nm}]^\top \in \mathbb{R}^{nm \times 1}$.

We now have that $\mathbf{ABC} = \mathbf{0} \iff \mathbf{D}\mathbf{b} = \mathbf{0}$. Moreover, unravelling the matrices in \mathcal{B} does not change the dimension of the space, thus:

$$\mathcal{H}\text{-dim}[\mathcal{B}] = \mathcal{H}\text{-dim}[\{\mathbf{b} \in \mathbb{R}^{nm \times 1} : \mathbf{D}\mathbf{b} = \mathbf{0}\}] = nm - \text{rank}(\mathbf{D}) \quad (10)$$

where the last equality comes from a application of the rank-nullity theorem.

Finally notice that the construction of \mathbf{d}^{ij} can be thought of as vertically stacking n copies of \mathbf{c}^j each scaled by a different entry in \mathbf{a}^i . We can also find scaled copies of \mathbf{a}^i by \mathbf{c}_k^j in \mathbf{d}^{ij} by selecting indices from the combined vector at regular intervals of m : $\mathbf{d}_{k+(\ell-1)m}^{ij} = \mathbf{c}_k^j \cdot \mathbf{a}_\ell^i$ for $\ell \in \{1, \dots, n\}$.

This means that scaled copies of both \mathbf{a}^i and \mathbf{c}^j can be found by selecting specific groups of indices in \mathbf{d}^{ij} . It follows that if $\mathbf{a}^1, \dots, \mathbf{a}^n$ are linearly independent then so are $\mathbf{d}^{1j}, \dots, \mathbf{d}^{nj}$ for any j . And similarly, if $\mathbf{c}^1, \dots, \mathbf{c}^m$ are linearly independent then so are $\mathbf{d}^{i1}, \dots, \mathbf{d}^{im}$ for any i . Hence if $\mathbf{a}^1, \dots, \mathbf{a}^n$ and $\mathbf{c}^1, \dots, \mathbf{c}^m$ are both linearly independent sets, then so is $\mathbf{d}^{11}, \mathbf{d}^{12}, \dots, \mathbf{d}^{nm}$. Since these \mathbf{a}^i and \mathbf{c}^j vectors form the rows and columns of rank n and m matrices: \mathbf{A} and \mathbf{C} , their corresponding sets of row and column vectors are linearly independent. Thus we have that $\text{rank}(\mathbf{D}) = \text{rank}(\mathbf{A}) \cdot \text{rank}(\mathbf{C})$, completing the proof. □

Lemma 2. For any c and $\mathcal{Y} + c = \{y + c : y \in \mathcal{Y}\}$ it follows that $\dim[\mathcal{Y} + c] = \dim[\mathcal{Y}]$.

Proof.

$$\dim[\mathcal{Y} + \mathbf{c}] = \min_{(\mathcal{V}, \mathbf{c}'): \mathcal{Y} + (\mathbf{c} + \mathbf{c}') \subseteq \mathcal{W}} \mathcal{H}\text{-dim}[\mathcal{W}] = \min_{(\mathcal{W}, \mathbf{c}') : \mathcal{Y} + \mathbf{c}' \subseteq \mathcal{W}} \mathcal{H}\text{-dim}[\mathcal{W}] = \dim[\mathcal{Y}]$$

□

Lemma 3. *If \mathcal{Y} is a vector-space then $\mathcal{H}\text{-dim}[\mathcal{Y}] = \dim[\mathcal{Y}]$.*

Proof. Recall the definition of $\dim[\mathcal{Y}]$:

$$\dim[\mathcal{Y}] = \min_{(\mathcal{W}, \mathbf{c}) : \mathcal{Y} + \mathbf{c} \subseteq \mathcal{W}} \mathcal{H}\text{-dim}[\mathcal{W}]$$

where \mathcal{W} is a vector-space. By choosing $\mathcal{W} = \mathcal{Y}$ and $\mathbf{c} = \mathbf{0}$ we see that $\dim[\mathcal{Y}] \leq \mathcal{H}\text{-dim}[\mathcal{Y}]$.

Suppose then that $\dim[\mathcal{Y}] < \mathcal{H}\text{-dim}[\mathcal{Y}]$. This implies that there is a vector space \mathcal{W} and offset \mathbf{c} with $d = \mathcal{H}\text{-dim}[\mathcal{W}] < \mathcal{H}\text{-dim}[\mathcal{Y}]$ and $\mathcal{Y} + \mathbf{c} \subseteq \mathcal{W}$. This means that for every $\mathbf{y} \in \mathcal{Y}$: $\mathbf{y} + \mathbf{c} = \sum_{i=1}^d \alpha_i^{\mathbf{y}} \mathbf{w}_i$ for some $\alpha_{1:d}^{\mathbf{y}}$ where $\mathbf{w}_{1:d}$ are a basis of \mathcal{W} . Since \mathcal{Y} is a vector space it must contain the $\mathbf{0}$ vector, hence $\mathbf{c} = \sum_{i=1}^d \alpha_i^{\mathbf{0}} \mathbf{w}_i$. Accordingly any $\mathbf{y} \in \mathcal{Y}$ can be written as $\mathbf{y} = \sum_{i=1}^d (\alpha_i^{\mathbf{y}} - \alpha_i^{\mathbf{0}}) \mathbf{w}_i$. However, this is a contradiction since $\mathcal{H}\text{-dim}[\mathcal{W}] < \mathcal{H}\text{-dim}[\mathcal{Y}]$. Hence $\dim[\mathcal{Y}] = \mathcal{H}\text{-dim}[\mathcal{Y}]$. □

Lemma 4. *If $\mathcal{X} \subseteq \mathcal{Y}$ then $\dim[\mathcal{X}] \leq \dim[\mathcal{Y}]$.*

Proof. If $\mathcal{X} \subseteq \mathcal{Y}$ then for any \mathbf{c} , $\mathcal{X} + \mathbf{c} \subseteq \mathcal{Y} + \mathbf{c}$. Because of the above, for any vector-space \mathcal{W} : $\mathcal{W} \supseteq \mathcal{Y} + \mathbf{c} \implies \mathcal{W} \supseteq \mathcal{X} + \mathbf{c}$, hence: $\{(\mathcal{W}, \mathbf{c}) : \mathcal{X} + \mathbf{c} \subseteq \mathcal{W}\} \supseteq \{(\mathcal{W}, \mathbf{c}) : \mathcal{Y} + \mathbf{c} \subseteq \mathcal{W}\}$. Notice that this last set-relation corresponds the set of vector-spaces that $\dim[\cdot]$ is minimizing over for \mathcal{X} and \mathcal{Y} respectively. Hence $\dim[\mathcal{X}] \leq \dim[\mathcal{Y}]$. □

Proposition 2. *Let Π be a set of m pointwise linearly independent policies $\pi_i \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ and let \mathcal{V} be a set of k linearly independent vectors $\mathbf{v}_i \in \mathbb{R}^{|\mathcal{S}|}$. Then,*

$$\dim[\mathbb{P}(\Pi, \mathcal{V})] \leq |\mathcal{S}| (|\mathcal{S}||\mathcal{A}| - mk).$$

Proof. First note that if $\pi_i \notin p\text{-span}(\Pi \setminus \{\pi_i\})$ then $\pi_i \notin \text{span}(\Pi \setminus \{\pi_i\})$. Hence, pointwise linear independence implies linear independence.

Since $|\mathcal{S}|$ and $|\mathcal{A}|$ are finite, we can assume that $\mathcal{A} = \{1, \dots, |\mathcal{A}|\}$ and $\mathcal{S} = \{1, \dots, |\mathcal{S}|\}$. For any transition probability kernel $\tilde{p}(s'|s, a)$ we can construct matrix $\tilde{\mathbf{P}} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}|}$ with $\tilde{\mathbf{P}}_{(a-1)|\mathcal{S}|+s, s'} = \tilde{p}(s'|s, a)$. Denote the constructed matrix corresponding to the true dynamics as \mathbf{P} . For any π_i we can construct a matrix $\mathbf{\Pi}_i \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}||\mathcal{A}|}$ with $(\mathbf{\Pi}_i)_{s, (a-1)|\mathcal{S}|+s} = \pi_i(a|s)$. Vertically stack these m $\mathbf{\Pi}_i$ matrices to construct $\mathbf{\Pi} \in \mathbb{R}^{m|\mathcal{S}| \times |\mathcal{S}||\mathcal{A}|}$. Additionally we construct $\mathbf{V} \in \mathbb{R}^{|\mathcal{S}| \times k}$ with $\mathbf{V}_{j,\ell} = (\mathbf{v}_\ell)_j$. Note that $\mathbb{P}(\Pi, \mathcal{V}) = \{\tilde{\mathbf{P}} \in \mathbb{P} : \mathbf{\Pi}(\tilde{\mathbf{P}} - \mathbf{P})\mathbf{V} = \mathbf{0}\}$. Define the sets $\mathcal{X} = \{\mathbf{X} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}|} : \mathbf{P}\mathbf{X}\mathbf{V} = \mathbf{0}\}$ and $\mathcal{Y} = \{\tilde{\mathbf{P}} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}|} : \mathbf{\Pi}(\tilde{\mathbf{P}} - \mathbf{P})\mathbf{V} = \mathbf{0}\}$.

Note the following three facts:

1. $\dim[\mathcal{X}] = \dim[\mathcal{Y}]$ since our notion of dimension is translation-invariant (Lemma 2).
2. $\dim[\mathcal{X}] = \mathcal{H}\text{-dim}[\mathcal{X}]$ since \mathcal{X} is a vector-space (Lemma 3).
3. $\mathbb{P}(\Pi, \mathcal{V}) \subseteq \mathcal{Y}$ which implies that $\dim[\mathbb{P}(\Pi, \mathcal{V})] \leq \dim[\mathcal{Y}]$ (Lemma 4).

Taken together this gives us that

$$\dim[\mathbb{P}(\Pi, \mathcal{V})] \leq \dim[\mathcal{Y}] = \mathcal{H}\text{-dim}[\mathcal{X}].$$

We can now apply Lemma 1 to obtain $\dim[\mathcal{X}] = |\mathcal{S}|^2|\mathcal{A}| - k \cdot \text{rank}(\mathbf{\Pi})$. Notice that $\text{rank}(\mathbf{\Pi}) = \min\{|\mathcal{S}||\mathcal{A}|, m|\mathcal{S}|\}$. Thus $\dim[\mathbb{P}(\Pi, \mathcal{V})] \leq |\mathcal{S}|(|\mathcal{S}||\mathcal{A}| - mk)$ as needed. □

Proposition 3. *The maximum-likelihood estimate of p^* in \mathcal{P} may not belong to a $\mathcal{P}(\Pi, \mathcal{V}) \neq \emptyset$.*

Proof. Suppose we are trying to estimate a transition matrix $\mathbf{P} \in \mathbb{R}^{n \times n}$ and choose to use one parameter $\theta_i \in \mathbb{R}$ per row. Specifically, we parametrize the distribution on the i -th row as

$$\tilde{p}_{ii} = \theta_i \text{ and } \tilde{p}_{ij} = (1 - \theta_i)/(n - 1), \text{ for } i \neq j, \text{ with } \theta_i \in [0, 1],$$

where $p_{ij} = p(s_j | s_i)$. We can then write the expected likelihood function for $\boldsymbol{\theta} \in \mathbb{R}^n$ as

$$\begin{aligned} m(\boldsymbol{\theta}) &= \sum_i \left[p_{ii} \ln \theta_i + \sum_{j \neq i} p_{ij} \ln(1 - \theta_i) - \sum_{j \neq i} p_{ij} \ln(n - 1) \right] \\ &= \sum_i \left[p_{ii} \ln \theta_i + (1 - p_{ii}) \ln(1 - \theta_i) - (1 - p_{ii}) \ln(n - 1) \right], \end{aligned}$$

which leads to the likelihood equation

$$0 = \frac{\partial m(\boldsymbol{\theta})}{\partial \theta_i} = \frac{p_{ii}}{\theta_i} + \frac{1 - p_{ii}}{\theta_i - 1} = \frac{p_{ii}(\theta_i - 1) + (1 - p_{ii})\theta_i}{\theta_i(\theta_i - 1)} = \frac{\theta_i - p_{ii}}{\theta_i(\theta_i - 1)}.$$

The MLE solution is thus to have $\theta_i = p_{ii}$ for $i = 1, 2, \dots, n$. This means that the solution provided by MLE will not be exact if and only if

$$p_{ij} \neq p_{ik} \text{ for any } (i, j, k) \text{ such that } i \neq j \neq k. \quad (11)$$

Now, suppose we have $\mathcal{V} = \{v\}$ with $v_i = 1$ for some i and $v_j = 0$ for $j \neq i$. In this case it is possible to get an exact value-equivalent solution—that is, $\mathbf{P}v = \tilde{\mathbf{P}}v$ —by making $\theta_i = p_{ii}$ and $\theta_j = 1 - (n - 1)p_{ii}$ for $j \neq i$, regardless of whether (11) is true or not. \square

Proposition 4. *Suppose $v \in \mathcal{V}' \implies \mathcal{T}_\pi v \in \mathcal{V}'$ for all $\pi \in \mathbb{I}$. Let $p\text{-span}(\mathbb{I}) \supseteq \mathbb{I}$ and $\text{span}(\mathcal{V}) = \mathcal{V}'$. Then, starting from any $v' \in \mathcal{V}'$, any $\tilde{m} \in \mathcal{M}(\mathbb{I}, \mathcal{V})$ yields the same solution as m^* .*

Proof. Denote the Bellman operator under a policy that always selects action a as \mathcal{T}_a , the greedy Bellman operator as $\mathcal{T}v = \max_a \mathcal{T}_a v$ and the Bellman operator under a policy π as \mathcal{T}_π , as before. Let $\mathcal{T}^{(n)}v$ represent n successive applications of operator \mathcal{T} on value v .

Note that for any $v \in \mathbb{V}$ we can construct a $\pi_v(s) = \arg\max_a (\mathcal{T}_a v)(s)$ such that $\mathcal{T}v = \max_a \mathcal{T}_a v = \mathcal{T}_{\pi_v} v$. This implies that the greedy Bellman operator is included in the assumption of our proposition:

$$v \in \mathcal{V}' \implies \mathcal{T}v \in \mathcal{V}'. \quad (12)$$

We now begin by showing that:

$$\mathcal{T}^{(n)}v = \tilde{\mathcal{T}}^{(n)}v \in \mathcal{V}' \implies \mathcal{T}^{(n+1)}v = \tilde{\mathcal{T}}^{(n+1)}v \in \mathcal{V}' \quad (13)$$

for any $v \in \mathbb{V}$ and any $n > 0$. Assume that $\mathcal{T}^{(n)}v = \tilde{\mathcal{T}}^{(n)}v \in \mathcal{V}'$. Since $\mathcal{T}^{(n)}v \in \mathcal{V}'$ and $\mathcal{V}' = \text{span}(\mathcal{V})$, we can use value equivalence to obtain:

$$\mathcal{T}_a \mathcal{T}^{(n)}v = \tilde{\mathcal{T}}_a \mathcal{T}^{(n)}v.$$

for any $a \in \mathcal{A}$. Next, since $\mathcal{T}^{(n)}v = \tilde{\mathcal{T}}^{(n)}v$ we can write:

$$\mathcal{T}_a \mathcal{T}^{(n)}v = \tilde{\mathcal{T}}_a \tilde{\mathcal{T}}^{(n)}v. \quad (14)$$

Since (14) holds for any $a \in \mathcal{A}$, we can write:

$$\mathcal{T}^{(n+1)}v = \max_a \mathcal{T}_a \mathcal{T}^{(n)}v = \max_a \tilde{\mathcal{T}}_a \tilde{\mathcal{T}}^{(n)}v = \tilde{\mathcal{T}}^{(n+1)}v.$$

We know from (12) that the fact that $\mathcal{T}^{(n)}v \in \mathcal{V}'$ implies that $\mathcal{T}^{(n+1)}v \in \mathcal{V}'$. Thus we have shown that (13) is true.

Finally, by choosing $v \in \mathcal{V}'$ and using analogous reasoning as above, we can show that $\mathcal{T}_a v = \tilde{\mathcal{T}}_a v$ and $\mathcal{T}v = \max_a \mathcal{T}_a v = \max_a \tilde{\mathcal{T}}_a v = \tilde{\mathcal{T}}v$, and since $v \in \mathcal{V}'$, $\tilde{\mathcal{T}}v = \mathcal{T}v \in \mathcal{V}'$. Thus $\mathcal{T}^{(n)}v = \tilde{\mathcal{T}}^{(n)}v$ for all $n \in \mathbb{N}$. This is sufficient to conclude that

$$\tilde{v}^* = \lim_{n \rightarrow \infty} \tilde{\mathcal{T}}^{(n)}v' = \lim_{n \rightarrow \infty} \mathcal{T}^{(n)}v' = v^*,$$

as needed. \square

A.1.2 Examples with a simple MDP

Consider the 3 state MDP with states s_1, s_2, s_3 and actions $\mathcal{A} = \{a_1, a_2\}$. Transitioning to state s_1 always incurs a reward of 1, taking any action in states s_2 and s_3 always results in transitioning to s_1 and taking action $a \in \mathcal{A}$ from s_1 transitions among the other states according to action-dependent distribution $(p_{11}^a, p_{12}^a, p_{13}^a)$. This MDP is depicted in Figure 5. We now use this MDP to illustrate several points made in the main text.

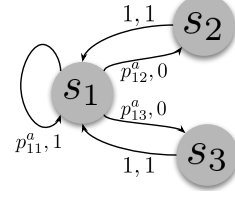


Figure 5

Closure under Bellman updates We now address the discussion surrounding Proposition 4 in the main text.

Consider a the following two-dimensional subspace of value functions $\mathcal{R} = \{[x, y, y]^\top : x, y \in \mathbb{R}\}$. We now show that, for the MDP described above, \mathcal{R} exhibits closure under arbitrary Bellman updates.

For an arbitrary policy $\pi : \mathcal{S} \mapsto \mathcal{P}(\mathcal{A})$ the Bellman update for a value function $v \in \mathbb{R}^3$ is given by $\mathcal{T}^\pi v = \mathbf{R}^\pi + \gamma \mathbf{P}^\pi v$ where

$$\mathbf{R}^\pi = \begin{bmatrix} \sum_{a \in \mathcal{A}} \pi(a|s_1) p_{11}^a \\ 1 \\ 1 \end{bmatrix}, \quad \mathbf{P}^\pi = \begin{bmatrix} \sum_{a \in \mathcal{A}} \pi(a|s_1) p_{11}^a & \sum_{a \in \mathcal{A}} \pi(a|s_2) p_{12}^a & \sum_{a \in \mathcal{A}} \pi(a|s_3) p_{13}^a \\ 1 & 0 & 0 \\ 1 & 0 & 0 \end{bmatrix}$$

Suppose $v \in \mathcal{R}$, then $v = [a, b, b]^\top$ for some $a, b \in \mathbb{R}$. Notice that for such a value function the following holds:

$$\mathcal{T}^\pi v = \begin{bmatrix} \mathbf{R}_1^\pi + \gamma[a\mathbf{P}_{11}^\pi + b(1 - \mathbf{P}_{11}^\pi)] \\ 1 + \gamma a \\ 1 + \gamma a \end{bmatrix} \in \mathcal{R},$$

thus we have illustrated that the two-dimensional subspace \mathcal{R} is closed under arbitrary Bellman updates in our 3 state MDP. This means that, once a sequence $v_1, v_2 = \mathcal{T}_\pi v_1, v_3 = \mathcal{T}_\pi v_2 \dots$ reaches a $v_i \in \mathcal{R}$, it stays in \mathcal{R} . We can then exploit this property finding value-equivalent models with respect to \mathcal{R} , as we show next.

A model class for which exact VE outperforms MLE We now provide an example of the scenario discussed around Proposition 3 in the main text by examining the setting where a model, from a restricted class, must be learned to approximate the dynamics of our MDP. We restrict our model class by requiring that for each action $a \in \mathcal{A}$ we represent $(p_{11}^a, p_{12}^a, p_{13}^a)$ as $((1 - \theta^a)/2, \theta^a, (1 - \theta^a)/2)$. Before continuing we note a few properties of value functions of our MDP. Notice that for any v^π we can write:

$$\begin{aligned} v_1^\pi &= \sum_{a \in \mathcal{A}} \pi(a|s_1) [p_{11}^a (1 + \gamma v_1^\pi) + (1 - p_{11}^a) (\gamma^2 v_1^\pi)], \\ v_2^\pi &= 1 + \gamma v_1^\pi, \\ v_3^\pi &= 1 + \gamma v_1^\pi, \end{aligned}$$

which illustrates that v^π *exclusively depends* on the value of $\mathbf{P}_{11}^\pi \equiv \sum_{a \in \mathcal{A}} \pi(a|s_1) p_{11}^a$.

First we consider the MLE solution to this problem: it can be easily shown (see the proof of Proposition 3) that, for the model class defined above, $\theta^a = p_{12}^a$ for all $a \in \mathcal{A}$ maximizes the likelihood. However notice that this implies that our approximation of p_{11}^a equals $(1 - p_{12}^a)/2$ which is clearly not true in general. Thus, there are settings of $(p_{11}^a, p_{12}^a, p_{13}^a)$ and policies for which the value function produced by MLE, \tilde{v}^π , is not equivalent to the true value function v^π .

Next we consider learning a value-equivalent model with the same restricted model class. Suppose we wish our model to be value equivalent to value $v = [1, 0, 0]^\top$ and all policies.

Note that any VE model with respect to $\mathcal{V} = \{v\} : \{\tilde{\mathbf{P}}^a\}_{a \in \mathcal{A}}$, must satisfy $\tilde{\mathbf{P}}^a v = \mathbf{P}^a v$. By requiring value equivalence with just v we have:

$$\tilde{\mathbf{P}}^a v = \begin{bmatrix} \tilde{p}_{11}^a \\ \tilde{p}_{21}^a \\ \tilde{p}_{31}^a \end{bmatrix} = \begin{bmatrix} p_{11}^a \\ 1 \\ 1 \end{bmatrix} = \mathbf{P}^a v$$

which implies that $\tilde{p}_{11}^a = p_{11}^a, \tilde{p}_{21}^a = \tilde{p}_{31}^a = 1$ and $\tilde{p}_{22}^a = \tilde{p}_{23}^a = \tilde{p}_{32}^a = \tilde{p}_{33}^a = 0$ for all $a \in \mathcal{A}$.

Taking these constraints together restricts the class of VE models to those of the form:

$$\tilde{P} = \begin{bmatrix} p_{11}^a & \tilde{p}_{12}^a & \tilde{p}_{13}^a \\ 1 & 0 & 0 \\ 1 & 0 & 0 \end{bmatrix}$$

where \tilde{p}_{1i}^a are “free variables” for all $i = 2, 3$ and $a \in \mathcal{A}$.

Notice that when $p_{11}^a \leq 0.5$ for all $a \in \mathcal{A}$, we can find a value equivalent model by setting: $(1 - \theta^a)/2 = p_{11}^a$. This means that the values produced by these value equivalent models exactly match those of the environment: $\tilde{v}^\pi = v^\pi$ for all π (and thus the solution of this model also coincides with the optimal value function, $\tilde{v}^* = v^*$).

A model class for which approximate VE outperforms MLE In the previous example we showed that it is possible to have an MDP and a restricted model class such that VE models are able to perfectly estimate v^* while MLE models fail to do so. Notice that in this example a value equivalent model *actually existed*, which is not guaranteed in general. We now show a related example where, in spite of an exactly value equivalent model not existing, an agent trained using an *approximate* value equivalent model will outperform its MLE counterpart.

We use our example MDP from before, shown in Figure 5, and denote its actions $\mathcal{A} = \{a, b\}$ for later notational convenience. We set our environment’s transition dynamics accordingly: $p^a \equiv (p_{11}^a, p_{12}^a, p_{13}^a) = (0.6, 0.4, 0.0)$ and $p^b \equiv (p_{11}^b, p_{12}^b, p_{13}^b) = (0.4, 0.2, 0.4)$. We also use the same model class as above: $(\tilde{p}_{11}^i, \tilde{p}_{12}^i, \tilde{p}_{13}^i) = (0.5(1 - \theta^i), \theta^i, 0.5(1 - \theta^i))$ for each $i \in \mathcal{A}$, being mindful of the boundary conditions $\theta^i \in [0, 1]$.

As we saw in the previous example, the MLE estimator for this problem will produce the following approximations: $p_{\text{MLE}}^a = (0.3, 0.4, 0.3)$, $p_{\text{MLE}}^b = (0.4, 0.2, 0.4)$.

We now consider what an approximate VE model will produce using the same value as before: $v = [1, 0, 0]^\top$ and all policies. Recall that we’re optimizing the following loss:

$$\begin{aligned} \sum_{j \in \{a, b\}} \sum_{i=1}^3 ((\tilde{P}^j v)_i - (P^j v)_i)^2 &= \sum_{j \in \{a, b\}} (\tilde{p}_{11}^j - p_{11}^j)^2 + ((\tilde{p}_{12}^j + \tilde{p}_{13}^j) - (p_{12}^j + p_{13}^j))^2 \\ &= \sum_{j \in \{a, b\}} (\tilde{p}_{11}^j - p_{11}^j)^2 + ((1 - \tilde{p}_{11}^j) - (1 - p_{11}^j))^2 \\ &= \sum_{j \in \{a, b\}} 2(\tilde{p}_{11}^j - p_{11}^j)^2 \\ &= 2(\tilde{p}_{11}^a - p_{11}^a)^2 + 2(\tilde{p}_{11}^b - p_{11}^b)^2. \end{aligned}$$

The form of this loss indicates that VE will attempt to minimize the MSE of \tilde{p}_{11}^a and \tilde{p}_{11}^b separately. Notice that for action a , we cannot perfectly estimate p_{11} due to the boundary conditions on θ^a . However, VE will still find the closest possible \tilde{p}_{11} that respects the boundary condition, giving: $\tilde{p}_{\text{VE}}^a = (0.5, 0.0, 0.5)$, $\tilde{p}_{\text{VE}}^b = (0.4, 0.2, 0.4)$.

We now display these models together in the following table:

	\tilde{p}_{11}^a	\tilde{p}_{12}^a	\tilde{p}_{13}^a	\tilde{p}_{11}^b	\tilde{p}_{12}^b	\tilde{p}_{13}^b
MDP	0.6	0.4	0.0	0.4	0.2	0.4
MLE	0.3	0.4	0.3	0.4	0.2	0.4
VE	0.5	0.0	0.5	0.4	0.2	0.4

Notice that when optimally planning on this MDP, an agent can obtain the most reward by transitioning from s_1 to s_1 as often as possible. The agent can do this taking the action among $\{a, b\}$ that is mostly likely to induce a self-transition each time it is at s_1 . In the true environment and the VE model this action is a . However, notice that the MLE model would instead prefer the sub-optimal action b , since $(\tilde{p}_{\text{MLE}}^b)_{11} > (\tilde{p}_{\text{MLE}}^a)_{11}$.

This is a concrete example where VE outperforms MLE even though there is no value-equivalent models in the model class considered (that is, VE can be enforced only approximately).

A.2 Experimental details

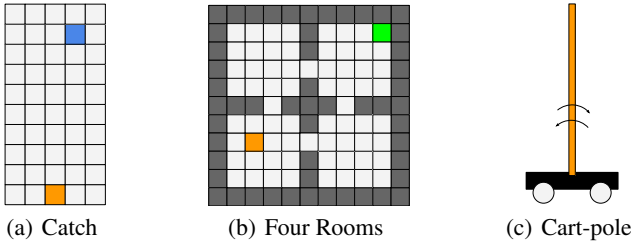


Figure 6: **(a) Catch:** the agent has three actions corresponding to moving a paddle (orange) left, right and staying in place. Upon initialization, a ball (blue) is placed at a random square at the top of the environment and at each step it descends by one unit. Upon reaching the bottom of the environment the ball is returned to a random square at the top. The agent receives a reward of 1.0 if it moves its paddle and intercepts the ball. **(b) Four Rooms:** the agent (orange) has four actions corresponding to up, down, left and right movement. When the agent takes an action, it moves in its intended direction with 90% of the time and in a random other direction otherwise. There is a rewarding square in the right top corner (green). If the agent transitions into this square it receives a reward of 1.0. **(c) Cart-pole:** In Cart-pole, the agent may choose between three actions: pushing the cart to the left, right or not pushing the cart. There is a pole balanced on top of the cart that is at risk of tipping over. The agent is incentivized to keep the pole up-right through a reward of $\cos(\theta)$ at each step where θ is the angle of the pole ($\theta = 0$ implies the pole is perfectly up-right). If the pole’s height drops below a threshold, the episode terminates and the agent receives a reward of 0.0. The cart itself is resting on a table; if it falls off the table, the episode similarly terminates with a reward of 0.0.

A.2.1 Environment description

The environments used in our experiments are described in depth in Figure 6. In both Catch and Four Rooms a tabular representation is employed in which each of the environment’s finitely many states (250 and 68, respectively) is represented by an index. In Cart-pole we have a continuous state space $S \subset \mathbb{R}^5$ (so $|S| = \infty$). Each state $s \in \mathbb{R}^5$ consists of the cart position, cart velocity, sine / cosine of pole angle, and pole’s angular velocity.

A.2.2 Experimental pipeline

As mentioned in the main text, a common experimental pipeline is used across all of our results, with slight variations depending upon the experiment type and environment. This pipeline is described at a high-level below:

- (i) **Data collection:** Data is collected using a policy which selects actions uniformly at random.
- (ii) **Model training:** The collected data is used to train a model.
- (iii) **Policy construction:** The model is used to produce a policy.
- (iv) **Policy evaluation:** The policy is evaluated to assess the quality of the model.

We now discuss steps (ii), (iii) and (iv) in detail.

(ii) Model training All of our experiments involve restricting the capacity of the class of models that the agent can represent: \mathcal{M} . In general we restrict the rank of the models in \mathcal{M} , but, depending upon the nature of the model, this restriction is carried in different ways.

1. **Tabular models:** On domains with $|S| < \infty$, we employ tabular models. In what follows, $n \times m$ matrices referred to as “row-stochastic” are ensured to be as such by the following parameterization:
 - (a) A matrix $F \in \mathbb{R}^{n \times m}$ is sampled with entries $F_{ij} \sim \text{Uniform}([-1, 1])$.

- (b) A new matrix \mathbf{P}_F is produced by applying row-wise softmax operations with temperature $\tau = 1$ to \mathbf{F} :

$$(\mathbf{P}_F)_{ij} = \frac{\exp(\mathbf{F}_{ij})}{\sum_k \exp(\mathbf{F}_{ik})}.$$

Here, \mathbf{F} can be thought of as the parameters of \mathbf{P}_F , which often will suppress as $\tilde{\mathbf{P}}$ for clarity.

That is, a model is represented by $|A| |S| \times |S|$ row-stochastic matrices: $\tilde{\mathbf{P}}^1, \dots, \tilde{\mathbf{P}}^{|A|}$. We ensure that each of these matrices has rank k by factoring it as follows: $\tilde{\mathbf{P}}^a = \mathbf{D}^a \mathbf{K}^a$ where $\mathbf{D}^a \in \mathbb{R}^{|S| \times k}$, $\mathbf{K}^a \in \mathbb{R}^{k \times |S|}$ and both are row-stochastic as well.

2. **Neural network models:** On domains with $|S| = \infty$ we instead use a neural network parameterized by θ : $f_\theta : (\mathcal{S}, \mathcal{A}) \mapsto (\mathcal{S}, \mathbb{R})$. f_θ takes a state and action as input and outputs an approximation of the expected next state and next reward. As an analogue to the rank restriction applied in the tabular case, we restrict the rank of weight matrices in all fully-connected layers in f_θ . Denote a fully-connected layer in f_θ as $L(x) = \sigma(Wx + b)$ where $\sigma(\cdot)$ is an activation function, W is a weight matrix and b is a bias term. We restrict f_θ by replacing each $L(x)$ with $L_k(x) = \sigma((DK)x + b)$ where $D, K \in \mathbb{R}^{|S| \times k}, \mathbb{R}^{k \times |S|}$.

The models with the restrictions above are trained based on data collected by a policy that selects actions uniformly at random. With a small abuse of notation, denote the collected data as $\mathcal{D} = (s_i, a_i, r_i, s'_i)_{i=1}^N$. We will now describe how this data is used to train models in different contexts.

1. **Tabular models:** When training a tabular model with capacity restricted to rank k , we use the following expressions:

- (a) **Reward:** In our experiments rewards are represented in the same way for both VE and MLE models:

$$\tilde{R}_{s,a} = \frac{\sum_{i=1}^N r_i \mathbb{1}\{s_i = s, a_i = a\}}{\sum_{i=1}^N \mathbb{1}\{s_i = s, a_i = a\}},$$

where $\mathbb{1}\{\cdot\}$ is the indicator function.

- (b) **Transition dynamics (MLE):** To learn the transition dynamics we first parameterize $\tilde{\mathbf{P}}^a = \mathbf{D}^a \mathbf{K}^a$ for all $a \in \mathcal{A}$, where \mathbf{D}^a and \mathbf{K}^a are row-stochastic matrices (see item 1 in the section ‘‘Restricting Model Capacity’’ above). Because we are assuming \mathcal{S} to be finite, we can identify each state $s \in \mathcal{S}$ by an index. Let $\delta(s) \in \{1, \dots, |S|\}$ be an index that uniquely identifies state s . We then compute $\tilde{\mathbf{P}}^a = \mathbf{D}^a \mathbf{K}^a$ by minimizing the following loss with respect to \mathbf{D}^a and \mathbf{K}^a :

$$\tilde{\ell}_{p,\mathcal{D}}(\mathbf{P}^a, \tilde{\mathbf{P}}^a) \equiv - \sum_{i=1}^N \mathbb{1}\{a_i = a\} \log \left[(\mathbf{D}^a \mathbf{K}^a)_{\delta(s_i)\delta(s'_i)} \right],$$

where $(\mathbf{D}^a \mathbf{K}^a)_{ij}$ is the element in the i -th row and j -th column of matrix $\mathbf{D}^a \mathbf{K}^a$. Note that the expression above is the empirical version of expression (5) in the paper [15].

- (c) **Transition dynamics (VE):** In the VE setting we have a set of value functions and policies: \mathcal{V} and Π . We have one transition matrix $\tilde{\mathbf{P}}^\pi$ associated with each policy $\pi \in \Pi$. As discussed in Section 5, in our experiments we used $\Pi = \{\pi^a\}_{a \in \mathcal{A}}$, where $\pi^a(a|s) = 1$ for all $s \in \mathcal{S}$. Thus, we end up with the same parameterized probability matrices as above: $\tilde{\mathbf{P}}^a = \mathbf{D}^a \mathbf{K}^a$. Let $\mathcal{D}_{ia} \subseteq \mathcal{D}$ be the sample transitions starting in state i where action a was taken, that is, $(s_j, a_j, r_j, s'_j) \in \mathcal{D}_{ia}$ if and only if $\delta(s_j) = i$ and $a_j = a$. We computed $\tilde{\mathbf{P}}^a = \mathbf{D}^a \mathbf{K}^a$ by minimizing the following loss with respect to \mathbf{D}^a and \mathbf{K}^a :

$$\ell_{\pi^a, \mathcal{V}, \mathcal{D}}(\mathbf{P}^a, \tilde{\mathbf{P}}^a) \equiv \sum_{i,a} \sum_{\mathbf{v} \in \mathcal{V}} \left(\frac{1}{|\mathcal{D}_{ia}|} \sum_{(s,a,r,s') \in \mathcal{D}_{ia}} v_{\delta(s')} - \sum_j (\mathbf{D}^a \mathbf{K}^a)_{ij} v_j \right)^2.$$

Note that the expression above corresponds to equation (7) when learning transition matrices associated with policies $\{\pi^a\}_{a \in \mathcal{A}}$ in an environment with finite state space \mathcal{S} (where states s can be associated with an index i) and $p = 2$.

2. **Neural network models:** When training a neural network model with capacity restrictions construct a network $f_\theta : (\mathcal{S}, \mathcal{A}) \mapsto (\mathcal{S}, \mathbb{R})$. The network is fully connected and takes the concatenation of \mathcal{S} with the one-hot representation of \mathcal{A} as input. For a given (s, a) pair we denote its output as $\tilde{s}'_{s,a}, \tilde{r}'_{s,a} = f_\theta(s, a)$. In all cases we train the neural network model by sampling mini-batches uniformly from \mathcal{D} . It is important to note that we only use these neural network models on deterministic domains (e.g., Cart-pole) meaning that the output of the model, \tilde{s}' represents a single state rather than an expectation over states.

(a) **Reward:** For both VE and MLE models we train our neural network models to accurately predict the reward associated with each state action transition:

$$\ell_{r, \mathcal{D}}(\theta) = \sum_{i=1}^N (\tilde{r}_{s_i, a_i} - r_i)^2.$$

(b) **Transition dynamics (MSE):** We learn models by encouraging f_θ to accurately predict the next state:

$$\ell_{s', \mathcal{D}}(\theta) = \sum_{i=1}^N (\tilde{s}'_{s_i, a_i} - s'_i)^2.$$

(c) **Transition dynamics (VE):** For VE models use (7), disregarding reward terms to give:

$$\ell_{\mathcal{V}, \mathcal{D}}(\theta) = \sum_{i=1}^n \sum_{v \in \mathcal{V}} (v(\tilde{s}_{s_i, a_i}) - v(s'_i))^2.$$

(iii) Policy construction In each experiment we present, after a model is constructed, we subsequently use it to construct a policy. The manner in which we do this varies based upon the type of the experiment and the nature of the environment. There are three mechanisms for constructing policies from models:

1. **Value iteration:** For experiments with $\mathcal{V} = \mathbb{V}$ (which are performed only with tabular models), we use the learned model $\tilde{m} = (\tilde{r}, \tilde{p})$ to perform value iteration until convergence, yielding \tilde{v}^* [30]. Here \tilde{v}^* represents the optimal value function of the model \tilde{m} . We then produce a policy according to $\pi(s) = \operatorname{argmax}_a (\tilde{r}(s, a) + \gamma \sum_{s'} \tilde{p}(s'|s, a) \tilde{v}^*(s'))$.
2. **Approximate policy iteration with least squares temporal-difference learning (LSTD):** For experiments on environments with finite \mathcal{S} and $\mathcal{V} = \mathcal{V}$ we used policy iteration combined with least square policy evaluation using basis $\{\phi_i\}_{i=1}^d$. Specifically, each iteration of policy iteration involved the following steps:
 - (a) Collect experience tuples using the previous policy, π , leading to $\mathcal{D} = (s_i, a_i, r_i, s'_i)_{i=1}^n$.
 - (b) Replace the reward and next-states with those predicted by the model: $\tilde{r}_i, \tilde{s}'_i = f_\theta(s_i, a_i)$, leading to $\mathcal{D}' = (s_i, a_i, \tilde{r}_i, \tilde{s}'_i)_{i=1}^n$.
 - (c) Learn $v_w(s) = \sum_{i=1}^d w_i \phi_i(s) \approx v_\pi$ using LSTD with \mathcal{D}' .
 - (d) Construct a new policy $\pi(s) = \operatorname{argmax}_a (\tilde{r}_{s,a} + \gamma v_w(\tilde{s}'_{s,a}))$ where $\tilde{r}_{s,a}, \tilde{s}'_{s,a}$ are sampled from the trained model conditioned on state s and action a .

This procedure is repeated for a fixed number of iterations.

3. **Deep Q-networks (DQN):** For experiments with $\mathcal{V} = \tilde{\mathcal{V}}$ and infinite \mathcal{S} we use Double Q-Learning to produce policies. We incorporate our learned model, f_θ , by replacing elements in the replay buffer (s, a, r, s') with $(s, a, \tilde{r}_{s,a}, \tilde{s}'_{s,a})$ where $\tilde{r}_{s,a}, \tilde{s}'_{s,a} = f_\theta(s, a)$.

(iv) Policy evaluation There are two methods to evaluate the policies resulting from the policy construction stage described above:

1. For policies produced using value iteration or policy iteration plus LSTD the ensuing policy, π , is exactly evaluated on the true environment, yielding $v_\pi(s)$. Then the average value of $v_\pi(s)$ over all states is reported.
2. For policies produced using DQN, the average return over the last 100 episodes of training is reported.

A.2.3 Classes of experiments

In addition to varying the capacity of \mathcal{M} , there are two primary classes of experiments that were run in our paper that assess different choices of \mathcal{V} . We distinguish between these two classes below:

$\text{span}(\mathcal{V}) \approx \tilde{\mathcal{V}}, \tilde{\mathcal{V}} = \mathbb{V}, \Pi = \mathbb{I}$: In these experiments we consider that there is no limitation on the agent’s ability to represent value functions, and focus on achieving value equivalence with respect to the polytope of value functions \mathbb{V} induced by the environment. We enable the agent to represent arbitrary functions in \mathbb{V} by restricting ourselves to tabular environments and using dynamic programming to perform exact value iteration in our Policy Construction step. We approximate the value polytope by randomly sampling deterministic policies: $\{\pi_1, \dots, \pi_n\}$ and evaluating them (again using dynamic programming) to produce $\{v_{\pi_1}, \dots, v_{\pi_n}\}$. We then choose $\mathcal{V} = \{v_{\pi_1}, \dots, v_{\pi_n}\}$. In this setting we vary the number of policies generated.

Corresponding experiments: the experiments in this class vary two dimensions: (1) the rank of the model and (2) the number of policies generated. In Figures 3(a) and 3(b) we depict plots for the Four Rooms environment that fix the number of policies while varying the rank of the model and plots that fix the rank of the model while varying the number of policies, respectively. Figures 3(c) and 3(d) are analogous plots for the Catch environment.

$\text{span}(\mathcal{V}) \approx \tilde{\mathcal{V}}, \Pi = \mathbb{I}$: In these experiments we explore the setting described in Remark 2. We assume that the agent has variable ability to represent value functions, $\tilde{\mathcal{V}}$, and attempt to learn a model in $\mathcal{M}(\tilde{\mathcal{V}}, \mathbb{I})$. From Proposition 1 we only need to find \mathcal{V} such that $\text{span}(\mathcal{V}) \supseteq \tilde{\mathcal{V}}$. Experiments in this class can further be broken down into two settings based upon the nature of $\tilde{\mathcal{V}}$:

- (a) **Linear function approximation:** In certain experiments our agent uses a class of linear function approximators to represent value functions: $\tilde{\mathcal{V}} = \{\tilde{v} : \tilde{v}(s) = \sum_{i=1}^d \phi_i(s)w_i\}$ where $\phi_i(s) : \mathcal{S} \mapsto \mathbb{R}$ and $w \in \mathbb{R}^d$. In this setting achieving $\text{span}(\mathcal{V}) \supseteq \tilde{\mathcal{V}}$ can be satisfied by choosing $\mathcal{V} = \{\phi_i\}_{i=1}^d$. For experiments using linear function approximation, we select our features $\{\phi_i\}_{i=1}^d$ to correspond to state aggregations. This entails the following procedure:
 - (i) Collect data using a policy that selects actions uniformly at random.
 - (ii) For tabular domains (e.g., Catch, Four Rooms), convert tabular state representations into coordinate-based representations. For Catch we convert each tabular state into the positions of both the paddle and the ball: $(x_{\text{paddle}}, y_{\text{paddle}}, x_{\text{ball}}, y_{\text{ball}})$. For Four Rooms we use the position of the agent: $(x_{\text{agent}}, y_{\text{agent}})$. Denote the function that performs this conversion as: $f : \mathcal{S} \mapsto \mathbb{R}^n$ where $n = 2$ and $n = 4$ for Four Rooms and Catch respectively.
 - (iii) Perform k-means clustering on these converted states to produce d centers $c_{1:d}$.
 - (iv) Define $\phi_i(s) = \mathbb{1}\{\text{argmin}_j \|f(s) - c_j\|_2 = i\}$, which corresponds to aggregating states according to their proximity to the previously calculated centers.

Corresponding experiments: the experiments in this class vary two dimensions: (1) the rank of the model and (2) the number of basis functions in $\{\phi_i\}_{i=1}^d$. In Figures 4(a) and 4(b) we depict plots of “slices” of this two-dimensional set of results on the Catch domain: 4(a) depicts fixing the number of basis functions while varying model-rank and 4(b) depicts fixing the model-rank while varying the number of basis functions.

- (b) **Neural network function approximation:** When Neural Networks are used to approximate the agent’s value functions we have $\tilde{\mathcal{V}} = \{\tilde{v} : \tilde{v}(s) = g_\theta(s)\}$ where g_θ represents a neural network with a particular architecture parameterized by θ . In our experiments we choose the architecture of g_θ to be a 2 layer neural network with a tanh activation for its hidden layer. Unlike the linear function approximation setting, it is less obvious how to choose \mathcal{V} such that $\text{span}(\mathcal{V}) \supseteq \tilde{\mathcal{V}}$. One option is to use randomly initialized neural networks in $\tilde{\mathcal{V}}$ as our basis. To randomly initialize a given layer in some network g_θ , we select weights from a truncated normal distribution where $\mu = 0$ and $\sigma = 1/\sqrt{\text{layer-input-size}}$ and initialize biases to 0.

However, we found in practice that a large number of these randomly initialized networks were required to achieve reasonable performance. Instead of maintaining a large set of

initializations in \mathcal{V} , we allow the elements of \mathcal{V} themselves to be stochastic. Every time we apply an update of gradient descent we sample a new set of randomly initialized neural networks to function as \mathcal{V} . This is equivalent to minimizing $\mathbb{E}_{\mathcal{V}}[\ell_{\Pi, \mathcal{V}, \mathcal{D}'}(m^*, \tilde{m})]$ where $\ell_{\Pi, \mathcal{V}, \mathcal{D}'}$ is defined in 7. We find that having more random elements in \mathcal{V} decreases the variance in the performance of VE models; $|\mathcal{V}| = 5$ in our experiments.

Corresponding experiments: the experiments in this class vary two dimensions: (1) the rank of the model and (2) the width of the neural networks in \mathcal{V} . In Figures 4(c) and 4(d) we depict plots of “slices” of this two-dimensional set of results on the Catch domain: 4(c) depicts fixing the network width while varying model-rank and 4(d) depicts fixing the model-rank while the network width varies.

A.2.4 Additional results

In the experimental section of the main text we showed that our theoretical claims about the value equivalence principle hold in practice through a series of bivariate experiments (e.g., varying model-rank and number of bases, varying model-rank and number of policies, varying model-rank and network width). We displayed our results as “slices” of these bivariate experiments, where one variable would be held fixed and the other would be allowed to vary. To keep the paper concise, we only displayed a subset of these slices where the “fixed” variable was selected as the median value over full set we experimented with. In what follows, we present the complete set of the experimental results we acquired. We indicate that a plot was included in the main text by printing its caption in bold font.

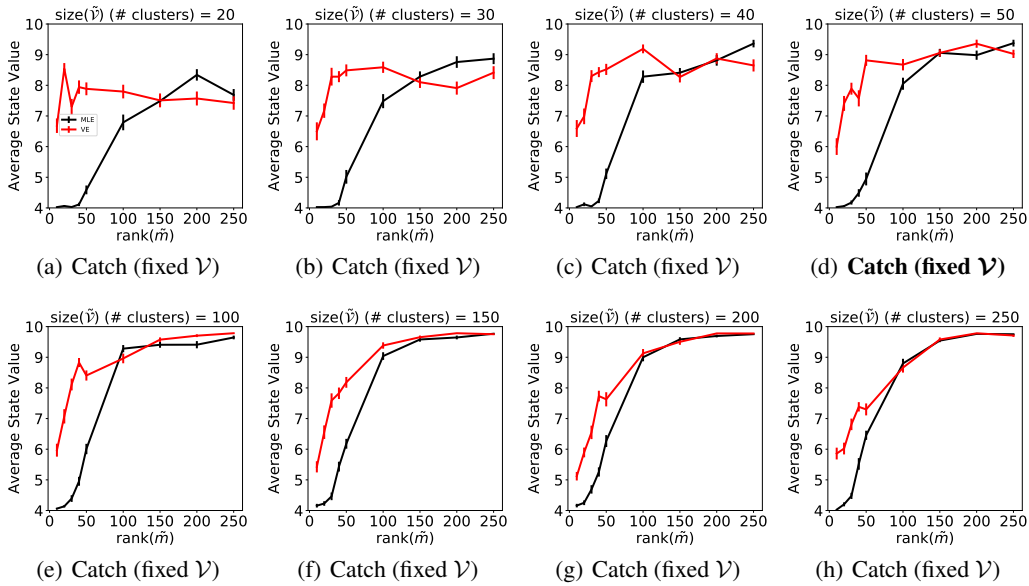


Figure 7: All Catch results with fixed \mathcal{V} and $\text{span}(\mathcal{V}) \approx \tilde{\mathcal{V}}$.

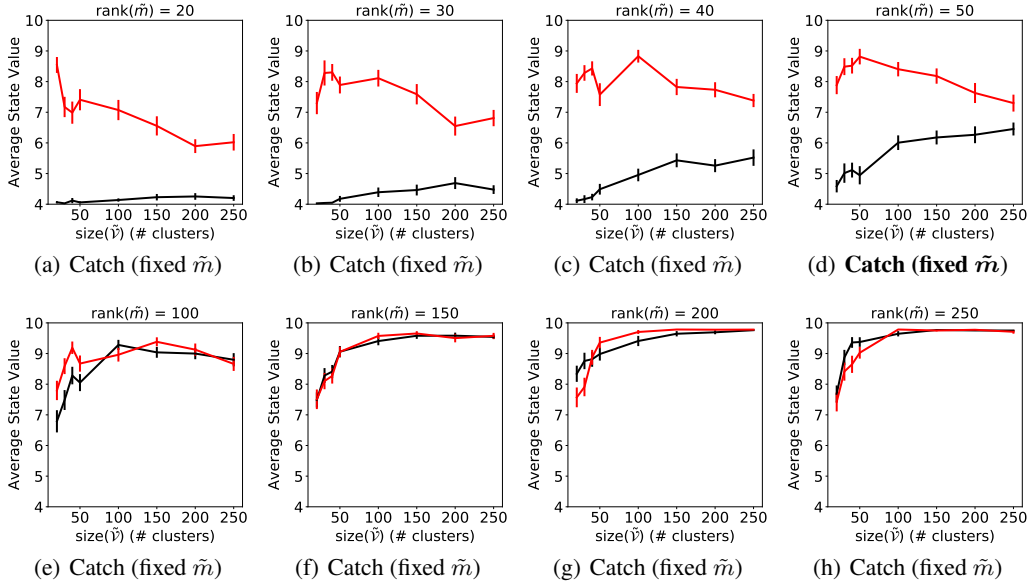


Figure 8: All Catch results with fixed \tilde{m} and $\text{span}(\mathcal{V}) \approx \tilde{\mathcal{V}}$.

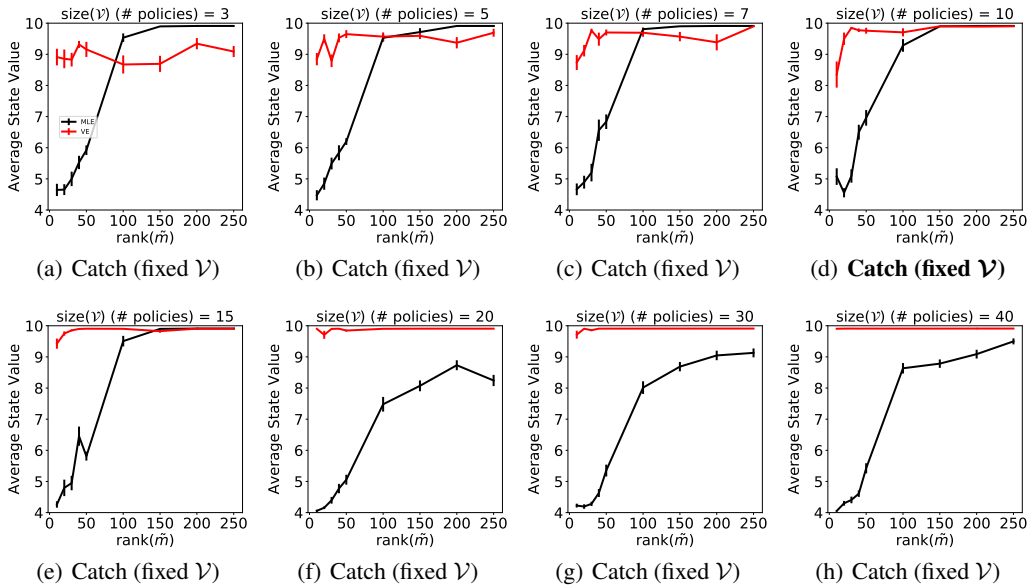


Figure 9: All Catch results with fixed \mathcal{V} and $\mathcal{V} = \{v_{\pi_1}, \dots, v_{\pi_n}\}$.

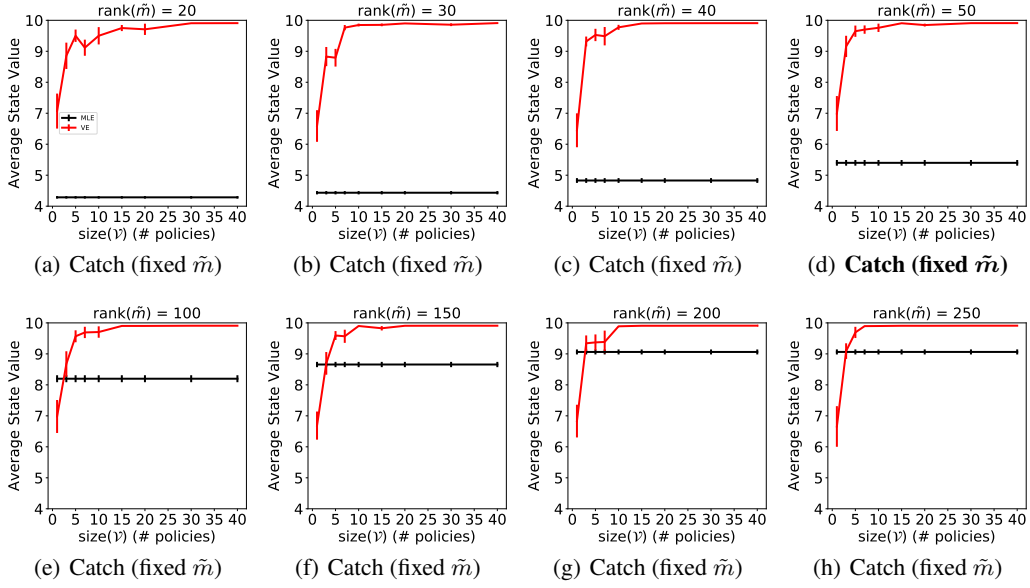


Figure 10: All Catch results with fixed \tilde{m} and $\mathcal{V} = \{v_{\pi_1}, \dots, v_{\pi_n}\}$.

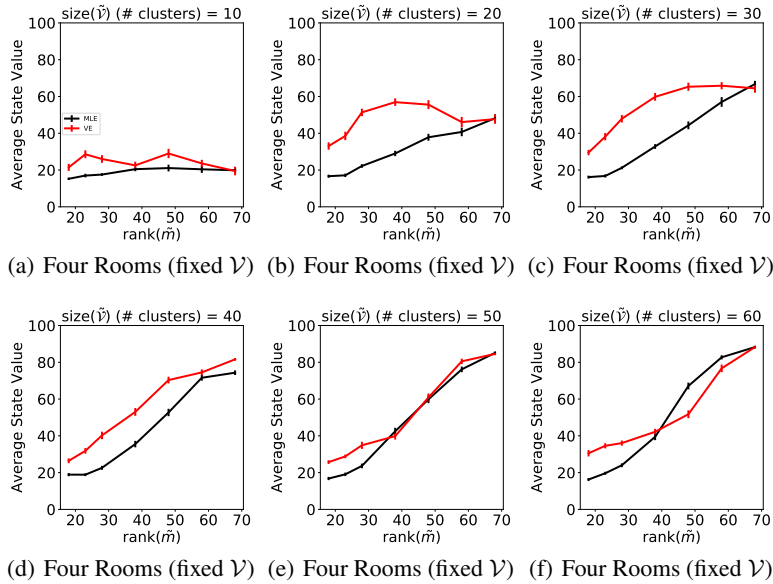


Figure 11: All Four Rooms results with fixed \mathcal{V} and $\mathcal{V} = \tilde{\mathcal{V}}$.

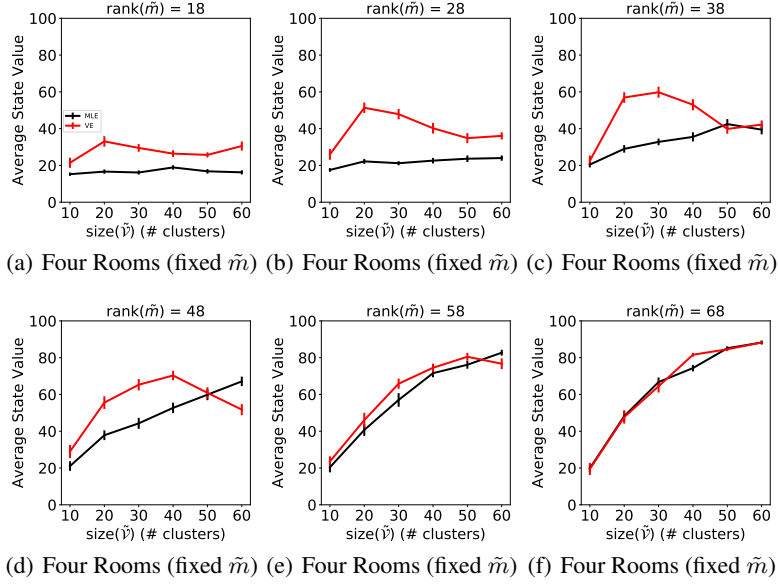


Figure 12: All Four Rooms results with fixed \tilde{m} and $\mathcal{V} = \tilde{\mathcal{V}}$.

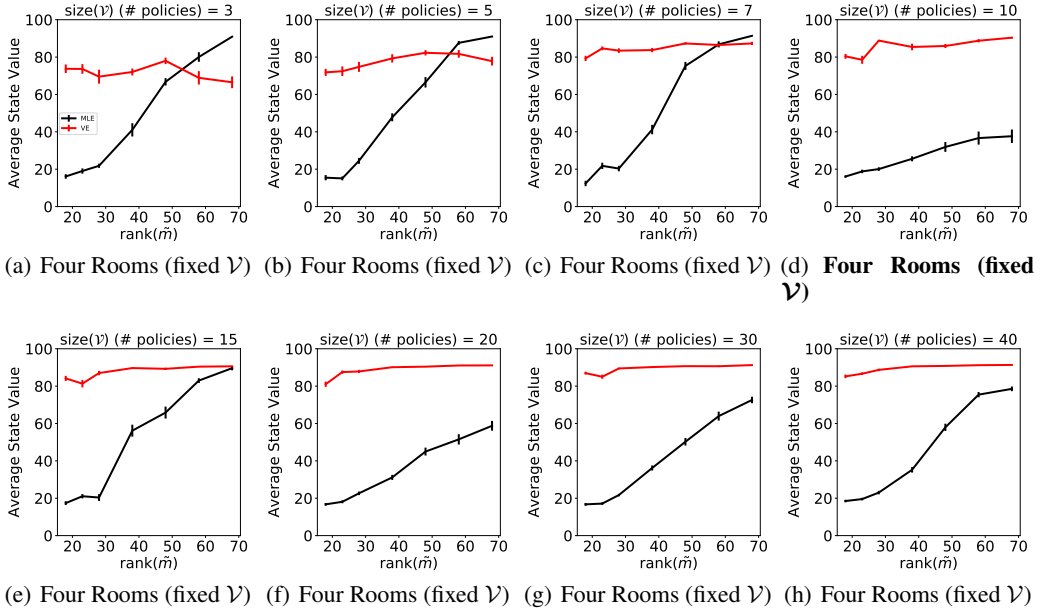


Figure 13: All Four Rooms results with fixed \mathcal{V} and $\tilde{\mathcal{V}} = \{v_{\pi_1}, \dots, v_{\pi_n}\}$.

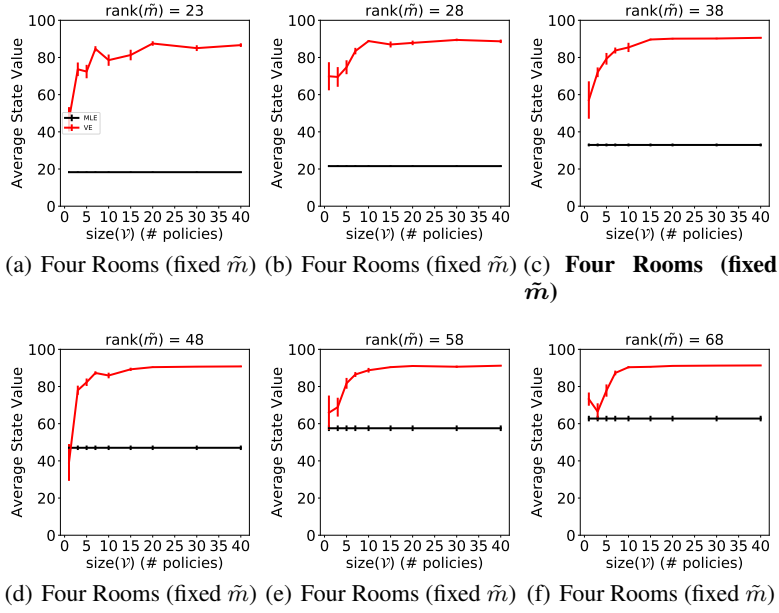


Figure 14: All Four Rooms results with fixed \tilde{m} and $\mathcal{V} = \{v_{\pi_1}, \dots, v_{\pi_n}\}$.

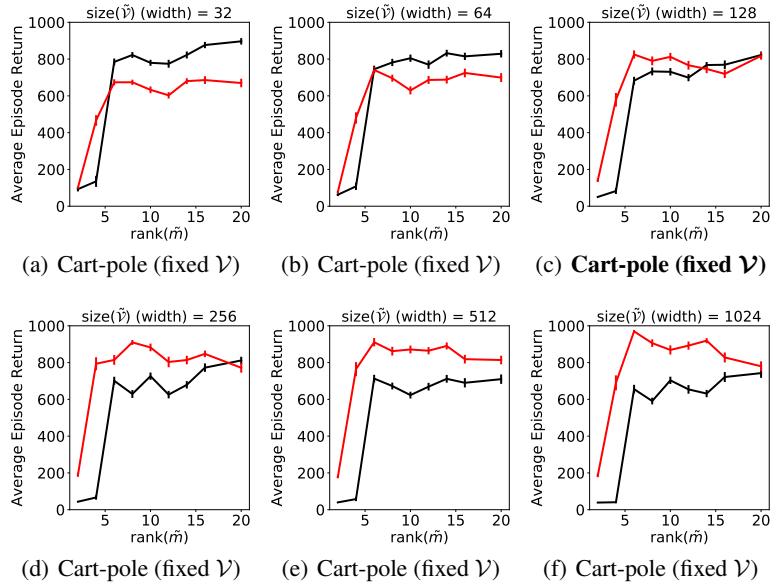


Figure 15: All Cart-pole results with fixed \mathcal{V} and $\text{span}(\mathcal{V}) \approx \tilde{\mathcal{V}}$.

A.2.5 Hyperparameters

Table 1 provides a list detailing the different hyperparameters used throughout our pipeline.

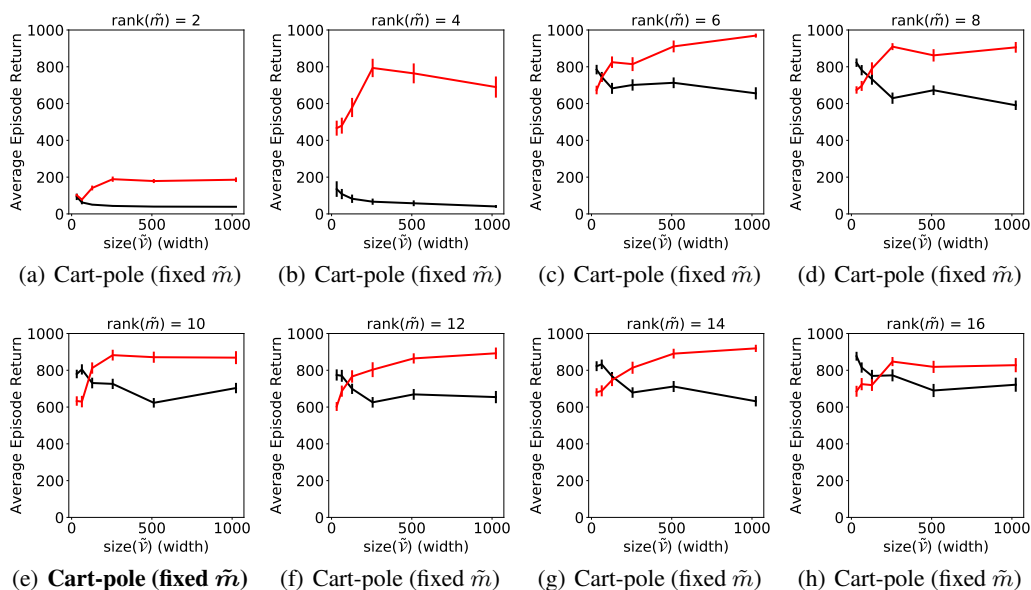


Figure 16: All Cart-pole results results with fixed \tilde{m} and $\text{span}(\mathcal{V}) \approx \tilde{\mathcal{V}}$.

Hyperparameter	Value	Description
minibatch size	32	Number of samples passed at a time during a training step of any learning method.
model learning rate	5e-5	Learning rate used to train all models.
# model samples	1,000,000	Number of transitions sampled by a random policy in the Data Collection phase.
model depth	2	Number of hidden layers in the model architecture.
model width	256	Number of units per hidden layer.
model activation	tanh	Activation function following a hidden layer.
model learning max steps	1,000,000	Maximum number of training iterations.
γ	0.99	Discount factor used across environments.
LSTD samples / policy	10,000	Number of samples collected for each phase of policy evaluation using LSTD.
# policy iteration steps	40	Number of steps of policy iteration in the policy construction phase, when applicable.
DQN learning rate	5e-4	Learning rate for DQN.
DQN # environment steps	2,500,000	Number of environment steps that DQN learns over.
DQN learning frequency	4	A learning update is applied after this many environment steps.
DQN depth	1	Number of hidden layers in the DQN.
DQN activation	tanh	Activation function following a hidden layer.
DQN target update	2000	Number of environment steps before the target network in the DQN is updated.
Tabular # eval episodes	20	Number of episodes to average performance over to assess a policy in the tabular setting.
DQN # eval episodes	100	Number of episodes to average DQN policy performance over at the end of training.
DQN ϵ	0.05	Chance of picking a random action during training.
Optimizer	Adam	Optimizer used for all learning operations. Default Adam parameters were used.

Table 1: List of hyperparameters used in the experiments.